

Dan Tufis, Florin Gh. Filip (coordonatori)

Limba Româna
în
Societatea Informatională - Societatea Cunoașterii

EDITURA
Expedit



Bucuresti, România

Editor: **Valeriu IOAN-FRANC**

Coperta si prezentarea grafica: **Nicolae LOGIN, Luminita LOGIN**

Toate drepturile asupra acestei editii apartin Editurii Expert. Reproducerea, fie si partiala si pe orice suport, este interzisa fara acordul prealabil al editorului, fiind supusa prevederilor legii drepturilor de autor.

ISBN 973- 8177- Aparut 2002

Dan Tufis, Florin Gh. Filip (coordonatori)

Limba Română
în
Societatea Informațională -
Societatea Cunoașterii

EDITURA
Expert

DEDICATIE

Acest volum este dedicat Academicianului Mihai Draganescu, Profesorul si mentorul unei întregi generatii de specialisti în stiinta si tehnologia informatiei în general si al problemelor societatii informationale si a cunoasterii în special. Marea majoritate a contributiilor din acest volum apartin unor experti ce fac parte din Comisia de Informatizare a Limbii Române, comisie a Academiei Române la a carei nastere un rol esential l-a avut Profesorul Draganescu, presedintele Sectiei de Stiinta si Tehnologia Informatiei. Savantul Mihai Draganescu are numeroase contributii în stiinta contemporana, binecunoscute atât în tara cât si în strainatate. Pentru cine îl cunoaste pare incredibila puterea sa de munca, debordanta creativitate si neostoita cautare a noului. Profesorul Draganescu este indiscutabil port-drapelul conceptului de societate informationala-societate a cunoasterii în România. În lucrarile sale din urma cu peste 25-30 de ani se regasesc cu claritate multe concepte foarte actuale în zilele noastre, previziuni curajoase atunci, acum realitati cotidiene. În lucrarile domniei sale din ultima vreme, apare un nou concept ce avem convingerea ca se va impune: Societatea Constiintei, o treapta superioara a societatii cunoasterii. Nu este de mirare deci ca în contextul societatii informationale si a cunoasterii profesorul Draganescu a sustinut cu consecventa si a afirmat cu claritate rolul Inteligentei Artificiale în devenirea noilor societati ale cunoasterii. Între domeniile Inteligentei Artificiale un loc de frunte în promovarea principiilor societatii cunoasterii îi revine Tehnologiei Limbajului Natural. Profesorul Draganescu a fost unul dintre putinii oameni de stiinta români care au înteles si au sprijin total aceste directii. Cu aproape douazeci de ani în urma (1983), Profesorul Draganescu edita (împreuna cu Adrian Davidoviciu si Ioan Georgescu) volumul „Inteligenta Artificiala si Robotica” pentru ca trei ani mai târziu (împreuna cu Corneliu Burileanu) sa editeze un alt volum de referinta „Analiza si sinteza semnalului vocal”. Astazi, cercetarile mondiale în domeniul tehnologiilor lingvistice au atins un nivel de maturitate ce permit sinergizarea eforturilor lingvistilor, informaticienilor, matematicienilor si a altor specialisti din sectorul academic sau industrial, sa abordeze proiecte mari, interdisciplinare având ca obiectiv prelucrarea automata, în mediile de comunicare electronica, a din ce în ce mai multe limbi naturale. Printre acestea, limba româna își face loc încet dar sigur. Volumul de fata este o marturie în acest sens. În acelasi timp, volumul se constituie într-o noua confirmare a realitatilor pe care Profesorul Mihai Draganescu le prefigura cu multi ani în urma.

Dr. Dan Tufis, m.c.A.R, Acad. Florin Gh. Filip

CUPRINS

INTRODUCERE	13
-------------------	----

SECTIUNEA I:

LINGVISTICA TEORETICA SI FORMALA; LEXICOGRAFIE

Resurse lingvistice elaborate la Institutul de Lingvistică „Iorgu Iordan” - <i>Ioana Vintilă-Rădulescu</i>	21
Contribuția lingvisticii la studiul terminologiilor științifice - <i>Angela Bidu-Vrănceanu</i>	35
Gramaticile generative nontransformaționale - <i>Emil Ionescu</i>	41
Către o teorie X-bar funcțională - <i>Neculai Curteanu</i>	53
Teoria HPSG. Studiu de caz: acordul încrucișat - <i>Ana Maria Barbu</i>	89
După 10 ani de experiență terminografică: noul model de date terminologice al TermRom - <i>Dan Matei</i>	111
Probleme de reprezentare a datelor terminografice într-o bază de date relațională - <i>Sorin Ghețaru</i>	123

SECTIUNEA II:

TEHNOLOGII ALE LIMBAJULUI SCRIS

RO-BALKANET - ontologie lexicalizată în context multilingv pentru limba română - <i>Dan Tufiș, Dan Cristea</i>	139
Algoritmi de segmentare a textului în unități de tip clauzal - <i>Dan Gălea, Niculai Curteanu, Constantin Linteș</i>	167
O metodă automată pentru inserarea diacriticelor în texte în limba română- <i>Rada F. Mihalcea, Vivi A. Năstase</i>	193
Contribuții privind structura statistică de cuvinte în limba română scrisă - <i>Adriana Vlad, Adrian Mitrea</i>	209
Dezambiguizarea semantică automată în corpusuri paralele - <i>Dan Tufiș</i>	237

Referențialitate și cursivitate în structura de discurs - <i>Dan Cristea</i>	271
DLIR - un sistem de căutare documentară multilingv - <i>Amalia Todirașcu</i>	305
Mediu hermenofor pentru asistarea învățării unor concepte dintr-o limbă străină - <i>Ștefan Trăușan-Matu</i>	319

**SECTIUNEA III:
TEHNOLOGII ALE LIMBAJULUI VORBIT**

Experimente în vederea recunoașterii vorbitorului - <i>Corneliu Burileanu, Luigi Bojan</i>	337
Prelucrarea inițială a textului de intrare în cadrul unui sistem de sinteză a vorbirii pornind de la text în limba română - <i>Dragoș Burileanu</i>	359
Utilizarea tehnicilor nuanțate (fuzzy) și de dinamică neliniară pentru sinteza adaptivă a vorbirii - <i>Horia N. Teodorescu</i>	381
Dicționarele multimedia ale limbii române. Secvențe de implementări și experimentări - <i>Dumitru Todoroi, Diana Micusa, Zinaida Todoroi, Ion Linga, Ion Covalenco, Nicolae Objeleanu, Ștefan Spătaru, Stela Lungu, Virginia Țurcanu, Elena Cozlov, Nadejda Ambrozii, Victor Slobodeanu, Igor Coșeru, Cătălina Suruceanu</i>	401
Mediu pentru editarea Transcrierilor Fonetice în Limba Română. Realizarea Atlaselor Lingvistice Românești Regionale - <i>Silviu Bejinariu, Vasile Apopei, Mariana Roma, Horia N. Teodorescu</i>	423

**SECTIUNEA IV:
DEZBATERI ȘI DISCUȚII**

Asupra a doi vectori funcționali ai Societății Cunoașterii: Managementul Cunoașterii și Învățarea Electronică. Cultura și Societatea Cunoașterii - <i>Mihai Drăgănescu</i>	441
Între lingvistica matematică și cea computațională - <i>Solomon Marcus</i>	471
Între lingvistica matematică și cea computațională: o altă perspectivă - <i>Dan Tufiș</i>	481

INTRODUCERE

Programul de cercetare aplicativa „Strategii si solutii pentru Societatea Informationala – Societatea Cunoasterii în România (SI-SC), din subprogramul A-strategic, al Programului National INFOSOC a avut ca principale obiective stabilirea unui program de veghe conceptuala pentru mentinerea pe linia tendintelor mondiale ale avansului SI-SC, sensibilizarea factorilor de decizie si a publicului larg, crearea unui cadru de reflectie prospectiva pe temele prioritare ale SI-SC: economice, sociale, culturale, tehnologice, ambientale, precum si operationalizarea unor solutii de interes prioritar pe plan national. În cadrul acestui proiect a fost elaborat volumul „*Societatea Informationala – Societatea Cunoasterii. Concepte, solutii si strategii pentru România*” (publicat la Ed. Expert în anul 2000), realizat sub coordonarea Academicianului Florin Gheorghe Filip. Acest volum avea ca scop construirea unei viziuni si continea o serie de studii si cercetari care au aprofundat rezultatele programului prioritar al Academiei Române privind *Societatea Informationala – Societatea Cunoasterii* si au identificat o serie de orientari strategice cerute de sustinerea unei dezvoltari de tip “salt” a SI-SC în România. Prin prisma obiectivelor proiectului, au fost analizate principalele aspecte conceptuale ale SI-SC, probleme legate de infrastructurile informatice si de comunicatii ale SI-SC, formarea profesionala si pregatirea generala a populatiei în si pentru SI-SC, rolul stiintei, cercetarii si inovarii, aspecte sociale si juridice, institutiile statului si relatia lor cu cetateanul, dezvoltarea economiei si afacerilor, dimensiunea culturala a SI-SC, actorii sociali ai crearii si difuzarii tehnologiei informatiei si comunicatiilor în contextul SI-SC. Studiile tematice, ancheta Delphi pentru consultarea opiniei expertilor privind tendintele globale si optiunile posibile de raportare la ele, scenariile de evolutie elaborate au sustinut functia prospectiva a proiectului.

Functia operativa a acestui proiect, respectiv identificarea de solutii tehnice privind rezolvarea principalelor prioritati identificate în faza analizei prospective urma sa se manifeste în perioada imediat urmatoare, printr-o dintr-o serie de cercetari/dezvoltari tehnologice ce vor trata pe larg problematica specifica a fiecaruia dintre directiile amintite anterior. Aceasta serie este deschisa prin prezentul volum ce înglobeaza contributiile ale unor specialisti români reprezentativi în domeniul prelucrării automate a limbajului natural si a resurselor lingvistice necesare utilizarii limbii române în mediile de comunicare electronica.

În [1] este definit conceptul de “Societate Informationala – Societate a Cunoasterii” (SI-SC) precum si principalii sai vectori tehnologici si functionali. În acest context „internetul dezvoltat” (ca vector tehnologic) si ”managementul utilizarii morale a cunoasterii la nivel global” (ca vector functional) sunt prezentati ca factori motrici esentiali ai Societatii Cunoasterii, si în perspectiva, a Societatii Constiintei. „Din momentul în care intervine Internetul cu marile avantaje pe care acesta le aduce (e-mail, comert electronic si tranzactii electronice, piata Internet, distributia de ‘continut’) prin cuprinderea în sfera informatiei electronice a unui numar cât mai mare de cetateni se trece la societatea

informationala. Cunoasterea este informatie cu înțeles și informație care acționează. De aceea societatea cunoașterii nu este posibilă decât greșită pe societatea informațională și nu poate fi separată de aceasta. În același timp, ea este mai mult decât societatea informațională prin rolul major care revine informației–cunoaștere în societate.” [1]

În 1984, William Gibson, un dizident cognitiv - după cum se auto-caracterizează, publică volumul SF „*Neuromancer*” (Ace Book, July 1984, ISBN: 0-441-56959-5), carte care pe lângă o mulțime de premii literare i-a adus notorietatea și pentru crearea termenului „cyberspace”: „the total interconnectedness of human beings through computers and telecommunication without regard to physical geography... A consensual hallucination experienced daily by billions of legitimate operators, in every nation, by children learning mathematical concepts...a graphical representation of data abstracted from the banks of every computer in the human system. Unthinkable complexity. Lines of light ranged in the nonspace of the mind. Clusters and constellations of data. Like city lights receding...” (op. cit).

Termenul a făcut cariera, actualmente fiind o noțiune care din punct de vedere tehnic subsumează conceptul „Internet”(scris cu majusculă): "cyberspace: The impression of space and community formed by computers, computer networks, and their users; the virtual "world" that Internet users inhabit when they are online ... The term internet (spelled with a lower case "i") is distinguished from the Internet (spelled with the "I" capitalized). The Internet refers to a specific, historic, ubiquitous worldwide digital communication network.” (cf. Glossary of Telecommunications, American National Standard T1.523-2001, www.atis.org/tg2k/cyberspace.html, 05.08.2002).

Dimensiunea tehnică (evocată mai sus) a noțiunii de „*ciberspațiu*” este complementată de dimensiunea socio-culturală și din această perspectivă problemele „satului global” previzionate de Societatea Informațională – Societatea Cunoașterii. Ideea atenuării schismei dintre specialiștii din domeniul tehnic și cei din zona științelor umaniste în contextul SI-SC este susținută puternic și de M. Derouzos [5], cel care a propus conceptul de „piața informațională”, pe care îl considera mai realist decât cel de „ciberspațiu”. De altfel, dimensiunile socio culturale ale SI-SC au fost evocate în capitolele 2,3,4 și 6 ale volumului „*Societatea Informațională – Societatea Cunoașterii. Concepte, soluții și strategii pentru România*”

Printre componentele socio-culturale ale SI-SC, utilizarea limbii materne și a accesului universal la ciberspațiu [2, 3, 4] constituie o prioritate.

În contextul actual, al comunicării mediate de tehnologia informației și de telecomunicații, limba devine obiect al investigației tehnice. Tehnologia limbajului impune metodologii specifice de cercetare/dezvoltare, dezvoltarea sau adaptarea resurselor lingvistice fundamentale cum ar fi dicționarele, tezaurele, corpusurile și gramaticile computerizate, în conformitate cu standardele sau recomandările existente. În funcție de resursele lingvistice disponibile, de volumul și calitatea lor, de compatibilitatea codificării lor în raport cu recomandările și standardele internaționale etc., se poate vorbi de *nivelul de tehnologizare* al unei limbi naturale. Nivelul de tehnologizare al unei limbi naturale este în corespondență directă cu statutul de limbă de *circulație electronică*. Aceasta sintagma, o parafrază la expresia *limba de circulație internațională*, încearcă să elimine antinomia, pe

cât de cunoscuta pe atât de goală în conținut spiritual și cultural, „limbi mari/limbi mici”. Conceptul de „limba de circulație electronică”, pe lângă semnificația lui directă, are profunde implicații culturale, sociale și nu în ultimul rând economice implicând dreptul fiecărui cetățean de a avea acces în propria limbă la cunoștințele, informațiile și serviciile ciberspaciului.

Promovarea limbii române în SI-SC presupune informatizarea limbii române ca factor infrastructural fundamental (vector funcțional) și precum și stimularea utilizării curente (prin vectori tehnologici) a limbii române în utilizarea tehnologiilor și a serviciilor informatice. Acestui obiectiv presupune un eforturi umane și materiale substanțiale și de dimensionarea lor se leagă orizontul de timp al realizării sale.

Volumul de față reunește lucrări ce tratează aspecte specifice prelucrării limbajului natural, în marea lor majoritate cu aplecare directă asupra limbii române. Inerent, volumul de față nu poate acoperi întreaga arie problematică a domeniului după cum nici reprezentarea specialistilor români în domeniul tehnologiei limbajului nu este completă, dar cititorul va găsi un larg evantai de direcții de cercetare, în care specialiștii români au obținut rezultate importante.

Volumul este structurat în patru părți (aspecte teoretice și probleme de terminologie, prelucrarea limbajului scris, prelucrarea limbajului vorbit, dezbateri și discuții) care pot fi citite în mod independent, în funcție de interesul specific al cititorului.

Prima parte „Lingvistică teoretică și formală; lexicografie” cuprinde 7 lucrări din domeniul lexicografiei, sintaxei și terminologiei.

În lucrarea „Resurse lingvistice elaborate la Institutul de Lingvistică *«Jorgu Iordan»*” Ioana Vintila Radulescu face o trecere în revista a celor mai importante resurse lingvistice realizate în cei peste 50 de ani de activitate la Institutul de Lingvistică *«Jorgu Iordan»*.

Angela Bidu-Vrâncăanu prezintă în lucrarea „Contribuția lingvisticii la studiul terminologiilor științifice” concluziile a trei contracte de cercetare științifică având ca obiect studiul terminologic al limbajului folosit în diverse domenii (matematică, filozofie, mineralogie, arte plastice).

Articolul „Gramaticile nontransformationale” al lui Emil Ionescu face o prezentare generală a gramaticilor bazate pe unificare și constrângeri precum și a principalelor realizări, în contextul acestei paradigme, în cercetarea lingvistică din România.

Nicolai Curteanu propune în lucrarea „Catre o teorie X-bar funcțională” o reconsiderare a teoriei clasice X-bar prin perspectiva modelului propriu SCD (Segmentare - Coeziune-Dependentă).

Ana-Maria Barbu prezintă în lucrarea sa „Teoria HPSG: studiu de caz: acordul încrucișat” principalele caracteristici ale teoriei HPSG și discută în acest context un caz de dependentă încrucișată specific limbii române, respectiv clauzele relative în care pronumele relativ este precedat de articolul genitival.

O serie de probleme legate de terminologia computationala sunt prezentate în ultimele doua lucrari ale primei sectiuni. În articolul „Dupa 10 ani de experienta terminologica: noul model de date terminologice al TERMROM” Dan Matei prezinta modelul dezvoltat în conformitate cu noile tendinte si standarde în domeniu si adoptat de Asociatia Româna de Terminologie – TERMROM.

Lucrarea lui Sorin Getaru „Probleme de reprezentare a datelor terminografice într-o baza de date relationala” aduce în discutie aspecte specifice reprezentarilor standardizate necesare realizarii dezideratului de interschimb si interoperabilitate între diverse tezaure terminologice si discuta elementele distinctive ale standardului ISO-12200 MARTIF (Machine-Readable Terminology Interchange Format).

Sectiunea a doua a volumului („Tehnologii ale limbajului scris”) este deschisa de lucrarea lui Dan Tufis si Dan Cristea „RO-BALKANET – ontologie lexicalizata în context multilingv pentru limba româna” care descrie stadiul dezvoltarii unui dictionar, pentru limba româna, structurat ca o retea semantica, de tip EuroWordNet, rezultat al unui program european ce-si propune extensia EuroWordnet (în prezent implementat pentru 10 limbi europene) cu înca 5 limbi.

Articolul lui Dan Gâlea, Nicolai Curteanu si Constantin Lintes „Algoritmi de segmentare a textului în unitati de tip clauzal” trateaza o problema delicata a prelucrării limbajului natural, respectiv cea a identificării, în raport cu un anumit criteriu functional, a structurilor „clauzale” si prezinta contrastiv doi algoritmi diferiti (unul dintre ei aparținând autorilor), atât prin prisma modelării lingvistice cât si al performantei computationale.

Rada Mihalcea si Vivi Nastase prezinta în articolul lor o metoda de inserare automata a caracterelor diacritice în texte scrise (cu studiu de caz pentru limba româna) fara diacritice si comenteaza rezultatele proprii în comparatie cu cele ale altor metode dezvoltate pentru rezolvarea aceleiasi probleme.

Adriana Vlad si Adrian Mitrea prezinta în lucrarea lor „Contributii privind structura statistica de cuvinte în limba româna scrisa” rezultate recente în caracterizarea statistica a limbii române scrise prin aproximarea ei ca un lant Markov ergotic multiplu cu ordin de multiplicitate mai mare decât 30, rezultate obtinute prin analiza riguroasa a unui corpus foarte mare de texte.

Articolul „Dezambiguizarea semantica automata în corpusuri paralele” al lui Dan Tufis prezinta o alternativa la spinoasa problema a dezambiguizării cuvintelor polisemantice, bazându-se pe extragerea cunostintele implicite existente într-un corpus multilingv (creat de traducatori profesioniști) si apelând la tehnici si euristici ale lingvisticii corpusului.

Dan Cristea prezinta în articolul „Referentialitate si cursivitate în structura discursului” elementele definiției ale teoriei sale asupra structurii discursive a textelor (teoria nervurilor) si își exemplifica argumentatia prin analiza dihotomiilor structura-referentialitate si structura-coerenta.

În lucrarea „DLIR - un sistem de cautare documentara multilingv” Amalia Todirascu prezinta o abordare bazata pe logici terminologice, ontologii si tehnici de

prelucrare a corpusurilor în implementarea unui sistem de regasire documentara bilingv (romtina si franceza).

Partea a doua a volumului se încheie cu articolul lui Stefan Trausan-Matu „Mediu hermenofor pentru asistarea învățării unor concepte într-o limba straina” care după o prezentare a notiunilor cu care operează în lucrare, descrie un modul de prelucrare a metaforelor utilizate în limbaje specializate (studiu de caz: limbajul financiar) încorporat într-un sistem de instruire inteligenta în învățarea conceptelor într-o limba straina, sistem distribuit dezvoltat în cadrul unui proiect european.

Sectiunea a treia a volumului este dedicata problemelor de prelucrare a vorbirii. Corneliu Burileanu si Luigi Bojan se opresc asupra tehnicilor de recunoastere a vorbitorului ca etapa distincta si strict necesara pentru recunoasterea automata a vorbirii si prezinta o parte a rezultatelor obtinute de catre autori.

Lucrarea lui Dragos Burileanu „Prelucrarea initiala a textului de intrare în cadrul unui sistem de sinteza a vorbirii pornind de la un text în limba româna” abordeaza problemele sintezei limbajului vorbit pornind de la un text în format electronic si detaliaza etapa de preprocesare a textului ca etapa primara în procesul transformarii sale în semnal vocal inteligibil si coerent.

Tot în domeniul sintezei vorbirii se plaseaza si lucrarea lui Horia Nicolai Teodorescu „Utilizarea tehnicilor nuanțate (fuzzy) și de dinamic? neliniar? pentru sinteza adaptiv? a vorbirii” ce subliniaz? rolul esențial al prozodiei și al model?rii sale algoritmice în realizarea unor sinteze vocale de calitate, purt?toare de informație emoțional?.

Un proiect de anvergura, este prezentat de Dumitru Todoroi, Diana Micusa, Zinaida Todoroi, Ion Linga, Ion Covalenco, Nicolae Objeleanu, Stefan Spataru, Stela Lungu, Virginia Turcanu, Elana Cozlov, Nadejda Ambrozii, Victor Slobodeanu, Igor Coseru si Catalina Suruceanu în lucrarea „Dictionarele multimedia ale limbii române. Secvente de implementari si experimentari ”.

Sectiunea a treia a volumului se încheie cu lucrarea elaborata de Silviu Bejinariu, Vasile Apopei si Mariana Roman „Mediu pentru editarea transcrierilor fonetice în Limba Româna. Realizarea Atlasului Lingvistic Român pe Regiuni” ce prezinta un instrument ce permite realizarea facila a transcrierilor fonetice într-un limbaj standardizat (IPA), ofera extensii specifice de adnotare fonetica (realizate pâna acum manual) si prefigureaza realizarea variantei computerizate a atlaselor lingvistice românești.

Ultima sectiune a volumului (Dezbateri si discutii) contine trei contributii. Prima dintre ele, elaborata de Academician Mihai Draganescu, „Asupra a doi vectori functionali ai Societatii Cunoasterii: Managementul Cunoasterii si Învatarea Electronica. Cultura si Societatea Cunoasterii” reprezinta liantul dintre volumul precedent (*Societatea Informatiionala – Societatea Cunoasterii. Concepte, solutii si strategii pentru România*, coordonator Academician Fl. Gh. Filip) si volumul de fata, rafinând clasificarea din lucrarea anterioara si adâncind o serie de probleme ridicate în [1].

Ultimele doua contributii reprezinta doua puncte de vedere asupra problematicii prelucrării limbajului natural, prima pozitie „Între lingvistica matematica si cea

computatională” fiind susținută de Academician Solomon Marcus, iar cea de a doua „Între lingvistica matematică și cea computațională: o altă perspectivă” fiind prezentată de Dan Tufis.

Mulumiri

Coordonatorii acestui volum, mulțumesc tuturor celor care au participat la realizarea proiectului „Strategii și soluții pentru societatea informațională-societatea cunoașterii în România” derulat în cadrul programului național INFOSOC. Mulțumiri speciale se cuvin directorului programului INFOSOC, Profesor Doina Banciu, care a susținut și a manifestat un interes deosebit față de desfășurarea acestui proiect.

Referințe bibliografice

- [1] M. Drăgănescu ”Societatea informațională și a cunoașterii. Vectorii societății cunoașterii” în *F. G. Filip (coord.) Societatea Informațională – Societatea Cunoașterii. Concepte, soluții și strategii pentru România*. Academia Română, Editura Expert, ISBN 973-8177-42-1, 2001, pp. 43-112
- [2] *** The Multilingual Information Society, Report of Commission of the European Communities, COM(95) 486/final, Brussels, November 1995.
- [3] *** Multilingualism in an Information Society, International Symposium organized by EC/DGXIII, UNESCO and Ministry of Foreign Affairs of the French Government, Paris 4-6 December 1997.
- [4] *** Promotion and Use of Multilingualism and Universal Access to Cyberspace, UNESCO 31st session, November 2001.
- [5] M. Dertouzos. “What It will Be “. Harper Edge. New York, 1977 (trad. în lb. Română “Ce va fi”, ed. Tehnica București, 2000).

SECTIUNEA I

**LINGVISTICA TEORETICA
SI FORMALA;
LEXICOGRAFIE**

Resurse lingvistice pentru limba româna elaborate la Institutul de Lingvistica „Iorgu Iordan”

Ioana VINTILA-RADULESCU
Institutul de Lingvistica „Iorgu Iordan – Al. Rosetti”
Bucuresti, Calea 13 Septembrie 13
e-mail: ioanar@fx.ro

1. Consideratii generale

Întelegând prin *resursa* în general o „rezerva sau sursa de mijloace (materiale sau spirituale) susceptibile de a fi valorificate într-o împrejurare data”¹, înțelegem prin *resurse lingvistice* pentru limba româna izvoarele fundamentale de informații cu privire la aceasta, stocate convenabil (chiar dacă încă preponderent în maniera tradițională) și care, în calitate de componente ale *culturii* în sensul cel mai larg, sunt susceptibile de a fi valorificate pentru studierea limbii române, precum și în diverse scopuri conexe, inclusiv aplicative, în cadrul *societății informatice* actuale.

Cât privește Institutul de Lingvistica „Iorgu Iordan”², acesta nu mai există formal ca atare, deoarece la începutul anului 2002, printr-o hotărâre de guvern adoptată la propunerea conducerii Academiei Române, s-a produs re-unirea sa și a Institutului de Fonetica și Dialectologie „Al. Rosetti”. (Spunem reunire întrucât cercetările de fonetica și de dialectologie formaseră inițial obiectul unui sector, respectiv al unei secții a Institutului de Lingvistica din București al Academiei Române (înființat în 1949), devenită din 1961 centru și apoi institut independent.) Întrucât în 1998 fusese oficializată, tot prin hotărâre de guvern, propunerea celor două institute, aprobată de Prezidiul Academiei, de a-și adăuga fiecare în titulatură numele fostului său director, institutul în cadrul căruia cele două nuclee care au fuzionat acum își continuă de fapt activitatea poartă numele dublu de Institutul de Lingvistica „Iorgu Iordan – Al. Rosetti”.

Fără îndoială, cele mai numeroase și mai importante resurse lingvistice pentru limba română s-au realizat la acum fostul Institut de Lingvistica „Iorgu Iordan”, înglobând, până în 1961 direct și apoi numai indirect, și contribuția colegilor foneticieni și

¹ [] *** (1975). Dictionarul limbii române (DLR). Serie nouă. Tomul IX, Litera R, București, s.v.

² [] Pentru o imagine de ansamblu asupra activității acestui institut și a istoriei sale v. Mioara Avram, Marius Sala, Ioana Vintila-Radulescu (coordonatori) (1999). Institutul de Lingvistica „Iorgu Iordan”. 50 de ani de existență (1949-1999), București.

dialectologi³, precum și, în unele cazuri, în colaborare cu alte institute de specialitate din țara ale Academiei – Institutul de Lingvistică și Istorie Literară „Sextil Puscariu” din Cluj și Institutul de Filologie Română „Alexandru Philippide” din Iași – și cu cadre didactice de la facultățile de profil mai ales ale Universității din București. Această activitate este continuată și în noul cadru organizatoric de sectoarele fostului institut, pe care în cele ce urmează îl vom numi, pe scurt, *Institutul*.

2. Resurse lexicografice

Dintre resursele lingvistice tradiționale dezvoltate până în prezent de Institut, cele mai importante din punctul de vedere care interesează aici sunt cele **lexicografice** – **dictionarele** (mono- și bilingve) –, activitatea lexicografică din Institut, începută încă de la înființarea sa, desfășurându-se din 1959 în cadrul unui sector specializat cu acest profil, condus până în 1985 de Mircea Seche, iar de atunci încolo de Ion Danaïla⁴.

2.1. Dictionare monolingve

2.1.1. Dintre dictionarele românești monolingve se distinge, prin anumite trăsături ale sale, dictionarul „explicativ general academic” intitulat pur și simplu **Dictionarul limbii române** – dar mai cunoscut ca „Dictionarul Academiei” –, a cărui realizare se apropie de sfârșit și care va cuprinde o mare parte a „tezaurului” lexical al limbii române – fără a putea și nici a intenționa să includă însă ansamblul cuvintelor românești folosite în toate epocile, în toate regiunile și în toate domeniile⁵. În ciuda marilor sale calități, care sunt bine cunoscute și asupra cărora nu credem deci că mai este nevoie să insistăm aici, acest dictionar prezintă un dezavantaj major din punctul de vedere al utilizării sale ca resursă de bază (pe lângă faptul că nu se prezintă și sub forma unei variante electronice, care nici nu putea fi imaginată până nu de mult) și anume caracterul său fatalmente neunitar, datorat faptului că a fost elaborat pe parcursul a aproape un secol⁶, de unde marile deosebiri dintre cele două părți ale sale: cea publicată între 1907 și 1949 sub conducerea marelui lingvist Sextil Puscariu și cea care a început să apară din 1965 și a cărei publicare se apropie, în fine, de sfârșit. „Seria veche” a dictionarului academic, desemnat de aceea prin sigla DA,

³ *Aceștia au produs mai ales „resurse” de un tip specializat, concretizate în principal în atlase lingvistice și în arhiva fonogramică a limbii române, de care nu ne vom ocupa în mod direct aici, dar care, ca și contribuțiile similare ale altor institute, au avut și un aport indirect la resursele fundamentale despre care vorbim, printre izvoarele cărora s-au numărat.*

⁴ [] Pentru detalii cu privire la lucrările acestuia v. Ion Danaïla (1999). Sectorul de lexicologie și lexicografie, în Mioara Avram, Marius Sala, Ioana Vintila-Radulescu, *op. cit.*, p. 98-113.

⁵ [] *Ideea, relativ utopică și controversată, a înregistrării și chiar a descrierii semantice a întregului inventar lexical al limbii române (ILEX) din toate timpurile, incluzând atât numele comune, cât și cele proprii (v. Ion Danaïla (1993)). Pentru un inventar general al limbii române, în „Limba română” XLII, nr. 2, p. 61-68), nici nu a început a fi pusă în practică.*

⁶ [] *V., printre altele, Marius Sala (1999). Institutul de Lingvistică „Iorgu Iordan” la 50 de ani, în Mioara Avram, Marius Sala, Ioana Vintila-Radulescu, op. cit., p. 35-37.*

cuprinde literele *A-C* (inclusiv putinele neologisme scrise acum cu *k-*, iar în *DA* cu *ch-*) și *F-J* complet, iar literele *D* și *L* partial (pâna la cuvântul *de*, respectiv *lojnita*), totalizând 3.142 de pagini de tipar, format mare, dintre ele lipsind în întregime, după cum se observa, litera *E*. Aceasta prima jumătate a dicționarului se distinge prin lista de cuvinte, bogată mai ales sub aspectul fondului tradițional, prin tratarea amanunțită a semantismului, bazată pe numeroase citate, prin dimensiunile și valoarea comentariului etimologic, precum și prin traducerea sensurilor în limba franceză⁷. Desigur, nu aveau cum figura în aceste prime volume numeroasele neologisme încetățenite în românește după elaborarea ei, ilustrarea sensurilor prin utilizarea lor de către autori mai noi și în general toate aspectele care sunt rodul evoluției ulterioare a limbii române, al cercetărilor dialectale, etimologice, filologice etc. mai recente și al dezvoltării lingvisticii și metodelor ei, în general. Din 1965 dicționarul și-a reînceput apariția, în format asemănător, ca *Serie nouă* (de data aceasta sub o sigla diferită, menționată în titlu, *DLR*), cu litera *M*, sub conducerea, la început, a lui Iorgu Iordan, Alexandru Graur și Ion Coteanu, iar actualmente a lui Gh. Mihaila și Marius Sala. Noua serie pastrează, în mare, principiile lui Sextil Puscariu, dar beneficiază de toate avantajele elaborării sale mai aproape de zilele noastre: ea include modificări și amplificări reflectând evoluția limbii române, a lexicografiei românești și a studiului limbii române în ansamblu, precum și a lingvisticii în general, dar nu mai cuprinde, în schimb, traducerea sensurilor (în anii '60 nefiind considerat oportun acest lucru, deși era util mai ales pentru cunoașterea limbii române de către străini, fără a fi, este drept, uzual într-un dicționar monolingv explicativ), iar secțiunea etimologică a fost redusă, dicționarul păstrându-și însă caracterul istoric (sensurile sunt date în ordinea atestării lor în texte și în alte surse)⁸. Institutul bucurestean a redactat literele *M, N, P, S* și *Z*⁹ și este pe cale de a încheia reluarea și terminarea literei *D* absentă din prima parte (trei volume); numai primele patru litere elaborate la București însumează 51.847 de cuvinte și variante, totalizând 5.839 p. Institutului din Cluj i-au revenit literele *O, R, T, T* (totalizând 2.044 de pagini de tipar), *U* (aflată sub tipar) și, din prima parte, reluarea și terminarea unei părți din litera *L*, iar celui din Iași – literele *S, V* (*S* și prima parte din cele trei ale literei *V* – singura dintre acestea apărută până acum – totalizând 599 de pagini de tipar), *W, X, Y*, precum și, din prima parte, elaborarea literei lipsa *E* și reluarea și terminarea unei părți din litera *L*; pentru etimologii au fost consultați specialiști din mai multe centre universitare. Majoritatea literelor au

⁷ [] Pentru o descriere amanunțită a *DA* v. *Mircea Seche (1969)*. Activitatea lexicografică a lui Sextil Puscariu, în *Schita de istorie a lexicografiei române, vol. II*, De la 1880 până astăzi, *București*.

⁸ [] *V. și Mircea Seche (1969)*. *Seria nouă a Dicționarului academic general, în Schita de istorie a lexicografiei române, vol. II*, De la 1880 până astăzi, *București*, p. 72-79.

⁹ [] Iorgu Iordan, Al. Graur, I. Coteanu (red. resp.) et al. (1965-2000). *Dicționarul limbii române (DLR)*. *Serie nouă*, București: T. VI, *Litera M*, 1965-1968 (apărut inițial în fascicule); VII, *Partea 1, Litera N*, 1971; *Partea a 2-a, Litera O*, 1969; VIII, *Litera P*. *Partea 1, P-PAZUI*, 1972; *Partea a 2-a, PE-PÎNAR*, 1974; *Partea a 3-a, PÎNA-POGRIBANIE*, 1977; *Partea a 4-a, POGRIJENIE-PRESIMTIRE*, 1980; *Partea a 5-a, PRESIN-PUZZOLANA*, 1984; IX, *Litera R*, 1975; X. *Litera S*. *Partea 1, S-SCLABUC*, 1986; *Partea a 2-a, SCLADA-SEMÎNTARIE*, 1987; *Partea a 3-a, SEMN-SÎVEICA*, 1990; *Partea a 4-a, SLAB-SPONGHIOS*, 1992; *Partea a 5-a, SPONGIAR-SWING*, 1994; XI *Partea 1, Litera S*, 1978; *Partea a 2-a, Litera T, T-TOCALITA*, 1982; *Partea a 3-a, TOCANA-TWIST*, 1983; XII, *Partea 1, Litera T*, 1994; XIII, *Partea 1, Litera V, V-VENI*, 1997; XIV, *Litera Z*, 2000.

aparut, unele pe sarite (*M* între 1965 și 1968, *N* în 1971, *O* în 1969, *P* între 1972 și 1984, *R* în 1975, *S* între 1986 și 1994, *S* în 1978, *T* în 1983, *T* în 1994, prima parte din *V* (pâna la *a veni*) în 1997 și *Z* în 2000) – în total 20 de volume –, cu excepția literelor *D, E, K, L, U*, a putinelor cuvinte începând cu litera *Q* și a ultimelor parti ale literei *V* (începând cu *venin*), la care se adauga literele *W, X* și *Y*. Deosebirea cea mai importanta consta în tipurile de cuvinte reprezentate în cele doua serii: la majoritatea primelor litere ale alfabetului (cu excepția celor care s-au redactat abia acum), neologismele sunt slab reprezentate, nu numai din cauza faptului ca foarte multe nici nu se încetatenisera înca în limba româna la vremea elaborării volumelor respective, dar și din cauza reticentei lui Puscariu cu privire la acest sector al vocabularului; într-o situatie asemanatoare se afla termenii regionali, deoarece cercetarile dialectale se aflau în acea vreme abia la început. Prima parte prezinta în schimb avantajul de a putea servi ca baza pentru o prelucrare bilingva, întrucât includea și traducerea sensurilor în limba franceza, la care a trebuit sa se renunte în perioada comunista. Reluarea și completarea acestui dictionar, absolut necesara, nu ni se mai pare astazi recomandabil și nici posibil de realizat prin mijloace traditionale (fise etc.), ci exclusiv pe baze informatizate. Ea ar trebui sa valorifice, printre altele, și banca de texte și cea de inovatii a limbii române, despre care va fi vorba mai departe. Ar fi necesar ca partea publicata înainte de 1949 sa fie reluata și adusa la zi, cu atât mai mult cu cât putine persoane și chiar biblioteci posedau dictionarul în întregime (chiar în cazul seriei noi, tirajele diferitelor litere au fost diferite și în continua scadere), iar îmbatrânirea hârtiei în cazul seriei vechi o face fragila și greu de consultat. Având în vedere ca pentru noua serie a dictionarului s-au adunat, manual, peste sase milioane de fise cu extrase și atestari (dintre acestea, în DLR au fost incluse cca 3.200.000 de citate¹⁰, reprezentând aproximativ 88% din totalul textului), este de sperat ca la reluarea, într-un viitor mai mult sau mai puțin apropiat, se va putea uza de avantajele elaborării computerizate, valorificându-se bancile de date în curs de elaborare în institut, despre care va fi vorba mai departe.

Având în vedere diferentele semnalate (dintre care unele se regasesc și între primele și ultimele litere din seria noua), este foarte binevenita ideea actualilor responsabili ai DLR de a se publica, pentru operativitate, un *Supliment* – care se poate realiza relativ mai lesne – „care sa înregistreze neologismele adoptate de limba literara de la începutul secolului” 20 „pâna în prezent, precum și o serie de cuvinte regionale incluse în atlasele lingvistice și în culegeri de pe teren sau termeni vechi extrasi din documente ale secolelor al XVI-lea – al XVIII-lea, editate în ultimele decenii”¹¹.

2.1.2. Din motivele expuse mai sus, la care se adauga și faptul ca DA/DLR este accesibil mai ales specialistilor și mai puțin publicului larg, institutul bucurestean pregateste între timp, la sugestia conducerii Academiei Române, o sinteza a marelui dictionar academic, fara citate și izvoare și cu un sistem foarte economic de prezentare a

¹⁰ [] În legatura cu reflectarea noilor norme ortografice ale limbii române în volumele DLR elaborate dupa 1993, semnalăm faptul ca forma sânt, reflectând un fonetism real, vechi și popular, este pastrata în citatele în care nu era folosit sunt.

¹¹ [] Marius Sala, G. Mihaila (2000). Cuvânt înainte, în Dictionarul limbii române (DLR). Serie noua. Tomul XIV. Litera Z, Bucuresti, p. VI.

informatiilor lexicografice. Acest *Mic dictionar academic* (MDA)¹² (care va avea totusi patru volume), inclus, alaturi de DLR, printre lucrarile fundamentale ale Academiei Române, va avea cca 175 000 de intrari (cc125 000 de cuvinte si cca 50.000 de variante); primul volum (A-C) a fost publicat în anul 2001 de editura Univers Enciclopedic. Proiectul *Micului dictionar academic*, numit astfel în opozitie cu „marele” dictionar academic, si-a propus sa reduca decalajul dintre cele doua serii ale acestuia, îmbogatinind primele litere pe baza unor surse lexicografice mai noi. La rândul sau, acest nou dictionar prezinta însa dezavantajul de a fi fost obligat, prin dimensiuni, sa renunte la citatele ilustrative, ceea ce limiteaza posibilitatea utilizarii lui ca sursa de informatii morfologice, gramaticale si stilistice; numarul neobisnuit de mare de abrevieri ne transparente, utilizate din acelasi motiv de economie, constituie un argument suplimentar în favoarea realizarii unei versiuni electronice a MDA care sa permita regasirea automata a informatiilor.

2.1.3. Spre deosebire de DA/DLR, o reflectare în general unitara a vocabularului limbii române ofera *Dictionarul explicativ al limbii române*¹³, despre a carui sigla, DEX, se afirma, pe drept cuvânt, ca a devenit un apelativ; denumirea, care ar fi trebuit protejata prin înregistrare, a fost preluata abuziv de *Noul dictionar explicativ al limbii române* publicat pe CD-Rom de firmele Litera în sigla NODEX, sugerând ca ar fi „un nou DEX”. Prima editie, un volum de 1.049 de pagini, cuprinzând 56.569 de cuvinte si variante, a fost urmata de un *Supliment la Dictionarul explicativ al limbii române (DEX-S)*¹⁴. Editia a doua a DEX¹⁵ totalizeaza 1.204 pagini; aceasta editie, care se publica în continuare în tiraje succesive, totalizase numai în primii patru ani de la aparitie 65.000 de exemplare vândute, dupa un calcul sumar rezultând ca la 42 de locuitori ai României revenea un DEX. Actualmente, se poate într-adevar afirma ca, prin DEX, *best-sellerul* lingvisticii românesti, Institutul a intrat în marea majoritate a caselor din România. Se preconizeaza ca DEX sa fie realizat, în fine, într-un viitor relativ apropiat, si în format electronic. El a fost deja supus, de catre Centrul de Cercetari Avansate în Învatarea Automata, Prelucrarea Limbajului Natural si Modelarea Conceptuala al Academiei Române, codificarii conform TEI¹⁶. Se estimeaza ca editia a III-a a DEX, conceputa sub conducerea lui Ion Danaila, va avea în plus fata de precedenta cca 30 000 de cuvinte. Sub conducerea lui Ion Coteanu si Ion Danaila, la sectorul de specialitate al Institutului a fost conceput si un *Nou dictionar explicativ al limbii române* (NEX), cu caracteristici diferite de cele ale DEX: inventar de

¹² [] V. I. Danaila (1994). De ce este nevoie de un MDA?, în „Limba româna” XLIII, p. 397-406 si Marius Sala (2001). Prefata, în *Micul dictionar academic (MDA)*, vol. I, A-C, Bucuresti.

¹³ [] I. Coteanu, Luiza Seche, M. Seche (conducatorii lucrarii) et al. (1975, 1996). *Dictionarul explicativ al limbii române (DEX)*, Bucuresti.

¹⁴ [] Ion Coteanu, Ion Danaila, Nicoleta Tiugan (conducatorii lucrarii) et al. (1988). *Supliment la Dictionarul explicativ al limbii române (DEX-S)*. Bucuresti.

¹⁵ [] Ion Coteanu, Lucretia Mares (sub conducerea) et al. (1996). *Dictionarul explicativ al limbii române (DEX)*, editia a II-a, Bucuresti.

¹⁶ [] Dan Tufîs (2000). Cercetare si colaborare internationala în ingineria lingvistica la RACAI, în „Terminologia în România si în Republica Moldova”, Cluj-Napoca, p. 34-36 si Recherche et collaboration internationale en industries de la langue r l'Académie Roumaine, în „Terminometro Hors-série n° 4. La terminologie en Roumanie et en République de Moldova”, p. 38-40.

cca 100.000 de cuvinte si variante (deci aproape de doua ori mai multe decât prima editie a DEX), definitii mai concise, prin eliminarea sinonimelor si – din pacate!–, neincluderea etimologiei cuvintelor; revizuit de cei doi responsabili, el asteapta introducerea în calculator, în vederea efectuării corelărilor semantice definicionale si sinonimice.

2.1.4. DEX a scos practic din circulatie dictionarele explicative mai vechi, limitate la limba româna literara, DLRLC si DM¹⁷. Prima sigla reprezinta **Dictionarul limbii române literare contemporane**¹⁸, elaborat de institutele din Bucuresti si Cluj pornind de la „baza manuscrisa” a DA si aparut între 1955 si 1957 în patru volume. El se mai foloseste si astazi – desi din el lipsesc cuvintele, sensurile si citatele neconforme cu ideologia vremii – pentru citatele cu care, spre deosebire de dictionarele de dimensiuni comparabile mai noi, sunt ilustrate sensurile cuvintelor (chiar daca, pentru unele neologisme, citatele provin, asa cum era obligatoriu în epoca, din traducerile „operelor clasicilor” marxism-leninismului!). Dintre acestea, primul mai merita însa atentie în virtutea faptului ca, spre deosebire de DEX si de MDA, include citate ilustrative, care din pacate au fost eliminate din dictionarele urmatoare.

2.1.5. O versiune prescurtata a acestui dictionar, cu un inventar putin marit si cu adaugarea etimologiei cuvintelor, dar cu eliminarea citatelor, a fost publicat de Institutul din Bucuresti în 1958 sub titlul **Dictionarul limbii române moderne**¹⁹ (abreviat DM).

2.1.6. Un dictionar de un tip special, cu o utilitate mult mai larga decât aceea care i se recunoaste de obicei, elaborat de data aceasta de colectivul de gramatica al Institutului (condus pâna de curând de Mioara Avram²⁰), este **Dictionarul ortografic, ortoepic si morfologic al limbii române (DOOM)**²¹. Este singurul dictionar al limbii române (mai bogat decât DEX₁) care contine ample informatii cu privire la formele flexionare ale cuvintelor variabile incluse, putând servi astfel (chiar daca aceste informatii nu sunt exhaustive) ca sursa pentru studii si aplicatii de morfologie. Institutul are în prezent în lucru, sub conducerea subsemnatei, o a doua editie, partial revazuta si adaugita, a DOOM (care va cuprinde si cuvinte neînregistrate în nici un dictionar românesc pâna în prezent). Aceasta va aparea în anul 2003, inclusiv pe CD-Rom, si va trebui sa serveasca drept baza unui nou corector ortografic si morfologic, care sa tina seama de modificarea unor recomandari oficiale în raport cu cele înca în vigoare.

¹⁷ [] Pentru detalii cu privire la aceste doua dictionare v. *Mircea Seche (1969)*. Dictionarele explicative ale limbii române literare, în *Schita de istorie a lexicografiei române, vol. II, De la 1880 pâna astazi, Bucuresti, p. 135-147*.

¹⁸ [] *D. Macrea, E. Petrovici (sub directia) et al. (1955-1957)*. Dictionarul limbii române literare contemporane (DLRLC), *Editura Academiei, Bucuresti, vol. I, A-C; II, D-L, 1956; III, M-R, 1957; IV, S-Z, 1957*.

¹⁹ [] *D. Macrea (sub directia) (1958)*. Dictionarul limbii române moderne, *Bucuresti*.

²⁰ [] Pentru activitatea acestuia v. *Mioara Avram (1999)*. Colectivul de gramatica, în *Mioara Avram, Marius Sala, Ioana Vintila-Radulescu, op. cit., p. 113-125*.

²¹ *Mioara Avram (red. resp.) et al. (1982)*. Dictionarul ortografic, ortoepic si morfologic al limbii române (DOOM), *Bucuresti, 1982*.

2.1.7. În fine, un dicționar mai puțin obisnuit, *Dicționarul invers*²², în care cuvintele sunt ordonate alfabetic pornind dinspre sfârșitul lor, este deosebit de util specialistilor pentru studierea terminațiilor, a desinentelor și a sufixelor, dar și poetilor, fiind utilizabil și ca dicționar de rime. Aceasta lucrare – care, spune „legenda”, a valorificat experiența din copilarie a uneia dintre autoare, care folosise în joacă o *pasareasca* de acest fel – ar merita și ea o nouă elaborare, pe baza unui inventar mai bogat și actualizat de cuvinte și a unui program care să permită „rasturnarea” lor automată.

2.1.8. Institutul a publicat, încă din 1968, un dicționar al lexicului unui autor, primul ales neputând fi altul decât Eminescu – *Dicționarul limbii poetice a lui Eminescu*²³, care însă, la acea vreme, nu se putea baza, evident, pe stabilirea concordanțelor așa cum se realizează ea în zilele noastre.

2.1.9. Institutul a elaborat de asemenea o serie de dicționare ale limbii române pe epoci sau pe probleme, cum sunt *Dicționarul limbii române literare vechi*²⁴ și *Dicționarul împrumuturilor latino-romance în limba română veche*²⁵, publicate de sectorul de limbă literară, filologie și poetică²⁶, condus de Ion Ghetie, iar în prezent de Alexandru Mares – și *Dicționarul elementelor românești din documentele slavo-române*²⁷, elaborat la sectorul de slavistică²⁸ – dicționare destinate în primul rând specialistilor.

2.1.10. Un cercetător din institut, Constant Maneca, a publicat, împreună cu Florin Marcu, un extrem de util, cu toate criticile care i s-au adus, *Dicționar de neologisme*²⁹, reluat și dezvoltat, după moartea celui dintâi, de Florin Marcu, în numeroase variante, de diverse dimensiuni, la diferite edituri, inclusiv pe CD-Rom.

2.1.11. Se afla în lucru și *Dicționarul etimologic al limbii române* (DELR) – coordonator: Marius Sala –, altă lucrare fundamentală a Academiei Române, la care colaborează cercetători din toate sectoarele Institutului, cercetători din Cluj și Timișoara și cadre didactice de la universitățile din București, Cluj și Timișoara.

²² [] *** (1957). *Dicționar invers*, București. V. și Mircea Seche (1969). *Schita de istorie a lexicografiei române*, vol. II, De la 1880 până astăzi, București, p. 254-255.

²³ [] Tudor Vianu (sub redacția) et al. (1968). *Dicționarul limbii poetice a lui Eminescu*, București.

²⁴ [] Mariana Costinescu, Magdalena Georgescu, Florentina Zgraon (1987). *Dicționarul limbii române literare vechi (1640-1780). Termeni regionali*, București.

²⁵ [] Gh. Chivu, Emanuela Buza, Alexandra Roman Moraru (1992). *Dicționarul împrumuturilor latino-romance în limba română veche (1421-1760)*, București.

²⁶ [] V. Ion Ghetie (1999). *Colectivul de limbă literară și filologie, în Mioara Avram, Marius Sala, Ioana Vintila-Radulescu, op. cit., p. 132-143.*

²⁷ [] Gheorghe Bolocan (redactor responsabil) et al. (1981). *Dicționarul elementelor românești din documentele slavo-române. 1374-1600*, București.

²⁸ [] *Cu privire la care v. Virgil Nestorescu (1999). Sectorul de lexicografie bilingvă. Fostul sector de slavistică, în Mioara Avram, Marius Sala, Ioana Vintila-Radulescu, op. cit., p. 165-174.*

²⁹ F. Marcu, C. Maneca (1961-1978). *Dicționar de neologisme*, București, 1961; ed. II revizuită și adăugită, 1966; ³1978. V. și Mircea Seche (1969). *Schita de istorie a lexicografiei române*, vol. II, *De la 1880 până astăzi*, București, p. 154-159.

2.1.12. Pe lângă resursele privitoare la numele comune, Institutul a elaborat și importante lucrări consacrate numelor proprii³⁰.

Astfel, în domeniul toponimiei, după clasică lucrare a lui Iorgu Iordan³¹, s-a realizat în Institut *Dictionarul toponimic al României*, partea I, *Oltenia*³², elaborat sub conducerea lui Gh. Bolocan în colaborare cu cadre didactice de la Universitatea din Craiova, din care au apărut în perioada 1993-2001 primele trei volume, precum și al doilea dictionar din serie, consacrat *Munteniei* și aflat în curs de definitivare.

În domeniul onomasticii, de asemenea urmând altei lucrări clasice a lui Iorgu Iordan³³, Institutul colaborează și la proiectul internațional PatRom, care realizează un dictionar istoric de antroponomie romanică, în care este reprezentată și limba română, și din care până acum a fost publicat un prim volum de prezentare³⁴.

2.2. Dictionare bilingve și multilingve

2.2.1. Pe lângă dictionarele monolingve ale limbii române, Institutul a realizat și unele din cele mai importante dictionare bilingve³⁵ (englez-român³⁶, german-român³⁷, rus-român³⁸, ceh-român³⁹ și sârb-român⁴⁰ – perechea sa, dictionarul român-sârb, fiind în curs de redactare; un dictionar francez-român a ramus nepublicat) și frazeologice (spaniol-român, sub tipar, și român-spaniol, în curs de elaborare), cărora li se adaugă dictionare bilingve⁴¹ – care au început a fi transpuse și pe CD-Rom – și dictionare frazeologice

³⁰ [] Pentru activitatea în acest domeniu v. Gheorghe Bolocan, Ecaterina Mihaila (1999). Colectivul de onomastica și Domnita Tomescu (1999). Grupul de lucru PatRom, în *Mioara Avram, Marius Sala, Ioana Vintila-Radulescu, op. cit., p. 125-132*.

³¹ [] Iorgu Iordan (1952-1963). *Nume de locuri românești în Republica Populăra Română* București, 1952; *Toponimia românească*, București, 1963.

³² [] Gh. Bolocan (sub redacția) et al. (1993-2001). *Dictionarul toponimic al României. Oltenia* (DTRO), vol. I-III, Craiova, Editura Universitaria.

³³ [] Iorgu Iordan (1983). Dictionar al numelor de familie românești, București, *Editura Academiei. DE VERIF.*

³⁴ *** (1997). Dictionnaire historique d'anthroponymie romane (PatRom). Présentation d'un projet, Tübingen.

³⁵ [] V. și Ilinca Constantinescu. (1999). Fostul sector de germanistică, în *Mioara Avram, Marius Sala, Ioana Vintila-Radulescu (coordonatori) (1999), op. cit., p. 174-179*.

³⁶ [] L. Levitchi (red. resp.) et al. (1974). *Dictionar englez-român*, București. Suplimentul la acest dictionar, care nu a mai apărut, coordonat de Ilinca Constantinescu, va fi inclus într-o nouă ediție, mult marită, a dictionarului, aflata sub tipar și care va reprezenta cel mai bogat dictionar englez-român.

³⁷ [] M. Isbăscu, Maria Iliescu (coord. și revizie) et al. (1966, 1988). Dictionar german-român, București, 1966; ediția a II-a revăzută și îmbogățită, București, 1988.

³⁸ [] Gheorghe Bolocan (redactor responsabil) (1964). Dictionar rus-român, București.

³⁹ [] S. Stati (red. resp.) et al. (1967). Dictionar ceh-român, București.

⁴⁰ [] M. Tomici (1998-2000). Dictionar sârb-român, 3 vol., Timisoara.

⁴¹ [] Gh. Bolocan (1972). *Dictionar bulgar-român*, București – Sofia; Gh. Bolocan et al. (1980). *Dictionar român-rus*, București – Moscova; Al. Calciu, C. Duhaneanu, D. Munteanu (1979). *Dictionar român-spaniol*, București; Ana Canarache (coord.) (1967, 1978). *Dictionar român-francez*, București, 1978; M. Isbăscu (red. resp.) (1963), *Dictionar român-german*, București; Valeria

românești⁴² și bilingve⁴³ elaborate de unii membri ai Institutului; *Dictionarul elen-român*, lucrare colectivă, se apropie și el de sfârșit.

2.2.2. Institutul a colaborat și la mai multe dicționare multilingve⁴⁴, dintre care ase distinge în mod deosebit un lexicon multilingv de un tip special – o adevărată premieră internațională – este *Dictionarul elementelor latinești savante din limbile romanice*, elaborat la sectorul de romanistică (condus inițial de marele romanist Iorgu Iordan, apoi de Marius Sala și în prezent de subsemnata)⁴⁵, în colaborare cu cadre didactice de la Facultatea de Limbi și Literaturi Straine a Universității din București și în coordonarea prof. dr. Sanda Reinheimer Rîpeanu, decanul Facultății: Negăsindu-și un editor „clasic” din cauza costurilor prea ridicate, acest dicționar va fi publicat direct pe Internet, sub auspiciile Universității din București.

3. Banci de date

3.1. Institutul a avut în proiect încă din anii 1978-80 realizarea primei banci computerizate de date lingvistice din România (*Banca de date fono-morfo-semantice a limbii române – BANDASEM*)⁴⁶, cel dintâi modul fiind cel de semantică, proiectat pentru un *Dictionar confruntativ de sinonime, de analogii și de asociații al limbii române (DCSAAs)*. Redactarea acestuia, care a ajuns la litera S, s-a făcut însă cu mijloace tradiționale, deși prin colaborarea cu Centrul de Calcul al Universității din București se elaborase un modul de program în sistemul Socrate pentru recunoașterea și selectarea, ca probă, a analogiilor și a asociațiilor cuvântului *blitz*. Elaborarea DCSAAs a fost întreruptă pentru un timp în favoarea lucrărilor prioritare al Academiei, iar reluarea lui se va putea face, sperăm, cu mijloacele informatice disponibile actualmente⁴⁷.

3.2. O minibanca inițiată în cadrul sectorului de gramatică al Institutului, a cărei alimentare a fost din păcate întreruptă în favoarea concentrării forțelor pentru realizarea

Neagu (2001). *Dictionar român-spaniol* (cu transpunere pe CD-Rom), București.

⁴² [] V. Breban et al. (1969). *Dictionar de expresii și locuțiuni românești*, București.

⁴³ [] Gh. Bolocan et al. (1968). *Dictionar frazeologic rus-român*, București; H. Mantsch et al. (1979). *Dictionar frazeologic român-german*, București.

⁴⁴ *** (1981). *Dictionnaire de la presse écrite et audiovisuelle. Espagnol-français-italien-portugais-roumain*, Paris; *** (2001). *Usage Dictionary of Anglicisms in Selected European Languages (UDASEL)* Oxford s.a.

⁴⁵ [] *Cu privire la activitatea acestuia v. Marius Sala (1999). Sectorul de limbi romanice și clasice, în Mioara Avram, Marius Sala, Ioana Vintila-Radulescu (coordonatori) (1999), op. cit., p. 147-164.*

⁴⁶ [] Ion Dănilă (2000). *Proiecte de prelucrare electronică a vocabularului limbii române, în „Terminologia în România și în Republica Moldova”, Cluj-Napoca, p. 36-37.*

⁴⁷ [] *Partea de fonetică/grafematică și de morfologie a BANDASEM a fost cedată institutului omolog din Cluj, pentru care v. Felicia Serban et al. (2000). Baza de date a limbii române, în „Terminologia în România și în Republica Moldova”, Cluj-Napoca, p. 37-38 și La base de données de la langue roumaine, în „Terminometre Hors-série n° 4. La terminologie en Roumanie et en République de Moldova”, p. 40-42. VERSIUNEA INTEGRALĂ ÎN TUFIS...*

editiei a doua a „Gramaticii Academiei”, este **Banca de inovatii a limbii române**, bazata pe monitorizarea presei scrise si audiovizuale actuale.

3.3. Având în vedere ca în DOOM informatia este atomizata, în folosul cititorului neprofesionist, în cadrul fiecarui cuvânt-titlu în parte, dar este greu de sistematizat de catre specialist, Institutul are în proiect, începând din 2003, realizarea unui baze de date care sa permita nu numai elaborarea unui **Nou dictionar ortografic, ortoepic si morfologic al limbii române** si a unor dictionare specializate de un tip asemanator, precum si aducerea lor permanenta la zi, ci si gruparea cuvintelor în clase în functie de caracteristicile lor fonetice, grafice si morfologice⁴⁸.

3.4. Institutul are în proiect si elaborarea sau definitivarea unor resurse terminologice⁴⁹ (dictionare terminologice bi- si multilingve, valorificând cele elaborate în cadrul proiectului PRACTEAST din cadru programului COPERNICUS al Comisiei Europene⁵⁰ si un dictionar al termenilor oficiali); de altfel, mai multi membri ai Institutului au colaborat la realizarea **Bancii de date terminologice** (BDT) multilingve a Asociei Române TermRom⁵¹, care, cu sprijinul Directiei de terminologie si inginerie lingvistica a Uniunii Latine, este accesibila pe site-ul TermRom gazduit de CIMEC (<http://www.cimec.ro/tr>) si, de curând, si pe CD-Rom. Reprezentarea României (prin subsemnata) în Reteaua Panlatina de terminologie (Realiter)⁵² si în Reteaua Francofona de

⁴⁸ [] *Clasificarea cuvintelor românesti conform modului lor de flexiune, realizata de Alf Lombard, Constantin Gădei (1981). Dictionnaire morphologique de la langue roumaine, Lund – Bucuresti, bazata pe inventarul DEX₁, prezinta unele inexactitati din cauza insuficientei cunoasteri de catre autori a limbii române actuale; ea constituie una din bazele realizarii, în Republica Moldova, a unui pachet de programe destinat elaborarilor de nivel morfologie, pentru care v. Elena Boian et. al. (2000). Instrumentar pentru aplicatii lingvistice, în „Terminologia în România si în Republica Moldova”, Cluj-Napoca, p. 38-40 si Instruments pour applications linguistiques, în „Terminometro Hors-série n° 4. La terminologie en Roumanie et en République de Moldova”, p. 42-44; o grupare pe tipuri a unui numar limitat de cuvinte ale limbii române a fost realizata de Flora Suteu, Elisabeta Sosa (1999) în Îndreptar ortografic si morfologic, Bucuresti.*

⁴⁹ [] *V. Ioana Vintila-Radulescu (1999). Institutul de Lingvistica „Iorgu Iordan” din Bucuresti, în „Terminologia în România si în Republica Moldova”, Cluj-Napoca, p. 13-15, si L’Institut de Linguistique Iorgu Iordan de Bucarest, în „Terminometro Hors-série n° 4. La terminologie en Roumanie et en République de Moldova”, p. 22-13.*

⁵⁰ [] *Nicoleta Petuhov. (2000). Colaborarea româneasca la proiectul PRACTEST, în „Terminologia în România si în Republica Moldova”, Cluj-Napoca, p. 58-59 si La collaboration roumaine au projet Practeast, în „Terminometro Hors-série n° 4, La terminologie en Roumanie et en République de Moldova”, p. 64-66.*

⁵¹ [] *Dan Matei (2000). Banca de date terminologice a TermRom, în „Terminologia în România si în Republica Moldova”, Cluj-Napoca, p. 29-30 si La banque de données terminologiques de TermRom, în „Terminometro Hors-série n° 4, La terminologie en Roumanie et en République de Moldova”, p. 32-33.*

⁵² [] *Dan Matei (2000). Prezenta româneasca în reseaua panlatina de terminologie (Realiter), în „Terminologia în România si în Republica Moldova”, Cluj-Napoca, p. 56-58 si La présence roumaine dans le Réseau panlatin de terminologie Realiter, în „Terminometro Hors-série n° 4. La terminologie en Roumanie et en République de Moldova”, p. 63-64.*

Amenajare Lingvistica (Rifal)⁵³ vor constitui desigur un sprijin în dezvoltarea resurselor terminologice pentru limba română în conformitate cu normele și recomandările internaționale.

În afara numelor comune, și numele proprii au constituit obiectul preocupărilor institutului și ale unor membri ai săi.

4. Corpusuri

O altă categorie importantă de resurse lingvistice o constituie corpusurile, la Institut fiind în curs de realizare o **Banca de texte românești**, care cuprinde texte din secolele al XVI-lea – al XVIII-lea, introduse integral în calculator, și în care se prevede introducerea câtorva sute de texte din toate epocile. Inițiată de directorul institutului, acad. Marius Sala. Banca, a fost deja valorificată în elaborarea unor teze de doctorat, printre altele la aceea a Janei Balacciu-Matei. Pentru exploatarea ei deplină în vederea identificării primelor atestări ale cuvintelor limbii române din fondul vechi, necesare MDA și *Dictionarului etimologic al limbii române* (DER) (în curs de elaborare sub conducerea acad. Marius Sala), a îmbogățirii dictionarelor limbii române în general și a dezvoltării studiilor privind istoria limbii române literare și a limbii noastre în ansamblu este necesară achiziționarea unor programe de ultimă oră, precum și specializarea unor persoane pentru utilizarea lor eficientă. Sperăm de asemenea ca într-un viitor nu prea îndepărtat se va realiza și dorita joncțiune cu Banca de texte din faza modernă și contemporană a limbii române, proiectată a se realiza la Centrul de Studii Românești de pe lângă Universitatea din Anvers, inaugurat în primăvara anului 2000 sub conducerea cunoscutei romaniste și românești Liliane Tasmowski.

5. Resurse bibliografice

Amintim pe scurt și principalele resurse bibliografice privitoare la limba română elaborate de Institut sau de membri ai acestuia. *Bibliografia limbii române*, inițiată de Al. Rosetti și definitivată de Aurel Nicolescu, a rămas nepublicată.⁵⁴ *Bibliografia românească de lingvistică (BRL)* referitoare la lucrările de lingvistică apărute în țară începând din 1944 apare anual în revista „Limba română”; în 1999, ea totalizase deja 64.340 de titluri, în peste 3.300 de pagini de tipar; se preconizează introducerea în calculator a tuturor numerelor din BRL în vederea publicării unui volum cu itemurile ordonate pe autori și pe domenii

⁵³ [] Ioana Vintila-Radulescu (2000). Colaborarea în cadrul ACCT/Agentiei Interguvernamentale a Francofoniei și al Rifal, în „*Terminologia în România și în Republica Moldova*”, Cluj-Napoca, p. 51-52 și La coopération dans le cadre de l'ACCT (Agence intergouvernementale de la Francophonie), în „*Terminometre Hors-série n° 4. La terminologie en Roumanie et en République de Moldova*”, p. 57-58.

⁵⁴ [] I. Coteanu, I. Danaïla (1970). Introducere în lingvistică și filologia românească. Probleme. Bibliografie, București; T. Vianu (red. resp.) et al. (1972). Bibliografia analitică a limbii române literare. 1780–1866, București; Gh. Chivu, Mariana Costinescu (1974). Bibliografia filologică românească. Secolul al XVI-lea, București.

(descrie și separate mai amanunțit decât în forma aparută, cronologic, cu indice de domenii, materii, cuvinte, autori etc.).

Pentru domeniul terminologiei s-au realizat bibliografiile ale dictionarelor terminologice, respectiv ale studiilor de terminologie⁵⁵ și ale standardelor românești de/cu terminologie⁵⁶, precum și un repertoriu bio-bibliografic al terminologilor din România⁵⁷, care va fi inclus în repertoriul internațional al terminologilor din domeniul neolatin pregătit de Uniunea Latina, fiind în curând accesibil pe Internet.

6. Concluzii

Nu ne vom referi aici la alte tipuri de lucrări (gramatici⁵⁸, tratate⁵⁹, enciclopedii⁶⁰ etc.) elaborate de Institut sau de cercetători ai acestuia ori la alte tipuri de resurse care ar merita să fie elaborate de noul institut, pentru a înlocui lucrări mai vechi și a valorifica posibilitățile oferite culturii de societatea informațională, de exemplu un nou dicționar de frecvență al limbii române ș.a.

Deși dictionarele pe CD-Rom și cele pe Internet sunt solicitate de tot mai mulți utilizatori din țară și din străinătate, care cer tot mai des informații cu privire la eventuale dictionare românești on-line, până în prezent a existat la noi o anumită reticență a editurilor proprietare ale drepturilor asupra edițiilor pe suportul tradițional de hârtie față de acest nou mod de difuzare. Nu trebuie însă să existe temerea că folosirea și a noilor suporturi ar diminua vânzarea cartilor, în condițiile în care, în ciuda tuturor eforturilor, un procent încă infim din populația României are acces la PC-uri. De altfel, practica altor țări arată că, în mod neașteptat, difuzarea și în format electronic chiar a sporit desfacerea cartilor, carora le-

⁵⁵ [] *Anca Fezi et al. (2000). Bibliografia lucrărilor de terminologie (1990-1999). România, în „Terminologia în România și în Republica Moldova”, Cluj-Napoca, p. 103-113 și pe discșeta anexată revistei „Terminometro Hors-série n° 4. La terminologie en Roumanie et en République de Moldova”, 2000.*

⁵⁶ [] *Aurora Petan, Edy Savescu (2000). Standarde românești de/cu terminologie (1990-1999). România, în „Terminologia în România și în Republica Moldova”, Cluj-Napoca, 2000, p. 117-126 și pe discșeta anexată revistei „Terminometro Hors-série n° 4. La terminologie en Roumanie et en République de Moldova”, 2000.*

⁵⁷ [] *Adriana Marinescu (2000). Repertoriul bibliografic al terminologilor. România, în „Terminologia în România și în Republica Moldova”, Cluj-Napoca, 2000, p. 128-139 și pe discșeta anexată revistei „Terminometro Hors-série n° 4. La terminologie en Roumanie et en République de Moldova”, 2000.*

⁵⁸ [] **** (1954, 1963). Gramatica limbii române, București, ed. I, 1954; ed. a II-a, revazută și adăugită, 1963; Mioara Avram (1986, 1997, 2001). Gramatica pentru toți, București, 1986; ²1997; ³2001.*

⁵⁹ [] *Al. Rosetti (redactor responsabil) et al. (1965, 1969). Istoria limbii române. București, vol. I. Limba latină, vol. al II-lea; Al. Graur, Mioara Avram (1970-1989). Formarea cuvintelor în limba română, București: I. Fulvia Ciobanu, Finuta Hasan (1970). Compunerea; II. Mioara Avram et al. (1978). Prefixele, 1978; III. Laura Vasiliu (1989). Sufixe, I. Derivarea verbală etc.*

⁶⁰ [] *Marius Sala, Ioana Vintila-Radulescu (1981). Limbile lumii. Mica enciclopedie, București; (1984). Les langues du monde. Petite encyclopédie, București – Paris; Marius Sala (coord.) et al. (1989). Enciclopedia limbilor romanice, București; (2001), Enciclopedia limbii române, București.*

a facut în felul acesta reclama si care prezinta, la rândul lor, alte avantaje în utilizare în raport cu CD-Romurile, cele doua tipuri specializându-se si în functie de necesitati. Astfel, având în vedere culegerea lor computerizata, atât DEX, cât si MDA si DOOM₂ ar putea fi primele dictionare ale Institutului difuzate în viitor si pe CD-Rom.

Credem ca si diverse lucrari valoroase ale Institutului, care, exclusiv din motive financiare, nu-si gasesc editori de ani de zile, nici în tara, nici în strainatate (ca *Bibliografia limbii române*, *Dictionarul spaniolei americane* s.a.), ar putea fi valorificate prin aducerea lor la cunostinta celor interesati pe aceasta cale, tot mai utilizata în societatea informationala actuala. O conditie pentru viitor este realizarea din capul locului a lucrarilor institutului pe calculator, care a devenit posibila prin tot mai buna dotare tehnica a Institutului, realizata prin eforturile directorului sau, precum si prin însusirea, de catre un numar tot mai mare de cercetatori din Institut, în special din generatiile tânara si mijlocie, a cunostintelor de operare pe calculator, inclusiv, în unele cazuri, a lucrului cu baze de date.

Prin realizarea proiectelor de editare pe CD-Rom si pe Internet vom recupera relativa întârziere în acest domeniu fata de difuzarea în România, de catre Grupului Editorial Litera din Republica Moldova si firma Litera International, cu sediul în Bucuresti, a unor CR-Romurile cuprinzând, în diverse combinatii, mai multe titluri⁶¹. Speram ca CD-Romurile consacrate unor dictionare ale Institutului vor fi, desi tot protejate, mai usor de instalat decât cele de la Litera si ca vor oferi mai multe facilitati în utilizare decât acestea, care nu sunt foarte practice, mai ales pentru cercetatori, în ciuda structurii lor modulare si a interfetei lor comune, despre care în reclama se spune ca permit activarea simultana a tuturor dictionarelor.

Pentru progresul cercetarilor si dezvoltarea si prelucrarea resurselor la nivelul exigentelor pe plan mondial, credem ca în viitor se impune o mai buna colaborare, în interes reciproc, între lingvisti si informaticienii preocupati de probleme asemanatoare.

⁶¹ [] Corectorul electronic ORTO 2001 ROM SP, Dictionarul ortografic al limbii române, Gramatica uzuala a limbii române, Noul dictionar explicativ al limbii române, Marele dictionar de neologisme de Florin Marcu, Dictionarul de dublete etimologice ale limbii române de Marcu Gabinschi si un Dictionar de termeni de afaceri englez-român.

Contributia lingvisticii la studiul terminologiilor științifice

Angela BIDU-VRĂNCEANU
Universitatea din București, Edgar Quinet nr. 5-7
vrancean@gpsnet.ro

1. Se admite „laicizarea” științelor [1] sau importanța lor socio-culturală, economică și pedagogică tot mai mare în societățile moderne. Aceasta înseamnă că *limbajele specializate* și *terminologiile* lor nu mai reprezintă coduri total inaccesibile vorbitorilor obișnuiți, nespecializați sau de altă specialitate. În direcția deschiderii, chiar și parțială a codurilor științifice, *dictionarele generale* [2], care includ un număr destul de mare de termeni științifici joacă un rol deosebit pentru a asigura accesul la sensul specializat oricărui vorbitor insuficient informat, pentru a-l ajuta să rezolve ambiguitățile de diferite tipuri și chiar să utilizeze adecvat o terminologie. Permanentă raportare la dictionarele generale ca forme instituționalizate de reglare a uzului nu numai al cuvintelor din limba comună, ci și a termenilor specializați constituie premisa de la care pornim pentru a susține importanța lingvisticii în descrierea terminologiilor științifice, în receptarea și utilizarea lor adecvată chiar și de către nespecialiști.

Pe aceste poziții s-a situat activitatea în cadrul a trei contracte de cercetare științifică pe anii 1997, 1999 și 2000, finanțate de CNCSIS (Consiliul Național de Cercetare Științifică). Au fost studiate limbajul **filozofic**, terminologiile **matematică**, **mineralogică** și din **artele plastice** și, dintr-o perspectivă mai limitată **medicină**, **lingvistică** și **științele politice**. Rezultatele cercetărilor au fost publicate în două volume: *Lexic comun, lexic specializat* [3], care conține studii cu caracter monografic și *Lexic științific interdisciplinar* [4], reprezentând o sinteză a lexicografiei generale și specializate pentru termenii din fiecare dintre domeniile studiate care apar mai mult decât într-o terminologie științifică.

În toate cercetările întreprinse s-a urmărit adoptarea *unei grile metodologice comune* atât pentru clase de cuvinte din limba comună (*abstractele*), cât și pentru termenii specializați din orice domeniu. S-a obținut atât caracterizarea fiecărei terminologii studiate în parte, cât și desprinderea unor trasaturi generale ale terminologiilor științifice, relevante din punct de vedere lingvistic. S-au avut în vedere aspecte *paradigmatice* privind diferitele modalități de definire a sensului, relațiile semantice (*monosemie/polisemie, hiponimie, sinonimie*) din perspectiva necesității ca termenii științifici să fie monoreferențiali, univoci din punct de vedere semantic și să nu aibă sinonime. Analiza **sintagmatică** a gradului de non-determinare contextuală ca o condiție de exprimare a sensului specializat a indivi-

dualizat terminologiile științifice studiate, de la o *libertate contextuală* mai mare (terminologia **matematică, mineralogică**) sau relativă (terminologia **filozofică**) până la o *strictă determinare contextuală* (terminologia **politică** și din **artele plastice**). Acolo unde independența contextuală e mai mare, determinările contextuale exprimă în mod similar în diferite terminologii (**matematică, filozofică, lingvistică**) subcategoriile științifice care dezambiguizează lexicul științific interdisciplinar. Caracterizarea termenilor științifici prin marci diastratice în dicționarele generale și enciclopedice ca tipuri de informații sintagmatice reprezintă un aspect foarte important pentru uzajul adecvat de către specialiști, aspect deficitar, inegal rezolvat.

De pe poziția receptorului nespecializat care decodează sensul total sau parțial, un rol important îl are *definiția lexicografică* care, spre deosebire de cea *terminologică* trebuie să fie mai mult sau mai puțin *naturală* și prin aceasta accesibilă. Existența celor două tipuri de definiții ale termenilor specializați este în general admisă și compararea lor este favorizată de prezentarea sintetică, sinoptică propusă de noi [4]. Chiar și în cazul definițiilor strict terminologice, Em. Vasiliu [5] a susținut și demonstrat prin diferite exemple relevanța diferită a unor componente de sens pentru vorbitorul specialist sau non-specialist. Pornind de la aceste constatări de principiu, ar fi justificat ca termenii științifici să aibă *definiții alternative*, științifice și pre-științifice [6], condiționate atât de o interpretare semantică, cât și de una pragmatică. Din această perspectivă, definițiile termenilor științifici în dicționarele generale ar trebui să difere de cele din dicționarele specializate pentru a facilita deschiderea codurilor științifice și pentru a dezambiguiza lexicul științific interdisciplinar (din principiu, de interes mai larg) sau tangentele cu limba comună. Din păcate, cu mici excepții (**matematică**) selecția termenilor științifici și definirea lor nu diferă aproape deloc în dicționarele generale și în cele specializate.

2. Din perspectiva lingvistică, terminologiile investigate prezintă o serie de particularități:

Matematică se caracterizează prin cel mai mare grad de abstractizare și de ermetism la nivelul sensurilor și definițiilor lor. Compararea definițiilor specializate cu cele din dicționarele generale arată că acestea din urmă definesc diferit și mai accesibil termenii, fără a afecta precizia lor semantică. Sensurile univoce, fără sinonime nu sunt condiționate contextual; sintagmele mai mult sau mai puțin fixe diferentiază subcategoriile conceptuale (de ex. *sistem de ecuații, ~ de curbe, ~ de numeratie, ~ de referință*) și nu afectează independența semantică a acestora. Această terminologie dispune de cea mai bună marcare diastratică în DEX, chiar dacă există numeroase situații în care apartenența la matematică rezultă numai din definiție (manieră de caracterizare practică sistematică și nu întotdeauna convenabil de DEX în cazul altor terminologii). **Matematică** are cel mai bogat lexic științific interdisciplinar, cei mai numeroși termeni comuni fiind cu **fizică, filozofia, logică**, dar și cu **lingvistică, biologie, arhitectură** s.a.; termenii interdisciplinari își pastrează aproape neschimbat sensul, indiferent de domeniul în care se utilizează. Dacă în unele cazuri (relația cu **fizică, logică, filozofia**) punctul de plecare pentru lexicul interdisciplinar nu se poate stabili cu certitudine, în destule alte situații, **matematică** este sursa „împrumutului” făcut de alte științe (**arte plastice, arhitectură, lingvistică** s.a.)

Mineralogia reprezintă și ea un grad mare de ermetism sau închidere a codului, majoritatea termenilor fiind univoci semantic, monoreferențiali și implicit, independenți contextual. Determinările contextuale reprezintă subtipuri, ca și în alte terminologii (**matematica, filozofie** de ex.: *acvamarin brazilian, ~ sintetic, ~ siamez*, etc.) Are un număr mai limitat de termeni comuni cu alte științe (**chimia, artele plastice, simbolistica**) și, cel puțin pentru ultimele două, **mineralogia** este punctul de origine al termenilor interdisciplinari. În ciuda caracterului strict specializat al acestei terminologii, marcarea diastratică din dicționarele generale este deficitară.

Terminologia **filozofică** se caracterizează printr-un grad oarecare de ambiguitate, determinat de variații de interpretare în funcție de curente și tipuri de texte, dar și de contactele cu alte științe sau cu limba comună. De aceea definițiile termenilor **filozofici** nu se pot limita la dicționare, fiind necesară analiza strategiilor argumentative și a figurilor textuale. Invers proporțional cu această necesitate de dezambiguizare, DEX-ul prezintă o marcă diastratică deficitară atât pentru termenii filozofici, cât și pentru celelalte terminologii cu care se stabilesc interdisciplinarități, cum ar fi **matematica, lingvistica** și alte domenii **umaniste**. O bună parte a lexicului științific interdisciplinar are ca punct de plecare **filozofia**, al cărui sens se păstrează ca o medie semantică în majoritatea disciplinelor. Ca și în alte științe, determinarea contextuală exprimă în general subtipuri (de ex. *sistem al științelor, ~ axiomatic, ~ filozofic*).

Terminologia **artelor plastice** prezintă aspecte paradoxale. Maniera de înregistrare și de definire echivocă, imprecisă a acestor termeni în dicționarele generale da impresia unui nespecialist de falsă accesibilitate, interpretare contrazisă categoric de definițiile precise, riguroase din dicționarele și textele specializate. Dependenta contextuală strictă a numeroși termeni din **artele plastice**, al cărui sens specializat e condiționat de sintagmele fixe în care apare (de ex. *acord cromatic, compoziție de gen, semn plastic*) reprezintă o altă caracteristică a acestei terminologii. **Artele plastice** au un lexic științific interdisciplinar bogat, în care se remarcă faptul că sunt preluați cu unele modificări semantice (privind interesul pentru acest domeniu) termeni din alte științe, cum ar fi **chimia, mineralogia, matematica, fizica**. DEX-ul nu utilizează decât marcele diastractice (pictura), (sculptura) dispuse nesistematic și rar, ceea ce contribuie la o tratare deficitară a acestei terminologii.

Lexicul **științelor politice** prezintă, din perspectiva analizei întreprinse de noi, o serie de particularități (unele asemănătoare cu **artele plastice**). Se remarcă dependenta contextuală strictă a acestei terminologii, nici unul dintre termeni nefiind total liber contextual. Sensul specializat în **științele politice** se exprimă, deci, aproape exclusiv pe cale sintagmatică, în contexte mai mult (*celula de criză, agregare de interese*, de ex.) sau mai puțin fixe (diverse combinații cu adjectivul **politic** în sintagme nominale: *capital politic, cartel ~, algoritm ~, contract ~, dialog ~, alternanță politică*). Preia (fără să fie niciodată punct de plecare termeni din numeroase și variate științe: **economia, filozofia, dreptul**, dar și **lingvistica, biologia, medicina, geografia, fizica, psihologia, sportul**). În majoritatea acestor cazuri nu există o motivare de conținut strictă (dincolo de întrebuintarea metaforică), ceea ce determină, în mare parte, mai curând un lexic științific interferent

decât unul interdisciplinar. Poate și din cauza modificărilor continue și rapide din domeniul politicii, DEX-ul înregistrează în mica măsură termeni și sensuri din acest domeniu diastratic, ceea ce constituie un dezavantaj în impunerea acestei terminologii.

3. Analiza lingvistică a limbajelor științifice (care ar putea fi extinsă) permite caracterizarea unor terminologii ca „*puternice*” (**matematica, mineralogia** de ex.), iar a altora mai „*slabe*” în diferite forme și grade (de ex. **științele politice, artele plastice**), cu dificultăți mai mari de deschidere a codurilor în cazul primei categorii.

Delimitarea componentelor de sens relevante diferit în funcție de vorbitori specializați și nespecializați ar putea constitui o bază obiectivă pentru rezolvarea mai eficientă a definițiilor alternative în dicționarele generale, foarte importante în „laicizarea” științelor necesară în grade diferite în epoca actuală. Exprimarea sensului specializat condiționat de dependențele contextuale mai mici (pentru terminologiile „*puternice*”) sau mai mari (pentru terminologiile „*slabe*”) constituie o caracterizare lingvistică relevantă. În schimb, în unele cazuri (ca pentru terminologia **politica**), determinările contextuale sunt mai favorabile, „transparentei” semantice sau deschiderii codurilor specializate.

Analiza lexicului științific interdisciplinar (LSI) poate contribui și ea la determinarea specificului unor terminologii. Științele care constituie sursa, punctul de plecare pentru o parte a LSI își susțin, și pe această cale, statutul de terminologie „*puternică*” (de ex. **matematica, fizica** și, din acest punct de vedere **filozofia**). Dimpotrivă, atunci când punctul de plecare nu se poate stabili aproape niciodată la nivelul unor terminologii (**științele politice, artele plastice**), aceasta constituie o modalitate de determinare specifică. Diferențierea interdisciplinarităților (cu o motivare de conținut determinată de considerarea referentului din diferite puncte de vedere sau de un transfer conceptual) de simplele interferențe (mai puțin sau deloc motivate, cu modificări de sens ale termenilor, multe metaforice) se bazează pe aprecierea distanței semantice, verificată obiectiv.

Dat fiind rolul dicționarelor generale în impunerea și extinderea terminologiilor științifice, de interes pentru diferite categorii de vorbitori, carentele constatate în tratarea sensului și în marcarea lor diastratică riguroasă conduc la concluzia necesității unei reconsiderări și remedieri a manierei de tratare din perspectiva „laicizării” științelor.

Referințe bibliografice

- [1] F. Rastier (1995) Le terme; entre ontologie et linguistique. *Banque des mots* 1995/7, p. 35-65
- [2] DEX - Dicționar explicativ al limbii române, (1996) ed. a 2-a sub coord. acad. I. Coteanu și Dr. Lucretia Mares, Ed. Univers Enciclopedic, București 1996

-
- [3] A. Bidu-Vranceanu – coordonator (2000). *Lexic comun, lexic specializat*, Editura Universitatii din Bucuresti, 2000, cu colaboratorii: Alice Toma (**matematica**), Silvia Savulescu (**mineralogie**), Claudia Ene (**filozofie**), Alexandra Vrânceanu (**arte plastice**)
- [4] A. Bidu-Vranceanu – coordonator (2001). *Lexic stiintific interdisciplinar*, Editura Universitatii din Bucuresti, 2001, cu colaboratorii: Silvia Savulescu (**stiinte politice si mineralogie**), Alice Toma (**matematica**), Claudia Ene (**filozofie**), Alexandra Vrânceanu (**arte plastice**)
- [5] Em. Vasiliu (1980). Sens si definitie lexicografica „Studii si cercetari lingvistice”, an XXXI, 465, 1980
- [6] Em. Vasiliu (1982/1983). Adevar analitic si definitie lexicografica „Analele stiintifice ale Universitatii „Al. I. Cuza” din Iasi”, sectiunea III, tom XXVIII/XXIX, 1982/1983

Gramaticile generative nontransformationale

Emil IONESCU
Universitatea Bucuresti, Facultatea de Litere
Str. Edgar Quinet nr. 5-7,
Email: eionescu@racai.ro

Acest articol este o prezentare generala a gramaticilor generative nontransformationale (GNT) si a prezentei lor în cercetarea lingvistica din România. În prima sectiune a articolului este descrisa geneza acestor gramatici. În sectiunea a doua, sunt prezentate pe scurt caracteristicile lor, în timp ce în partea treia si a patra se mentioneaza principalele realizari stiintifice si formele de existenta institutionala ale curentului. Partea a cincea este consacrata initiativelor si pasilor care au dus la patrunderea acestor gramatici în mediile stiintifice de la noi. Concluziile articolului se vor a fi o pledoarie în sprijinul eforturilor de dezvoltare a acestei directii în cultura stiintifica româneasca.

1. Gramaticile generative nontransformationale: aparitia lor

Gramaticile generative nontransformationale reprezinta, în interiorul lingvisticii formale contemporane, o directie extrem de influenta si de un remarcabil dinamism. Istoria acestei directii este, desigur, mai recenta decât istoria generativismului din care face parte. Este însa o istorie deja bogata si diversa. Printre altele, diversitatea se exprima si prin faptul ca suntem obligati sa vorbim despre *gramatici* si nu despre o gramatica nontransformationala, pur si simplu.

Putem plasa începuturile acestei istorii la cumpana dintre anii '70 si '80. Sunt anii când programul gramaticii universale al lui Noam Chomsky este pe punctul sa depaseasca starea de impas atinsa prin faza denumita de istoricii miscarii "teoria standard". Privita din perspectiva prezentului, lucrarea din 1981 a lui Chomsky ("Lectures on Government and Binding") tocmai acest lucru îl subliniaza: depasirea crizei prin propunerea unui model nou de gramatica universala.

Punctele în care gramatica universala este reformulata în cadrul modelului "Government and Binding" (GB) nu sunt putine si nici neînsemnate. Dar cea mai importanta modificare a fost operata într-una din componentele care nascuse initial cele mai mari sperante: componenta transformarilor. Formulata succint, regândirea conceptului de transformare în cadrul modelului GB înseamna doua lucruri: simplificare si îngradire.

Simplificare, deoarece marea varietate de transformari se reduce acum la o singura operatie: deplasarea unui constituent oarecare α . Si îngradire, pentru ca deplasarea nu se poate produce oricum, ci numai în conditiile în care anumite reguli foarte generale, numite *principii*, sunt respectate.

Nu toti adeptii generativismului au fost însa multumiti cu noua propunere. Ceea ce s-a reprosat a fost ca transformarile ramâneau mai departe mecanisme prea puternice - în ciuda îngradirilor si a simplificarilor – deoarece ele operau pe un domeniu prea larg: cel al structurilor sintactice. O alta obiectie viza temeiurile mentale ale operatiei de deplasare: în ciuda plauzibilitatii aparente a acestei ipoteze, nu exista dovezi - sustineau criticii - ca mintea implicata în utilizarea limbajului ar face uz de o astfel de operatie. În sfârșit, existau cercetatori care considerau ca noul model de gramatica universala era greoi din punct de vedere computational, tocmai din cauza operatiei de deplasare: anume, pentru fiecare deplasare de constituenți, este necesara o verificare a compatibilitatii dintre principii si deplasarea constituentului.

În ansamblu, divergentele legate de conceptul de transformare au pregatit cea mai mare ruptura pe care a cunoscut-o în istoria sa curentul gramaticii universale. Criticii radicali ai conceptului de transformare au propus renuntarea la acest mecanism, propunere pe care Chomsky si cei ce l-au urmat nu au acceptat-o niciodata. Începând cu anul 1981, ruptura se oficializeaza. Apar pe rând Gramatica Lexico-Functionala (LFG - Bresnan si Kaplan), Gramatica Sintagmatica Generalizata (GPSG - Gazdar, Klein Pullum si Sag), Gramatica Arborilor Adaugati (TAG - Joshi), Gramatica Centrilor de Sintagma (HPSG – Pollard si Sag), Gramaticile Categoriale de Unificare (CUG- Uzkoreit)

2. Caracteristicile GNT

Dincolo de varietatea lor, gramaticile nontransformationale au un set de trasaturi comune:

- Exploateaza în mod generalizat reprezentarile în termeni de trasaturi
- Fac recurs la mecanismul unificarii
- Se bazeaza pe constrângeri
- Sunt gramatici lexicaliste
- Au adecvare computationala

2.1. Reprezentari: structurile de trasaturi

Reprezentarile în termeni de trasaturi sunt bine cunoscute în lingvistica moderna, datorita fonologiei si semanticii structurale. GNT au meritul de a fi generalizat aceasta notatie la scara întregii teorii lingvistice. Prin perechea trasatura (atribut)–valoare, orice fel de informatie lingvistica – fonologica, morfologica, sintactica semantica, pragmatica – își gaseste o reprezentare adecvata. Câteva exemple: notatia [P(arte de)V(orbire): nume] spune ca o anumita entitate lingvistica este un nume. Reprezentarea [F(orma)V(erbala):

suficiente. Polona, de pilda, face la verbele de persoana I deosebirea între verbele folosite de un barbat și cele folosite de o femeie. Verbul are asadar gen în polona, dar nu și în româna. Pentru a face aceasta diferență între cele două limbi trebuie să se admită că unificarea informației de gen cu cea de verb se poate face în polona dar nu se poate face și în româna. Numai că de această dată constrângerea privind unificarile nu mai are temei formal. Nu se poate spune că în mod necesar verbul are sau nu gen. Unificarile acestor informații sunt prin urmare “contingente”, sau cu un alt termen, “empirice”, tocmai pentru că ele nu derivă din natura însăși a operației. Gramatica unei limbi se descrie mai ales în termenii unificarilor “contingente”.

2.4 Lexicalism

În teoriile contemporane ale gramaticii, lexicalismul este o opțiune privitoare la modul în care este concepută structura cuvintelor în relația lor cu sintaxa. Există teorii, precum GB, care consideră că procesul de constituire morfologică a cuvintelor are loc în sintaxa. În acest sens, GB este o morfosintaxă deoarece generalizează operația de deplasare la nivelul morfologiei însăși, prin mecanismul numit “deplasare centru-centru” (engl. “Head to Head Movement”). Gramaticile de unificare adoptă o strategie distinctă: ele consideră că procesele de constituire morfologică a cuvintelor sunt independente de sintaxa. În această perspectivă, rezultatul proceselor morfologice furnizează sintaxei inputul necesar: cuvintele gata formate. Modularizarea celor două componente ale gramaticii se dovedește preferabilă mai ales în cazul limbilor cu morfologie bogată.

Un alt aspect al lexicalismului asumat de GNT este ilustrat de modul în care sunt construite explicațiile de gramaticalitate. Explicațiile în GNT se sprijină în măsura posibilului (dar într-o măsură mult mai mare decât în alte teorii) pe proprietățile cuvintelor. În istoria generativismului, pasivul, de pilda, a fost considerat multă vreme o structură explicabilă *sintactic*, adică o construcție rezultată din transformări ale unei alte structuri sintactice. GNT afirmă însă că nu e nevoie să se recurgă la structuri sintactice anumite, deoarece toate elementele de care e nevoie pentru a explica o construcție pasivă pot fi codificate la nivelul cuvintelor⁶². Un tratament asemănător poate fi observat în cazul dependentelor la distanță, sau în cel al construcțiilor de ridicare (engl. “raising”), unde rolul unităților lexicale în determinarea proprietăților acestor construcții este de asemenea semnificativ.

⁶² *Preferința aceasta pentru un compartiment de limbă în defavoarea altui compartiment, atunci când se pune problema mecanismelor care justifică o anumită construcție nu e înțeleasă încă nici azi de unii lingviști. Este vorba de aceia care cred că a avansa o explicație lexicalistă atunci când există deja una sintactică pentru un fenomen oarecare înseamnă doar a propune variațiuni pe aceeași temă. Diferențele sunt în realitate cruciale și privesc mecanismele cognitive angajate în utilizarea limbajului. Este deja cunoscut că procesarea unităților lexicale este mai ușor de efectuat decât unele dintre procesările structurilor sintactice. Acest fapt oferă un criteriu valoros de judecare a plauzibilității unei gramatici privite din unghi cognitiv.*

2.5 Adecvare computationala

În lingvistica, o teorie este considerată adecvată, dacă teoria acoperă domeniul de fapte pentru care este construită ca o explicație. O morfologie a unei limbi, de pildă, este adecvată dacă prin regulile propuse se seamă de construcțiile morfologice corecte ale limbii supuse analizei.

Acest principiu foarte general a fost nuanțat de către Chomsky. Nuanțarea este deja celebră: pornind de la ideea că utilizarea limbajului este o proprietate a minții omenestii, Chomsky a susținut că o teorie trebuie socotită adecvată nu doar pentru că produce explicații ale cazurilor de corectitudine, ci și pentru că mecanismele utilizate sunt dovedite (sau cel puțin presupuse) a fi însușite de către mintea omenească. Quine afirmase că dacă avem două gramatici care cu mijloace diferite explică aceeași realitate lingvistică, nu există criterii suplimentare de alegere a uneia dintre ele. Chomsky a replicat că un astfel de criteriu există totuși, el fiind măsura în care fiecare dintre aceste gramatici se folosește de operații cunoscute ca aparținând minții în procesele ei cognitive.

Criteriul suplimentar formulat de Chomsky în evaluarea teoriilor lingvistice a apropiat comunitatea generativistilor de cea a psihologilor și a impulsionat cercetările de psiholingvistică. S-au obținut rezultate interesante și s-au construit ipoteze neașteptate. De pildă, regulile de constituente sînt socotite astăzi niste operațiuni cu mare probabilitate de a fi folosite de inteligența umană. Recursivitatea este și ea considerată a fi o proprietate de care inteligența umană face uz în utilizarea limbajului.

Criteriul lui Chomsky a condus însă și la cercetări cu rezultate greu de judecat. De pildă, despre realitatea psihologică a *urmelor*, concept cardinal al teoriei GB, s-a argumentat și pro și contra, și este foarte dificil chiar și azi să se poată lua o poziție.

Un lucru este cert totuși în evoluția raporturilor dintre teoria lingvistică și realitatea ei psihologică: comparativ cu faza de început, interesul psihologilor și al psiholingvistilor față de ipotezele venite din comunitatea "chomskyenilor" a scăzut semnificativ. A crescut însă interesul psiholingvistilor pentru ipotezele venite din lumea inteligenței artificiale. Este celebră în acest sens ipoteza de organizare a cunostintelor lexicale a lui Quillian, care a atras atenția în mod special colectivității de psihologi și de psiholingviști. Un al treilea factor intra astfel în joc, rezultatul fiind că unele teorii lingvistice au devenit atente la operațiile și mecanismele utilizate de inteligența artificială. Erau exact teoriile generative netransformationale. Consecința principală a acestei deplasări de interes a fost că teoriile în cauză au devenit accesibile utilizării automate. Cu alte cuvinte - și spre deosebire de gramaticile lui Chomsky - ele pot fi implementate computationally.

Vom numi adecvarea unei teorii la domeniul de fapte pe care îl abordează *adecvare lingvistică*. Măsura în care o teorie lingvistică aparține (sau poate fi presupusă a aparține) minții omenestii definește *adecvarea ei psihologică*. Iar gradul în care ea este livrabilă inteligenței artificiale indică *adecvarea ei computatională*. Direcția actuală a curentului de idei pare să fie următoarea: legăturile și dialogul dintre psihologia cognitivă și inteligența artificială sunt într-o continuă creștere, astfel încât adecvarea computatională a

unei teorii lingvistice are sanse mari sa-i confere si adecvare psihologica. Pe aceasta directie sunt plasate gramaticile generative netransformationale.

3. Realizari

Una dintre cele mai importante realizari ale gramaticilor nontransformationale îl reprezinta numarul mare de aplicatii. O enumerare a limbilor supuse analizelor nu este posibila aici, dar se poate preciza ca aproximativ doua treimi din familiile de limbi (considerate in esantioanele lor reprezentative) au fost analizate din perspectiva netransformationala. Este caracteristic acestor analize faptul ca refuza deosebirea chomskyana centru-periferie („core-periphery”). Ele se concentreaza asupra varietatii de date oferite de corpusuri.

Ceea ce este însa cel mai important sub aspectul realizarilor este faptul ca GNT au reusit sa produca replici viabile la analizele paradigmei dominante, cea chomskyana. O serie de fenomene gramaticale – privite de obicei ca fiind de la sine caracterizabile prin mecanismul deplasarii constituentilor – au primit in cadrul GNT analize alternative. Asa s-a întâmplat cu constructiile pasive, cu fenomenul de ridicare (si mai general cu fenomenele de depedenta limitata), cu constructiile nonlocale (precum topicalizarile, structurile relative si interogative). In aceasta privinta, GNT au continuat traditia fireasca, inaugurata de structuralism, traditie constând in regândirea fenomenelor de limba odata cu fiecare noua scoala lingvistica.

4. Forme institutionale de sustinere

GNT sunt bine reprezentate institutional. Ele si-au facut loc în primul rând în programele curriculare ale unor universitati de prestigiu, precum Universitatea Stanford, Universitatea Statului Ohio (Columbus), Universitatea Tuebingen, Universitatea Saarbruecken, Universitatea Groningen, King's College din Londra Universitatea Edinburgh, Universitatea Paris 7. Extensiile acestor programe curriculare sunt scolile de vara. O prestigioasa scoala de acest fel („European Summer School in Logic Language and Information” – ESSLLI) este organizata anual din 1989, cu rolul de diseminare a evolutiilor si curentelor formate în interiorul gramaticilor netransformationale. Este apoi de semnalat, în aceeasi linie a „didacticii” gramaticilor nontransformationale, nou înfiintata scoala de vara de la Konstaz (Germania).

În planul congreselor stiintifice, HPSG si LFG au de multa vreme propriile lor conferinte anuale. Iar un congres tinut o data la doi ani - cel de gramatici formale - urmareste sa adune sub acelasi acoperis toate scolile aceleiasi familii.

Pâna de curând, gramaticile nontransformationale nu au avut o revista proprie. Lucrarile însa au fost si sunt publicate in reviste de prestigiu, precum „Computational

Linguistics” „Natural Language and Linguistic Theory”, „Journal of Linguistics”, „Language” sau „Langages”. O revista orientata explicit spre aceste gramatici este editata de putina vreme la cunoscuta editura olandeza Kluwer. Este vorba despre revista „Grammars”. De asemenea, pe lânga Centrul de Studii asupra Limbajului si Informatiei de la Universitatea Stanford exista de mai multa vreme o deja celebra editura care publica lucrarile esentiale ale domeniului.

5. Gramaticile nontransformationale în România

Prezenta GNT în Romania poate fi discutata având în vedere doua coordonate: cea a contributiilor stiintifice si cea a programelor curriculare.

Din primul punct de vedere, întâia contributie (dupa cunostinta noastra, cel putin) a venit din partea Adrianei Costachescu ([14]). Adriana Costachescu este autorul unui studiu, din perspectiva GPSG (teorie care a precedat si inspirat HPSG), asupra relatiei dintre coordonarea adversativa si subordonarea concesiva. Studiul a fost elaborat in 1993 si publicat in 1996.

Lucrari de prezentare generala a diferitelor forme de GNT sau, dimpotriva, de prezentare a trunchiului comun – unificarea – au fost publicate în ultimii sase ani de Adrian Atanasiu, Verginica Barbu, Ana-Maria Barbu, Florentina Hristea, Emil Ionescu si Rodica Tatar.

Printre „pionierii” aplicatiilor acestor gramatici la limba româna trebuie mentionati Liviu Ciortuz si cercetatoarea italiana Paola Monachesi. Amândoi au folosit teoria HPSG. Rolul lui Monachesi în stimularea aplicatiilor de acest tip la limba româna trebuie în mod special subliniat. Studiile sale asupra cliticelor pronominale din româna au determinat o „mobilizare” a energiilor câtorva cercetatori români. Este vorba despre Ana-Maria Barbu, Emil Ionescu si Amalia Todirascu.

Ana-Maria Barbu a aplicat HPSG în analiza elementelor gravitând în jurul verbului – adverbul de negatie, semiadverbele, auxiliarele – si a ajuns la concluzia ca acestea sunt mai apropiate de afixe decât de cuvinte. Concluzia analizei se întâlnește cu concluzia exprimata în lucrarea Valeriei Gutu Romalo, „Morfologie structurala a limbii române”, în care formele compuse ale verbelor sunt considerate forme cu afix mobil.

O alta contributie a Anei-Maria Barbu priveste ordinea constituentilor in grupul nominal. Valorificând sugestiile de analiza ale lui Valerio Allegranza, Ana-Maria Barbu a propus o clasificare a constituentilor grupului nominal, care este relevanta pentru problema ordinii acestora. Analiza produce astfel solutii clare si eficiente într-o problema complicata de gramatica a limbii române.

Semnalând unele neajunsuri în analiza GB a fenomenului de anticipare clitica a complementului direct nominal în româna, Verginica Barbu si Emil Ionescu propun o abordare alternativa HPSG. Analiza poate fi extinsa si la alte limbi care prezinta fenomenul în cauza. Analiza sustine ca pronumele neaccentuate nu au un comportament uniform,

proprietatile lor depinzând de faptul dacă participa sau nu la structuri de dublare. Noutatea abordării vine din faptul că fenomenul anticipării obiectului direct este în mod ultim justificat prin proprietatile lexicale ale verbului tranzitiv.

Un fenomen care, în aparenta cel puțin, implică recursul la mecanismul deplasării – este vorba de prezenta pronomelor neaccentuate în acuzativ în contexte în care ele nu sunt subordonate față de vreun element din acel context – este tratat într-un alt studiu asupra cliticelor pronominale românești⁶³ (). Studiul arată că ipoteza deplasării constituenților nu este necesară în analiza fenomenului. Este propusă în alternativă o analiză fără deplasări care captează toate proprietățile fenomenului.

O analiză HPSG este propusă de asemenea pentru fenomenul negației duble și multiple în română (). În sfârșit, Amalia Todirascu abordează într-unul din studiile sale asupra limbii române, o categorie de dependente limitate (asa-numitele *tough-constructions*), din aceeași perspectivă HPSG.

În aceeași linie a contribuțiilor științifice, merita amintită o inițiativă instituțională: acreditarea de către CNCSIS, în anul 2001, a Centrului de Lingvistică Computațională de pe lângă Facultatea de Litere. Centrul este perechea universitară a Centrului de Studii Avansate în Inteligența Artificială. Apariția sa a fost semnalată în buletinul european ELSNEWS. Unul dintre programele de cercetare pe anul 2002 ale centrului are în vedere dezvoltarea aplicațiilor de gramatică netransformațională la limba română.

În planul programelor curriculare, GNT și-au făcut loc mai greu, și au fost întâmpinate uneori nu doar cu neîncredere, ci și cu ostilitate. A existat însă din fericire un sprijin substanțial și constant al factorilor de decizie. Ne referim la decanul Facultății de Litere, acad. prof. Dan Horia Mazilu, la rectorul Universității București, prof. dr. Ioan Mihăilescu, la prorectorul aceleiași instituții, prof. dr. Ioan Pânzariu, și la acad. Dan Ioan Tufis, directorul Centrului de Studii Avansate în Inteligența Artificială al Academiei Române, cărora autorul acestor rânduri le exprimă via și profundă sa grațitudine, pentru susținerea pe care a simțit-o mereu în inițiativele sale. Mulțumim acestui sprijin, au devenit realitate câteva proiecte care pot fi considerate succese:

- În programa cursurilor optionale de limbă pentru anul al IV-lea al Facultății de Litere a fost introdus în 1996 un curs introductiv de GPSG, iar din 1997 până în 2001 s-a ținut un curs introductiv de gramatică de unificare cu referire specială la HPSG.
- Din 1999, se predă la Facultatea de Matematică a Universității din București un curs opțional de un an de prelucrare automată a limbii naturale, în care un loc important îl ocupă gramaticile de unificare.
- Din 1997 până în prezent masteratul de lingvistică teoretică al Facultății de Litere din cadrul aceleiași universități găzduiește un curs de un semestru de teorie HPSG aplicată la limba română.

⁶³ În engleză, fenomenul este cunoscut sub numele de „clitic climbing”, și este ilustrat în română de structuri de tipul *Nu-l pot suferi pe Ion*.

- Din 1999, același masterat ofera un seminar de gramatici cu implementare computațională.
- În anul 2000, un proiect de dezvoltare a componentei de lingvistică computațională în cadrul masteratului de lingvistică teoretică a primit sprijin de finanțare din partea Bancii Mondiale și a Guvernului României, sprijin care a făcut posibilă printre altele organizarea unor cicluri de conferințe pe teme de GNT (în special HPSG) la Facultatea de Litere a Universității București. Au conferențiat Ivan Sag (Universitatea Stanford), Anne Abeille și Daniele Godard (Universitatea Paris 7), Stefan Müller (Universitatea din Jena), Robert Malouf (Universitatea Groningen), Howard Gregory (King's College, Londra), Erhard Hinrichs (Universitatea Tübingen), toți fiind personalități recunoscute ale domeniului. Mulțumită aceluși program, cercetătorii români au putut petrece stagii de specializare la universitățile din Lille și Stanford, sau au putut participa la manifestări reprezentative, cum ar fi colocviul UNESCO asupra spațiilor virtuale și multilingvismului de la Paris (aprilie 2001), colocviul de gramatici bazate pe constrângeri Trondheim (august 2001), sau congresul de prelucrare automată a limbilor naturale de la Tokyo, (noiembrie, 2001). Cea mai importantă realizare legată de acest program, a constat însă în posibilitatea unor mobilități studentești, concretizate în vizitele de studiu ale studenților masteratului de lingvistică teoretică, la universitățile din Darmstadt, Tübingen, Paris 7 și Siena.

6. Concluzii

Deși GNT au pătruns în mediile științifice din România mai târziu decât în alte țări, faptul că ele sunt prezente la noi este un lucru încurajător. Există tentația de a privi aceste eforturi de sincronizare cu mișcarea de idei din domeniul lingvisticii formale drept tentative mimetice și superficiale. Este o greșeală gravă. Diversele comunități de lingviști pot desigur ignora un curent, precum cel prezentat mai sus, dar aceasta este o atitudine, pentru a spune așa, pe proprie răspundere. GNT și teoria lingvistică pe care ele au inspirat-o și-au făcut deja loc în lingvistica zilelor noastre și au devenit una din paradigmele majore. În plus, dubla deschidere a acestor gramatici către psihologia cognitivă, pe de-o parte, și către inteligența artificială, pe de alta parte, recomandă această paradigmă drept cadrul *privilegiat* de dialog interdisciplinar din științele umaniste ale contemporaneității. Din acest triunghi, sunt așteptate să apară noi aplicații – unele au și apărut deja - care vor extinde într-un mod neașteptat conceptul de lingvistică aplicată. Pentru toate aceste motive, tentativele de a păstra un contact viu și de perspectivă cu comunitatea științifică a GNT reprezintă o investiție sigură pe termen lung.

Bibliografie

- [1] Abeillé, A. *Les nouvelles syntaxes. Grammaires d'unification et analyse du français*, Armand Colin, Paris, 1993
- [2] Atanasiu, A. *Curs de lingvistică matematică*, Editura Universității București, 1998
- [3] Barbu, A.M. *Gramatici categoriale. Studiu comparativ cu gramaticile de constituenți*, "Limba Română", XLVI, 4-6, p 239-252, Ed. Academiei, 1997
- [4] Idem, *Complexul verbal*, "Studii și Cercetări Lingvistice", Ed. Academiei, sub tipar.
- [5] Idem, *Romanian Determiners: Order and Classification*, "Revue Roumaine de Linguistique", Ed. Academiei, sub tipar
- [6] Idem, *Funcțiile sintactice în Teoria X-Bară*, "Studii și Cercetări Lingvistice", Ed. Academiei, sub tipar Barbu, A.M. și E. Ionescu *Teorii gramaticale contemporane: Gramatica Centrilor de Sintagmă*, "Limba Română", 1, 1996, 31-55
- [7] Idem, *Accusative Clitic Doubling in Romanian*, Liviu Ciortuz, Paola Monachesi, Hans Uszkoreit (editori) "Informal Proceedings of the GE&GL Workshop: Grammar Engineering and Grammar Learning", Tușnad, România, 1997
- [8] Barbu, V. *Despre gramaticile de unificare*, Analele Universității București, seria limbă și literatură română, 2001, p. 45-52
- [9] Barbu, V. și E. Ionescu *Anticiparea complementului direct în limba română în perspectiva HPSG*, Lucrările colocviului "Perspective moderne asupra limbii române", București, Editura Universității din București, (sub tipar)
- [10] Borsley, R. *Syntactic Theory: A Unified Approach*, Edward Arnold, London, 1991
- [11] Bresnan, J (editor) *The Mental Representation of Grammatical Relations*, MIT, Press, Ca. Mass, 1982
- [12] Ciortuz, L. *An HPSG Kernel for Romanian*, manuscris, 1996
- [13] Ciortuz, L, P. Monachesi, și H. Uszkoreit (editori) *Informal Proceedings of the GE&GL Workshop: Grammar Engineering and Grammar Learning*, Tușnad, România, 1997
- [14] Costăchescu, A. "Coordination" adversative et "subordination" concessive, Iliescu, M. și S. Sora, (editori), Rumänisch: Typologie, Klassifikation, Sprachcharakteristik, München, 1996, p. 121-134
- [15] Gazdar, G, E. Klein, G. Pullum și I. Sag, *Generalized Phrase Structure Grammar*, Cambridge, Harvard University Press, 1985
- [16] Gerlach, B. și J. Grijzenhout (editori) *Clitics in Phonology, Morphology and Syntax*, John Benjamins Publishing Company, Amsterdam / Philadelphia, 2000
- [17] Hristea, F. *Introducere în procesarea limbajului natural cu aplicații în PROLOG*, Editura Universității București, București, 2000

-
- [18] Iliescu, M. și S. Sora, (editori), *Rumänisch: Typologie, Klassifikation, Sprachcharakteristik*, München, 1996, p. 121-134
- [19] Ionescu, E. *A Type of SOV Construction in Romanian*, “Cahiers de Linguistique Théorique et Appliquée”, tomes XXXII-XXXIII, 1995-1996, 19-39
- [20] Idem, *Accusative Weak Pronouns in Romanian*, “Cahiers de Linguistique Théorique et Appliquée”, tomes XXXII-XXXIII, 1995-1996, 19-39
- [21] Idem, *Accusative Clitic Doubling in Romanian*, “Cahiers de Linguistique Théorique et Appliquée” tomes XXXII-XXXIII, 1995-1996, 53-73
- [22] Idem, *Accusative Clitic Climbing in Romanian*, “Cahiers de Linguistique Théorique et Appliquée”, tomes XXXII-XXXIII, 1995-1996, 74-87
- [23] Idem, *A Quantification-based Approach to Negative Concord in Romanian* in Geert-Jan M. Kruijff and Richard T. Oehrle (editori), *Proceedings of Formal Grammar Conference Utrecht, 1999*, p. 25-36
- [24] Idem, *pro-Drop: An HPSG Account without Lexical Rules*, “Bucharest Working Papers in Linguistics”, vol. I, nr.1, 1999, 117-124
- [25] Idem, *On the Status of PE in the Direct Object Construction in Romanian*, *Romanian Journal of Information Science and Technology*, volume 4, numbers 3-4, 2001, p. 293-310
- [26] Joshi, A. Introduction to Tree Adjoining Grammar, *Manaster Ramer, A. (editor) The Mathematics of Language, John Benjamins, Amsterdam, 1987*, p. 87-114
- [27] Kruijff, G-J. M. and R. T. Oehrle (editori), *Proceedings of Formal Grammar Conference, Utrecht, 1999*
- [28] Manaster Ramer, A. (ed.) *The Mathematics of Language*, John Benjamins Publishing Company, Amsterdam, 1987
- [29] Monachesi, P. *Clitic Placement in the Romanian Verbal Complex*, Gerlach and Grijzenhout (2000), p. 255-294.
- [30] Pollard, C. și I. A. Sag, *Information-based Syntax and Semantics*, CSLI, University of Chicago Press 1987
- [31] Idem, *Head-driven Phrase Structure Grammar*, The University of Chicago Press, Chicago, 1994
- [32] Shieber, St. *An Introduction to Unification-based Theories of Grammar*, CSLI, University of Chicago Press, 1986
- [33] Tătar, D. *Inteligență artificială*, Editura Albastră, Cluj, 2001
- [34] Todirașcu, A. *Romanian Tough-Constructions*, Ciortuz, L, P. Monachesi, și H. Uszkoreit (editori) *Informal Proceedings of the GE&GL Workshop: Grammar Engineering and Grammar Learning*, Tușnad, România, 1997
- [35] Wood, M. McGee, *Categorical Grammars*, Routledge London and New York, 1993

Catre o teorie X-bar functionala

Neculai CURTEANU
Institutul de Informatica Teoretica, Academia Româna, Filiala Iasi
curteanu@iit.tuiasi.ro

1. Teorii X-bar mai vechi si mai noi

Scopul prezentei lucrari este dublu: (*a*) de a propune o noua X-bar schema, numita X-bar schema *functionala si recursiva* (pe scurt, FX-bar schema), mai generala si mai adecvata decât cele existente, care sa satisfaca cerintele unei abordari functionale a limbajului natural (LN), în particular, ale strategiei lingvistice SCD (Segmentare-Coeziune-Dependentă) [1], [2], si (*b*) de a pune în evidenta faptul ca teoria FX-bar propusa poate reprezenta o posibila (si necesara) solutie la urmatoarea problema ridicata de Noam Chomsky în teoria Minimalist Program [3]: în doua capitole diferite, Chomsky afirma (în doua abordari diferite, aparent contradictorii, asupra structurii sintactice a LN) atât importanta crescânda a teoriei X-bar cât si posibilitatea ca teoria X-bar standard sa fie “*largely eliminated in favor of bare essentials*” (vezi sectiunea 5).

1.1. Teoria X-bar clasica

Printre (sub)teoriile care reprezinta substanta majora pentru câteva teorii formale importante asupra sintaxei (LN), un rol fundamental este jucat de catre asa-numita teorie X-bar. X-bar schemele propuse sunt de obicei însoțite de definitii, ipoteze, restrictii, principii si alte (sub)teorii gramaticale care specifica într-o cât mai mare masura modul concret în care X-bar schemele sunt utilizate pentru a construi structurile sintactice de baza ale LN. În general, teoria X-bar stabileste categoriile gramaticale principale, proiectiile lor lingvistice (minimale si maximale), relatiile de dominare dintre categorii în cadrul acestor proiectii, sub-, co-, sau supra-ordonarea lor. Toate aceste aspecte asigura numai coloana vertebrala (infrastructura) consistenta a structurii sintactice în reprezentarea LN. Un capitol de o importanta deosebita este relatia dintre teoria X-bar si alte sub(teorii) sintactice si semantice care formeaza întregul corpus al unei anumite teorii lingvistice.

Prima forma a X-bar teoriei este propusa de catre Noam Chomsky în lucrarea *Remarks on Nominalizations* (1970) [4]. Chomsky scoate în evidenta diferentele reale existente în urmatoarele sintagme nominale:

- (1.1) *John's criticism of the book;*
- (1.2) *John's criticizing the book;*

în special datorita sablonului verbal (similar cu al verbului "criticize") rezultat din gerunziul nominal (pentru engleza) "criticizing", în comparatie cu forma nominala derivata "criticism".

Teoria X-bar originala propusa de Chomsky identifica trei categorii lexicale *primitive*, N [Eng: *noun*], V [Eng: *verb*] si A [Eng: *adjective*], fiecare dintre ele cu câte doua categorii sintagmatice corespunzatoare. Mai exact, utilizând notatia $X = N, V, A$, categoria gramaticala X se întâlnește ca *nucleu* [Eng: *head*] într-o categorie intermediara X' (sau X1, sau X¹), traditional numita X-bar, precum si într-o categorie maximala X" (sau X2, sau X²), traditional numita XP, reprezentând *proiectia maximala* a categoriei gramaticale X (lexicala sau nelexicala). Categoria X este numita *nucleul sintagmelor X'* (sau X1) si X" (sau X2) care o contin. Sa mai notam ca prescurtarea pentru categoria *prepozitionala* este P.

Ulterior au fost considerate *patru* categorii lexicale, bazate pe urmatoarele combinatii ale celor doua trasaturi N si V (considerate ca fiind generice pentru categoriile lexicale):

N	este o categorie X cu trasaturile [+N, -V];
V	este o categorie X cu trasaturile [-N, +V];
A	este o categorie X cu trasaturile [+N, +V];
P	este o categorie X cu trasaturile [-N, -V].

Teoria X-bar poate fi înțeleasa si ca o specificare a modalitatii în care unele categorii gramaticale sunt *dominate* de catre altele, deci ca o teoriei a dominantei gramaticale (sau, asa cum spune Chomsky, a "*guvernarii*"), care arata cum un nucleu (sau o categorie lingvistica) X se proiecteaza (se extinde) catre categoriile mai complexe (structurile sintagmatice) X' (sau X1) si X" (sau X2, sau XP). Structurile sintactice X1 su X2 devin categorii esentiale ale organizarii si reprezentarii textului în LN.

Deci, X-bar teoria clasica considera ca X, împreuna cu o secventa de *complemente* (sau *argumente*, notate Arg_i) este imediat dominata de X1, în timp ce X1 împreuna cu o secventa de *specificatori* (notata Spec_j) este imediat dominata de catre X2 (sau XP). Utilizând binecunoscutele notatii din domeniul teoriilor lingvistice formale, (X' = X1, X" = X2 = XP), categoriile lexicale si gramaticale ale teoriei X-bar clasice a lui Chomsky sunt urmatoarele:

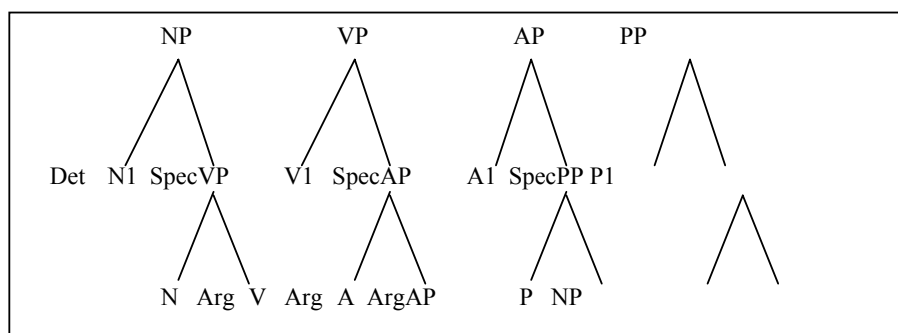


Figura 1.1. Proiectiile categoriilor lexicale din teoria X-bar clasica

1.2. Extinderea teoriei X-bar la categorii non-lexicale

Stowell [5] propune ca teoria X-bar clasica sa fie extinsa la categorii nelexicale sau *functionale*. În particular, categoria gramaticala S [Eng: *sentence*; Rom: *fraza*], care corespunde uneia sau mai multor propozitii gramaticale (clauze), este vazuta ca I2 sau IP, deci ca proiectia maximala a categoriei nelexicale "I", sau INFL [Eng: *Inflectional*]. Nucleul nelexical I (INFL) reprezinta multimea de trasaturi de flexionare atribuite nucleului lexical al clauzei-matrice (propozitia principala, sau chiar una regenta) dintr-o fraza, asa cum sunt timpul, aspectul etc. în clauza a unei fraze. Remarcam *categoria S*, care introduce un anumit grad de ambiguitate în analiza gramaticala, atât în engleza cât si în româna. Termenul adecvat pentru realitatea lingvistica codificata de categoria S ar trebui sa fie acela de "*clauza gramaticala*" pentru engleza [Eng: (*grammatical*) *clause*], si de "*propozitie gramaticala*" pentru limba româna, cu doua sorturi principale: *clauza finita*, prescurtata CLF sau mai simplu CL, si *clauza infinita*, prescurtata CLI.

Astfel în extensia nelexicala a teoriei X-bar, S este proiectia (lingvistica) maximala a categoriei virtuale (nelexicale) I, în timp ce S1 este vazuta ca fiind C2, sau CP, unde nucleul C este un *complementizator*, o categorie gramaticala ce corespunde unei expresii (unui delimitator) sau unei sintagme care introduce o clauza subordonata, e.g. pronume relativ, conjunctie, locutiune conjunctionala etc. Teoria X-bar extinsa acrediteaza urmatoarele structuri:

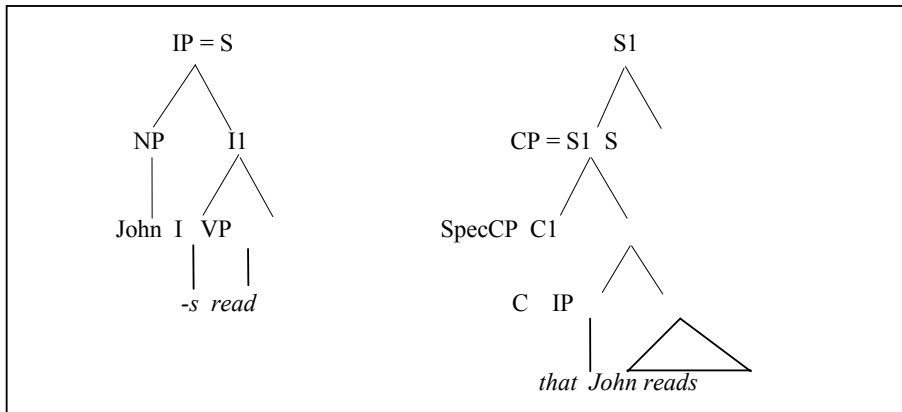


Figura 1.2. Teoria X-bar extinsa la categorii nelexicale

Sunt necesare câteva remarci:

(a) Teoria X-bar extinsa utilizeaza terminologia de "categorii nelexicale (sau functionale)", prin care Stowell, Chomsky si alti lingvisti definesc *noile nuclee* ale structurilor sintactice considerate. Categoria virtuala "I" este, desigur, una nelexicala, si sustine o anumita functionalitate depinzând de categoria lexicala careia îi este atribuita. Categoria C nu este, de obicei, nelexicala (exceptând situatia, posibila, când ea lipseste) deoarece C corespunde unor categorii gramaticale lexical nevide. În ceea ce priveste functionalitatea lui C, suntem de acord ca C corespunde într-adevar unor functii si relatii sintactice si semantice importante pe care le numim *marcheri de propozitie* (subordonate) [1], [2], [6], uneori incluse în clase mai largi cum sunt cea a *marcherilor de discurs* [7], reprezentând în acelasi timp si un element (deci o relatie) de co-referinta în cadrul fenomenului de *legare*, si/sau o "bariera" [8] în cadrul *teoriei limitarii* [9]. Aceste aspecte multi-functionale ale categoriei C nu sunt contradictorii ci doar complementare, întregind un tablou complex al functionalitatii lexical-semantice pentru o categorie lingvistica atât de speciala cum este C.

(b) A doua observatie este dedicata rolului unor categorii nelexicale în cadrul X-bar schemelor extinse. Din Fig. 2. reiese ca subiectul NP are rolul (nesigur) al unui specificator pentru S = IP, în timp ce VP reprezinta complementul categoriei virtuale I. De asemenea, S1 = CP se considera a fi proiectia maximala a categoriei C, în timp ce complementul sintagmei CP este IP. Admitând ca în engleza, din punct de vedere sintactic, aceasta supozitie are sens deoarece categoria C reprezinta nucleul acestor sintagme, în alte limbaje, inclusiv româna, acest lucru este nedecis, în special din perspective semantice si functionale. Unele abordari functionale ale acestor probleme sunt discutate în mai multe lucrari, dar ne vom restrânge sa mentionam aici la a mentiona solutiile oferite de catre *teoria gramaticii functionale* [10] si *strategia lingvistica* SCD [1], [2], [6]. Un interes

special prezinta abordarea *lexicala* (inclusiv *functionala*) a teoriei X-bar ca subteorie de baza în cadrul teoriei sintactice HPSG [Eng: *Head-driven Phrase Structure Grammar*] [11]. O analiza comparativa cu *FX-bar schema* propusa în aceasta lucrare va fi facuta într-o lucrare viitoare.

1.3. X-bar schemele din teoria GB

X-bar schemele propuse de teoria *Government and Binding* (GB) a lui Chomsky [5] sunt urmatoarele:

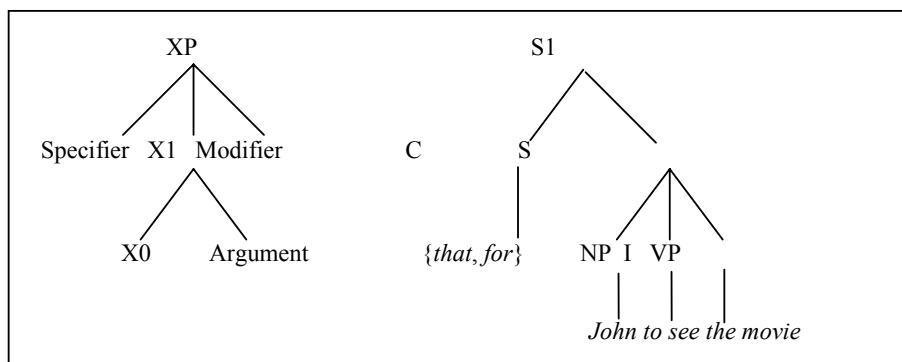


Figura 1.3. X-bar schema generala din GB, X = N, V, A, P, S

În teoria GB exista urmatoarele X-bar *echivalente* pentru proiectiile categoriilor gramaticale (lexicale si nelexicale).

Tabelul 1.3

Proiectii ale categoriilor lingvistice în GB

X	X1	X2
N	N1	NP
V	V1	VP
A	A1	AP
P	P1	PP
I	S	S1

În lucrarile GB [5] si cele care urmeaza, Chomsky considera categoria I ca fiind nucleul lui S, iar complementizatorul C ca fiind nucleul lui S1. În subsectiunea urmatoare, teoria sintactica GPSG a lui G. Gazdar [12] face un important pas înainte catre lexicalitate si catre utilizarea explicita a trasaturilor lingvistice atribuite categoriilor gramaticale.

1.4. Teoria X-bar în GPSG

În teoria lingvistică GPSG [Eng: *Generalized Phrase Structure Grammar*] [12], [13] etc., (sub)teoria X-bar joacă de asemenea un rol central, o sintagmă a LN fiind definită ca proiecția *trasaturilor lingvistice* atribuite *nucleului* [Eng: *head*] acelei sintagme. Informația cuprinsă în trasaturile nucleului determină caracteristicile principale ale comportamentului sintactic al sintagmelor LN. Reamintim că o categorie sintactică în GPSG se reprezintă ca o mulțime de perechi $\langle \text{trasatura}, \text{valoare} \rangle$. De exemplu, eticheta NP [Eng: *noun phrase*] (sau N2), prin care se notează o sintagmă nominală, reprezintă o abreviere pentru mulțimea $\{ \langle N, + \rangle, \langle V, - \rangle, \langle \text{BAR}, 2 \rangle \}$, unde BAR este *numele trasaturii* ce codifică *nivelul de proiecție* a categoriei sintactice $N = \{ \langle N, + \rangle, \langle V, - \rangle \}$. Trasatura BAR poate lua valorile 0, 1, 2. Teoria GPSG consideră N, V, A și P ca fiind *categorii sintactice majore*. Toate celelalte sunt considerate de GPSG ca fiind *categorii minore*: determinatori, complementizatori, marcheri, cuantificatori, alte particule etc. Categoriile majore sunt considerate de către teoria GPSG ca având întotdeauna o *valoare* pentru *trasatura* BAR. Valoarea BAR pentru categoriile minore nu este definită niciodată în GPSG.

Teoria sintactică a GPSG aduce câteva elemente noi și interesante comparativ cu teoria GB: **(a)** X-bar schemele au, ca și în GB, trei nivele de proiecție (valorile trasaturii BAR); **(b)** Pentru economia reprezentării, GPSG propune ca în X-bar schemele de bază, nivelul proiecției lingvistice să fie conservat când se trece de la nucleu către expresiile subcategorizate, mai puțin în cazul în care acest lucru se face prin (alte) reguli explicite; **(c)** Printr-un mecanism de *mostenire implicită*, nivelele BAR de proiecție a nodului-radacină și ale nodurilor-fiice rămân aceleași, mai puțin în cazul în care există o indicație contrară expresă.

O altă caracteristică este aceea că în GPSG nu se întâlnesc categorii abstracte, non-lexicale, cum ar fi "I" (INFL) din GB. Acest lucru este posibil deoarece în GPSG, pentru aceste categorii nelexicale, nu există un nivel de proiecție pe care ele să fie reprezentate (sub nivelul lexical BAR = 0). Consecința este aceea că, în GPSG, S este proiecția unei categorii V. Mai exact, proiecțiile maxime ale lui V sunt VP, S, și S1, depinzând de următoarele valori luate de către trasaturile SUBJ și COMP (= complementizator = C):

$$V[\text{BAR } 2][\text{SUBJ } -][\text{COMP NIL}] = \text{VP};$$

$$V[\text{BAR } 2][\text{SUBJ } +][\text{COMP NIL}] = \text{S};$$

$$V[\text{BAR } 2][\text{SUBJ } +][\text{COMP } \alpha] = \text{S1}; \text{ unde } \alpha \in \{ \text{that, for, whether, if} \}.$$

În sfârșit, trebuie să remarcăm că GPSG trebuie să rezolve problemele întâlnite în mod obișnuit în formalismele gramaticale bazate pe unificarea lingvistică (și/sau logică), de exemplu PATR-II [14], HPSG [15], [16] etc. O astfel de problemă este, în particular, transmiterea informației despre *timbul verbului* între forma flexionară codificată de verb și nodul S. Pentru teoriile lingvistice care permit inserarea în arborele de derivare a cuvintelor flexionare, așa cum este cazul cu GPSG, HPSG etc., informația despre forma flexionară trebuie să poată fi mutată în ambele direcții pe nivelele X-bar schemei. Din aceasta deriva,

în GPSG, condiția ca V să fie *nucleul structurii clazale* care corespunde categoriei S. Pe de altă parte, în GB, informația asupra timpului unui verb poate fi transmisă dinspre nodul I către proiecția sa în S înainte ca I să fie combinat cu forma flexionată a verbului din S. Aceasta situație poate produce potențiale dificultăți procedurale și de reprezentare.

Este important de menționat că proiecțiile categoriilor din Tabelul 4 rămân aceleași pentru GPSG și LFG [Eng: *Lexical Functional Grammar*] (vezi de exemplu [13]), cu diferența notabilă că prima celulă din ultima linie a Tabelului 4 este goală, deoarece în aceste două teorii lingvistice (ca și în altele), categoria virtuală I lipsește.

1.5. O formulare recursivă a X-bar schemelor din teoria Tbarr

Vom propune în această subsecțiune o *formulare recursivă* a teoriei X-bar avându-și originea în *teoria barierelor* (Tbarr) [8], [17] și fiind compatibilă cu teoria sintactică a *Programului Minimalist* (MinP) [4] și cu modelul sau gramatical din *Principii și Parametri* (P&P) [4]. În conformitate cu MinP și P&P, gramaticile concrete ale limbajelor naturale (LNs) reale pot fi modelate de mulțimi de parametri și valorile lor, care specifică principiile și teorii lingvistice universale valabile. Pentru o asemenea setare (asignare) a valorilor parametrilor, relațiile de precedentă (de ordonare liniară) dintre categoriile gramaticale sunt obținute din proprietăți ca marcarea cazuală, atribuirea de roluri tematice (θ -roluri și θ -marcheri), împreună cu alte relații și marcheri ce se aplică la nivelul sintagmelor, clauzelor, și discursului. Din acest motiv relațiile de precedentă pentru X-bar schemele propuse pot fi utilizate independent pe arborii sintactici considerați, informația de ordonare (liniară) a categoriilor fiind dată de următorii *parametri de precedentă*.

(OrdPar) *Un anumit parametru (depinzând de limbaj) specifică dacă secvența de specificatori precede sau succede nucleul, iar un alt parametru (depinzând de limbaj) precizează când secvența complementelor precede sau succede nucleul din X-bar schema.*

De exemplu, în *engleza*, specificatorii preced de obicei nucleele lor nominale, în timp ce în *româna*, în mod normal, ei succed nucleele lor. În general, complementele (argumentele) succed nucleele lor și în *engleza* și în *româna*. Un caz special al argumentului este *subiectul* (sintactic). Aceasta exprimare a (OrdPar) poate fi încă particularizată în funcție de categoriile lexicale concrete, din LNs concrete. De exemplu, atât în *româna* cât și în *engleza*, când o sintagma adjectivală (adverbială) este *predicatională activă*, fiind urmată de anumite argumente (complemente sau adjuncti), atunci este obligatoriu ca ea să succedă propriul nucleu și nu să îl precedă.

Consecința principală a parametrizării dependente de limbaj a precedentei categoriilor lingvistice este că în exprimarea teoriilor lingvistice se pot utiliza arbori neordonati, iar principiile propuse de teoria X-bar primesc un puternic caracter de independență relativ la regulile structurilor sintagmatice. Este important faptul că X-bar schemele obținute în cadrul teoriei X-bar considerate să asigure proiecții adecvate ale categoriilor lexicale, permițând inserarea *adjunctilor*, obținerea categoriilor de proiecție

maximala, si acceptarea faptului ca unele proiectii minimale sau maximale din *structura de adâncime* pot fi *vide* (deci noduri care sa domine *categoriile vide*), conform [9], [8], [17].

Fiind stabilit *principiul* (OrdPar), teoriile GB si Tbarr considera urmatoarele trei nivele ale proiectiei din teoria X-bar, sintetizate de urmatoarele *reguli (principii)* si de *X-bar schemele* corespunzatoare:

(PX0) Fiecare nod X0 dintr-o schema X-bar este fie *vid*, neavând nici o trasatura, fie este nodul-mama al unui element lexical a carei categorie gramaticala si trasaturi sunt specificate la nivelul lexiconului.

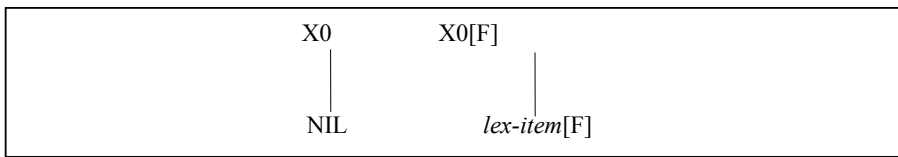


Figura 1.5.1. Nodul X0 în TBarr

(PX1) Fiecare nod X1 (X' sau X^1) având trasaturile lexicale F este fie nodul-radacina al exact unui nod X (care este *nucleu*) cu trasaturile F si al unei secvente de noduri XP (care sunt *complemente*, sau *argumente*), fie este radacina unui nod identic X1 împreuna cu exact un nod XP (care este *adjunct*).

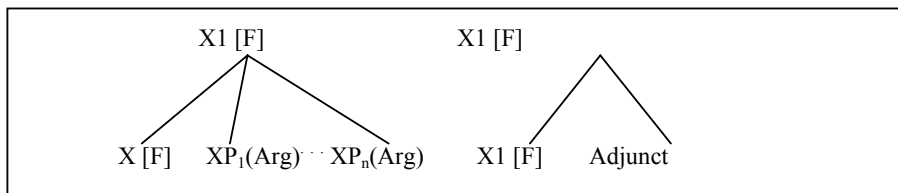


Figura 1.5.2. Nodul X1 în TBarr

(PX2) Fiecare nod XP care are trasaturile lexicale F trebuie sa satisfaca una si numai una din urmatoarele conditii: **(i)** XP este un nod-frunza (nu mai are nici un nod-fiica) si multimea F este *vida*; **(ii)** XP este radacina unei secvente de XPs (*specificatori*) si a exact unui nod X1 mostenind trasaturile F; **(iii)** XP este radacina unei secvente de XPs (*complemente*, sau *argumente*) si a exact unui nod X cu trasaturile F; **(iv)** XP este radacina unui alt nod XP mostenind trasaturile F si a exact unui nod XP.

O *observatie importanta* este aceea ca unele dintre *secventele* XP specificate în regulile (PX1) si (PX2) pot fi *vide*.

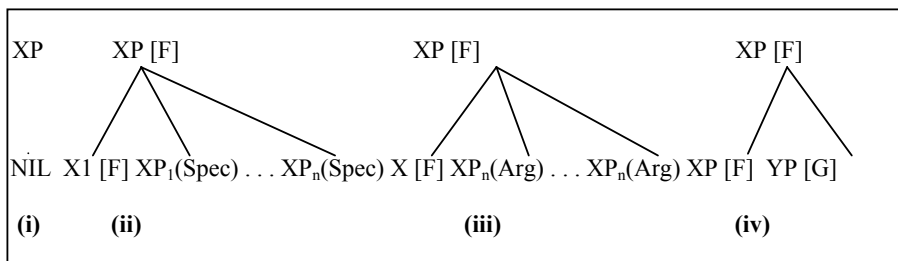


Figura 1.5.3. Nodul X2 în teoria TBarr

Combinând recursiv X-bar schemele rezultate din *regulile* (XP0)-(XP1)-(XP2) se pot obtine toate structurile sintactice întâlnite în *X-bar teoria clasica si extinsa*:

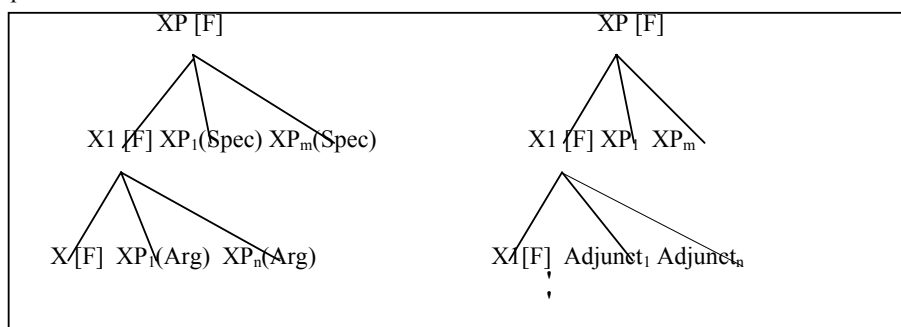


Figura 1.5.4. Formele generale (si recursive) ale X-bar schemelor din TBarr

2. X-bar teoria din modelul P&P al teoriei MinP

2.1 Sistemul Chomskyan al gramaticii universale

Aceasta subsecțiune conturează câteva aspecte implicate de către teoria X-bar în cadrul teoriilor MinP (*Minimalist Program*) și P&P (*Principles and Parameters*) [3]. Pentru a înțelege contextul, este necesar să schițăm teoria lui Chomsky a gramaticii universale UG [Eng: *Universal Grammar*] și a relațiilor sale cu abordarea MinP bazată pe P&P [3]. Sunt introduse următoarele concepte de UG.

Capacitatea utilizării și înțelegerii LN se bazează în esență pe proceduri care pot genera obiecte numite *descrieri structurale* (SDs). SDs sunt *expresii* de limbaj. Teoria unui LN particular constituie gramatica acestuia, în timp ce teoria tuturor limbajelor și a expresiilor pe care le generează ele reprezintă *Gramatica Universală* (UG).

Se consideră că UG specifică anumite *nivele lingvistice*, sau sisteme de reprezentare a informației lingvistice. UG a lui Chomsky [3] presupune că fiecare SD este o secvență (δ , σ , π , λ) de patru reprezentări pe următoarele nivele, respectiv: *structura de adâncime* (D-structură), *structura de suprafață* (S-structură), *forma fonetică* (PF) și *forma logică* (LF). O ipoteză constructivă pentru UG este aceea că limbajul este scufundat în *sisteme de*

performanta care permit ca exprimari în LN sa fie folosite pentru articulare, interpretare, referire, interogare, reflectie si alte actiuni, în timp ce SDs devin un complex de instructiuni pentru aceste sisteme de performanta.

O alta ipoteza standard pentru constructia UG este aceea ca un LN este format din doua componente: un *lexicon* si un *sistem computational*. Aceasta constructie este o inovatie esentiala comparativ cu teoria GB, care pretinde independenta sa fata de orice aspecte computationale sau de implementare. Lexiconul specifica elementele de intrare pentru sistemul computational, în timp ce acesta foloseste intrarile de lexicon pentru a genera derivari si SDs. Derivarea unei exprimari lingvistice particulare implica alegerea elementelor din lexicon si evaluarea construind perechea pe doua nivele de performanta, numite si *reprezentari de interfata*. Una din ipotezele de baza ale teoriei lui Chomsky *Minimalist Program* este aceea ca în constructia SD, utilizând lexiconul si sistemul de evaluare, sunt luate în considerare *numai doua* nivele de interfata, corespunzând lui PF (forma fonetica) si lui LF (forma logica), împreuna cu multimile de perechi (π, λ) rezultate din cele doua forme.

În abordarea P&P a teoriei lingvistice MinP, UG asigura un *sistem de principii* fixat, asociat cu un tablou finit de *parametri evaluati* (pe un numar finit de valori). Regulile pentru un LN particular se reduc la alegerea valorilor pentru acesti parametri. Notiunea de constructie gramaticala este eliminata, împreuna cu regulile particulare de constructie, specifice gramaticilor generative. Constructii ca VP, clauza relativa, pasivul etc. devin doar elemente ale unei taxonomii generale, sau colectii de fenomene explicate prin interactiunea principiilor de UG, legate (setate) cu anumite valori fixate ale parametrilor.

În sistemul computational al UG exista un set de *principii invariante*, fiecare cu un domeniu de *optiuni* restrânse la elementele functionale si proprietatile generale ale lexiconului. O *selectie* Σ printre aceste optiuni determina LN concret. În schimb, un limbaj determina o multime infinita de SDs lingvistice, fiecare pereche (π, λ) fiind obtinuta din nivelele de interfata (PF, LF), respectiv. *Achizitia de limbaj* implica fixarea multimii Σ , în timp ce *gramatica* limbajului se reduce la specificarea lui Σ . În fine, un *sistem de parsare* care este invariant si neantrenat (cum adesea se presupune) poate fi vazut ca o transformare a perechii (Σ, π) într-o schema structurata similara cu o SD. Conditiiile asupra reprezentarilor LN impuse pentru diferite principii si (sub)teorii, cum ar fi teoria *legarii*, teoria *cazurilor*, θ -teoria etc., sunt satisfacuate pe nivelele de interfata ale sistemelor de performanta. Toate aceste ipoteze fac parte din teoria MinP a lui Chomsky si din constructia sa pentru UG.

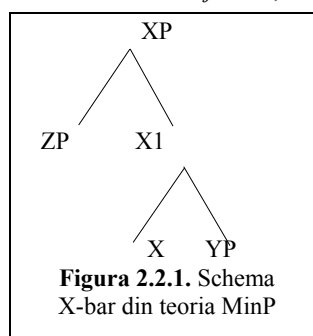
2.2 (Sub)teoria X-bar în contextul teoriei MinP

Sistemul computational al unui LN concret preia reprezentarile unei forme date si le modifica, în timp ce UG trebuie sa furnizeze mijloacele de a reprezenta o multime de elemente din *lexicon* într-o forma care sa poata fi accesata si procesata de catre sistemul computational. Forma sub care este accesat lexiconul de catre sistemul computational poate fi considerata ca fiind o anumita *versiune* a teoriei X-bar. Schemele X-bar pot fi asociate în mod natural cu *structuri de trasaturi lingvistice* [18], ca un *tip de date lingvistice* standard si invariant pentru a reprezenta si a procesa LN eficient. În strategia SCD, *schemele X-bar*

augmentate [19] considerate până acum nu sunt doar tipuri de reprezentare a datelor la nivelul lexiconului ci ele pot asigura structurile invariante fundamentale pentru a reprezenta și a procesa textul în LN la nivel sintactic [1], [2], [6].

În teoria *Minimalist Program* și modelarea P&P a UG, proprietățile și relațiile esențiale sunt formulate în termenii simpli și elementari ai *teoriei X-bar*. Astfel, o *structura X-bar* este compusă din *proiecțiile* lingvistice ale *nucleelor* selectate din lexicon. În *schema X-bar* a teoriei MinP reprezentată în Fig. 2.2.1. sunt prezente două relații locale: relația *Specificator-Nucleu* de la ZP la X, și relația *Nucleu-Complement* de la X și YP (ordinea categoriilor nu este esențială, fiind stabilită de către parametri P&P adecvați de ordonare). Relația Nucleu-Complement (Nucleu-Argument) nu este numai "locală" ci și fundamentală deoarece este asociată (θ -)relațiilor tematice.

Dacă, pentru moment, nu este luată în considerare *relația de adjunctie*, sau adjunctii se considera a se afla printre argumentele-complemente, X-bar structurile pot fi reduse la X-bar schema din Fig. 2.2.1, cu următoarele specificări: **(a)** Sunt considerate numai relațiile locale (deci nici o relație de proiectie între X și vre-o sintagma inclusă în proiectiile maximale YP sau ZP); **(b)** Relația *Nucleu-Complement* reprezintă *relația locală de nucleu* [Eng: *core relation*]; **(c)** O relație locală *admisibilă* a schemei X-bar din MinP este cea *Nucleu-Nucleu*. De exemplu, relația unui verb predicativ cu nucleul predicțional (deverbal) al unei sintagme nominale pe care o subcategorizează; **(d)** O alta relație în X-bar schema din MinP este *legătura de lant* [Eng: *chain link*], corespunzând unui *lant de dominare* sau de *guvernare*.



Guvernarea realizată de nucleu joacă un rol central în toate componentele teoriei MinP asupra UG. Una dintre problemele-cheie este asignarea corectă a trasaturilor nucleului. În HPSG și SCD, de exemplu, acest lucru este realizat la nivel de lexical (BAR = 0), după aplicarea flexionării, cât și la nivel de lexicon (nivel de proiectie notat convențional cu BAR = -1) pentru clasa categoriilor lingvistice cu *proprietăți funcționale (predicative, relationale)*, fie ele verbe, substantive, adjective, marcheri de sintagma, marcheri de discurs etc. care antrenează un comportament sintactic funcțional [2], [6]. În particular, pentru teoria MinP, subteorii ca θ -guvernarea și guvernarea de caz, corespunzând θ -marcării și Cazmarcării, sunt cele mai importante forme de dominare. Un studiu comparativ al guvernării categoriilor (dependentă, dominare), relație prezenta firesc în cele mai importante teorii sintactice formale existente în acest moment, este inclus în [20].

Structurile propuse de teoria X-bar trebuie "animate" de către (sub)teoriile (de asemenea complementare) continute în MinP și P&P, și care explicitează fenomenele de *guvernare*, *legare*, *limitare* etc. ce s-au dovedit a fi importante pentru orice teorie

lingvistica deoarece ele asigură reguli pentru organizarea *lexiconului* și a *sistemului computational* care generează și recunoaște SDs.

De exemplu, în funcționarea *teoriei cazurilor* în contextul schemelor X-bar din MinP, ipoteza standard din MinP este aceea că, într-o frază (propoziție), relația *Specificator-Nucleu* atrage după sine *cazul structural* pentru *poziția de subiect*, în timp ce *poziția de obiect* primește cazul sub guvernarea nucleului V, incluzând construcții în care obiectul marcat cazual de către un verb nu este complementul sau ci doar un adjunct (as-numita *marcare de caz excepțională*).

În continuare este prezentată structura X-bar de bază a *clauzei* în teoria MinP, cu următoarele notații uzuale: C = COMP = Complementizator, T = Timpul, Agr_S = acordul subiectului; Agr_O = acordul obiectului etc.

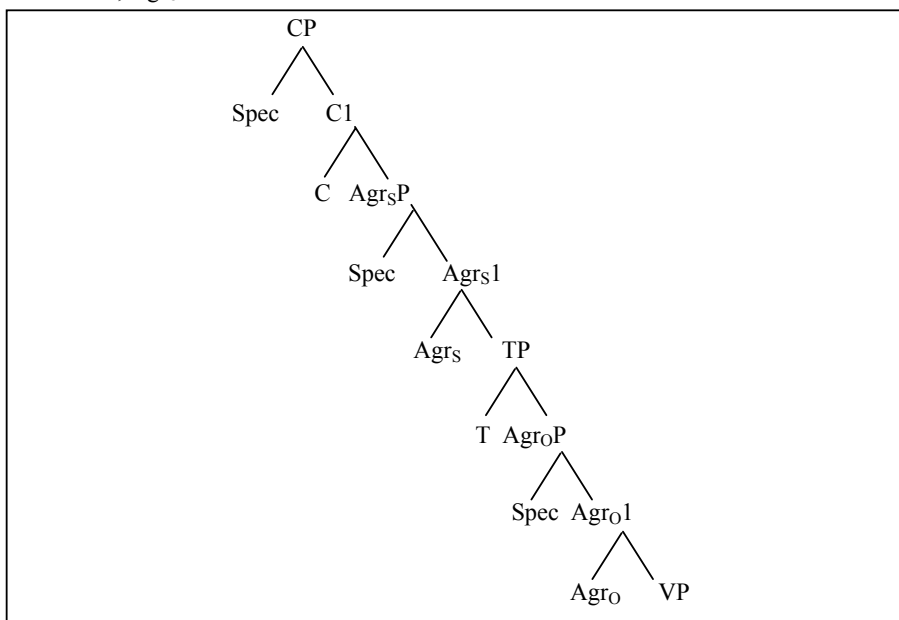


Figura 2.2.2. X-bar structura clauzei în teoria MinP

Schemele X-bar clauzale clasice din Fig. 1.2. și Fig. 1.3. sunt expandate în Fig. 2.2.2., cu următoarea posibilă interpretare funcțională: X-bar schema MinP are ca nucleu VP, care își selectează sintagma-Obiect (sau argument, mai general) prin acord și marcare, afectată apoi de Specificator. Un timp finit T aplicat sintagmei Verb-Obiect generează sintagma TP [Eng: *tensed phrase*], careia i se aplică apoi aceleași funcții de selecție a subiectului (acord, marcare, specificare), generând sintagma Verb-Obiect-Subiect, care este

de fapt clauza finită simplă (notată S). În fine, prin aplicarea asupra lui S (văzută ca sintagma Agr_SP) a unui complementizator C (sau marker clauzal, marker de discurs etc.) se obține o clauză “completă” ce poate, prin recursie, să ne ofere orice frază [Eng: *sentence*].

Alte exemple de X-bar scheme bazate pe MinP și P&P, ce pot fi discutate în contextul mai general al fenomenelor de guvernare sunt date de Fig. 2.2.3. care urmează.

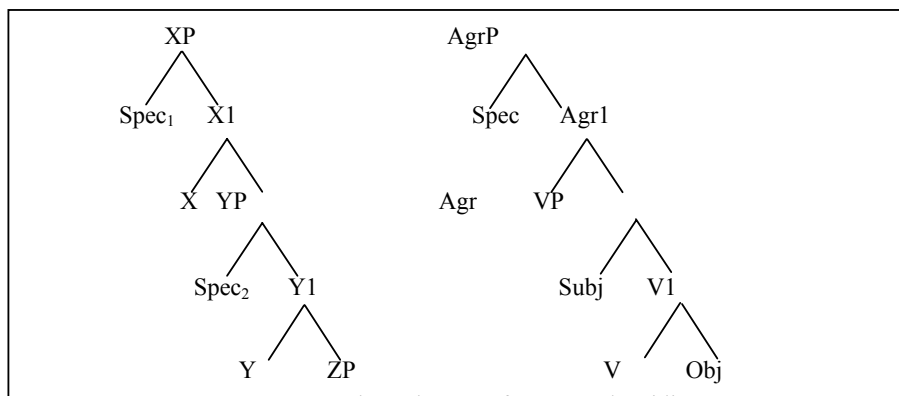


Figura 2.2.3. X-bar scheme în fenomene de "ridicare" la nivel de Spec în MinP

Concluzia este aceea ca teoria X-bar din MinP sintetizeaza relatiile fundamentale de dependenta, descrise de X-bar schemele propuse, si implicate în procesele de organizare a lexiconului si a sistemului computational din UG. X-bar teoria în abordarea MinP reflecta în principal *aspectele statice* întâlnite în fenomenele de *guvernare* (c-comanda, m-comanda, bariere, categorii de blocare etc.), în *teoria legarii* si în procesele de *referinta-coreferinta*, în stabilirea dependentelor la mare distanta (extra-clauzale) etc. Nu vrem sa intram în detalii (oricum complicate) si sa explicitam mecanismele de lucru ale X-bar schemelor considerate, ci mai curând sa atragem atentia asupra *teoriei X-bar* ca o *componenta fundamentala* a unei teorii lingvistice noi si elaborata cum este MinP si modelul sau P&P [3].

Teoriile MinP si P&P nu reprezinta un punct-terminus pentru evolutia teoriei X-bar. Dimpotriva, asigura o baza de pornire pentru o strategie radical diferita în care Chomsky examineaza cele mai serioase argumente pentru a abandona teoria X-bar [3; Cap. *Categorii si transformari*!] Aceasta alternativa si consecintele sale sunt discutate în sectiunea 5, si ar trebui sa reprezinte una dintre cele mai importante provocari prezente pentru domeniul analizei si proiectarii teoriilor lingvistice [21].

Unul dintre principalele scopuri ale sectiunii care urmeaza este de a introduce propunerea noastra de *scheme X-bar functionale* (*scheme FX-bar*) în cadrul strategiei lingvistice SCD. Propunerea noastra o consideram a fi o pozitie pragmatica si echilibrata în *directia* teoriei X-bar, atragând atentia asupra adevaratului sau rol si oportunitatilor computationale din lingvistica reala. Înțelegerea corecta a aspectelor *statice* si *dinamice* ale acestei versiuni a teoriei X-bar ar trebui sa fie de asemenea o consecinta a unei priviri cuprinzatoare a întregului context al teoriilor lingvistice care stabilesc principiile de dependenta, clasele de marcheri, categoriile si ierarhiile, regulile de referire si structurare, în strânsa relatie cu formele si regulile de constructie ale schemelor FX-bar.

3. Scheme X-bar functionale si strategia lingvistica SCD

În [19], în contextul *strategiei lingvistice SCD* (*Segmentare-Coeziune-Dependentă*) [22], [19], [1], [2], [6], este definită o clasă de *scheme X-bar augmentate* (*scheme AX-bar*), scheme destinate a reprezenta *invariantii sintactici generali* de reprezentare și operare cu structurile gramaticale ale LN, în particular pentru limba română, ca soluție la problemele de analiză și generare automată a LN. *Schemele FX-bar (functionale)* propuse aici completează și extind *schemele AX-bar* [19], și pot fi interpretate în mai multe moduri: (1) din punct de vedere *static*, schemele FX-bar pot furniza câteva de tipuri fundamentale de date pentru reprezentarea informației lingvistice în structuri de trasaturi lingvistice, standardizate și tipizate; (2) din punct de vedere *dinamic*, schemele FX-bar pot codifica informația lingvistică în forma procedurală ca *funcții și relații standard* ce sunt (recursiv) apelate în cadrul proceselor de analiză și generare a LN; (3) schema FX-bar generală poate fi de asemenea interpretată și utilizată ca un automat pe baza căruia să se realizeze o analiză *on-line* a textului unei fraze, cuvânt cu cuvânt.

3.1. Câteva preliminarii asupra SCD

Sunt necesare unele precizări asupra noțiunilor și notațiilor cu care lucrează *strategia lingvistica SCD*. Unul dintre elementele importante este că nivelul 2 ($\text{BAR} = 1$) în X-bar schema clasică joacă un *rol-cheie* în SCD pentru construcția structurilor sintactice, și este utilizat sub numele de *grup nominal* (NG), *grup verbal* (VG), *grup adjectival-adverbial* (AG), în general XG, pentru $X = N, V, A$. *Grupul XG* corespunde *proiecției lexicale* X_1 , cu $X = N, V, A$, și *clauzei minimale* CL_0 , în X-bar schema fundamentală propusă în Fig. 5.7.2.1.

Să menționăm că orice XG (X_1) este un XP (X_2), dar nu și invers, deoarece nivelul proiecției categoriei X în cazul XG lucrează numai pentru nivelul $\text{BAR} \leq 1$. SCD face de asemenea distincție între câteva tipuri de NGs (NGs elementare, predicative, non-predicative, etc.), VGs (VGs la un timp finit și la un timp non-finit) etc.

O altă trasatură esențială și specifică a SCD este un tratament adecvat al *proprietăților functionale* ale categoriilor lingvistice, ca și al tuturor *categoriilor relationale* și sintagmelor (expresiilor) de *discurs*. Mecanismul utilizat pentru a obține acest lucru se bazează pe *clase de marcheri lingvistici* și *ierarhiile* lor [1], [2], [6]. Câteva observații se impun:

(a) *Marcherii* din SCD, numiți *marcheri de structuri sintagmatice* (PS-Ms) [Eng: *phrase-structure markers*], sunt cu totul diferiți de ceea ce teoria lui Chomsky numește formal "*marcheri de sintagma*" [Eng: *phrase-markers*] în [17], sau *T(ree)-marcheri* în [3]. *Marcherii Chomsky* sunt definiți ca "*taieturi orizontale*" (sau "factorizari") în cadrul unui arbore de derivare, sau ca fiind arborele însuși. Mult mai apropiați de ceea ce sunt PS-Ms în HPSG [16], *marcherii de structuri sintagmatice* (PS-Ms) din SCD sunt acele categorii lexicale și nelexicale care se aplică cuvintelor și structurilor sintagmatice (PSs) cu scopul de *evidențiere*, de a *marca*, anumite funcții și relații sintactice și semantice pe care PSs

respective le joaca în cadrul unei exprimari. Punerea în evidenta a anumitor functii care se aplica PSs se refera la (cel puțin) câteva elemente: *tipul functiei* (sintactic, semantic, relational, logic, pragmatic, discursiv etc.), *locul*, în text, *unde începe aplicarea* functiei sau relatiei, si *domeniul* (*domeniile, conexe sau nu*) de aplicare a functiei sau relatiei (limitele textuale între care se aplica).

Exemple tipice de PS-Ms din SCD sunt: (a) trasaturile *predicative* generate de catre *categoriile predicationale* (de fapt, verbe, substantive, adjective si adverbe predicationale); (b) acele *mijloace gramaticale* prin care sunt introduse *noi* NGs (grupuri nominale în limbajul SCD), VGs, AGs (*Caz-marcarrea*, acordul, gradele de comparatie, etc.); (c) acele *categorii si expresii* (numite si *marcheri de discurs*) care introduc *noi* clauze; (d) PS-Ms care introduc *proprietati relationale* asupra PSs si clauzale (de exemplu de *marcheri* de tip logic cum sunt structurile *daca-atunci-altfel*, *deoarece*, etc., dar si *marcheri* de tip sintactico-semantic cum sunt aceia care introduc *categorii si clauze subordonate* etc.)

(b) SCD se aseamana din unele puncte de vedere cu abordarea [16] a HPSG si, partial, cu [15], care exploateaza, pentru prima oara în clasa teoriilor lingvistice bazate pe *gramatici* de PSs (PS-Gs), într-o mult mai mare masura, categoria lingvistica a *marcherilor* PS-Ms. În [16], Pollard & Sag "postuleaza o noua parte a *marcherilor* de discurs, ... ce se remarca ... printr-un nou atribut al categoriilor (în plus fata de NUCLEU si SUBCAT) numita MARKING, cu valori din sortul *marking*". Teoria HPSG enunta PRINCIPIUL MARCARII [16, p. 400] dupa cum urmeaza:

"Într-o sintagma cu nucleu, valoarea trasaturii MARKING este lexical-identica cu cea a trasaturii MARKER-DAUGHTER daca aceasta exista, si cu cea a trasaturii HEAD-DAUGHTER în caz contrar.

Modul în care HPSG [16] pune la lucru PS-Ms reprezinta un bun si esential pas înainte, desi credem ca nu exploateaza îndeajuns potentialul functional si relational al diferitelor clase de *marcheri* si ierarhiile acestora (asa cum face strategia SCD, vezi si [7]).

(c) Continuând si extinzând constructia limbajului, ca o expresie de convergenta între gramatica categoriala si *Minimalist Program*, Chomsky [3] considera *transformarile generalizate* (GTs) si concepe un demers de înlocuire a X-bar teoriei, ce explica în Programul Minimalist structura constituentilor (sintagmatici) complecsi, prin GT *Merge* care construiesc obiecte sintactice pornind de la obiecte sintactice simple (de exemplu, "*speaks*" si "*French*" sunt "reunite" într-un nou obiect sintactic "*speaks French*" etc.). Mai multe formalizari ale acestui nou curent al ideilor lui Chomsky pot fi gasite în cadrul gramaticilor logice multi-modale si de tipuri categoriale, e.g. [21], [23], [24] etc. (vezi si sectiunea 5).

(d) Dintr-o perspectiva diferita dar oarecum similara, *gramatica functionala* (FG) [25] a lui Simon Dik, orientata functional si semantic, încearca sa faca aceleasi lucruri. Ca si în SCD, FG gaseste *patru tipuri ierarhice* de baza ale categoriilor relationale, aceste tipuri corespunzând într-o buna masura cu *clasele de marcheri* PS-Ms si *ierarhiile* lor stabilite în

SCD [7], [2], [6]. PS-Ms reprezintă acele mijloace lingvistice de “suprafata” pe care le utilizează un limbaj natural pentru a organiza sintactic și semantic structurile codificate în construcții gramaticale. Se impune în viitor o analiză comparativă între cele *patru nivele* sau “*straturi*” din organizarea formală și semantică furnizată de FG [25], și cele *patru nivele* de proiecție lingvistică, împreună cu clasele de marcheri corespunzătoare, din SCD: (1) *cuvântul* (lexical); (2) *sintagma* XG (X = N, V, A) subclauzala; (3) *clauza* (finită și infinită); (4) *discursul* (una sau mai multe fraze, care să formeze un *segment* de discurs).

(e) În fine, privitor la utilizarea intensivă a *caracterului predicational* pe care categoriile lexicale majore (N, V, A) îl poartă (proprietate moștenită sau dobândită apoi de alte categorii gramaticale), *strategia lingvistică* SCD este comparabilă în special cu FG, cu accentul particular pe ierarhiile de delimitare și marcare aplicate structurilor sintactico-semantică. SCD porneste de la *lexicon* și stabilește la acest nivel o *taxonomie predicatională inițială* pentru categoriile lexicale majore. Un exemplu simplu al acestei taxonomii predicative este dat de către cele două categorii importante de substantive comune: *substantive existențiale* sau *obiectuale*, a căror *predicationalitate* (*trasatura* PRED) este EXIST (e.g. [Eng: *student, table*; Rom: *elev-student, masa*]) și a căror reprezentare funcțională reflectă categorii individuale sau personale, de exemplu predicatul uni-variabil *student(X), masa(X)* etc., și substantive de tip-predicational, a căror predicationalitate (*trasatura* PRED) are valoarea ACT, e.g. [Rom: *întâlnire, invidie, marcare* etc.], și ale căror reprezentări funcționale depind de mai multe variabile, de exemplu *întâlnire(X, Y,...), invidie(X, Y,...), marcare(X, Y)* etc. Substantivele proprii și/sau personificările sunt codificate prin constante ale variabilelor din predicatul de mai sus. Câteva din remarcile anterioare vor fi aprofundate în concluziile finale ale lucrării.

Schemele FX-bar, ca și precursorii lor *schemele AX-bar* [19], reflectă pentru SCD faptul că un XPG (grupul sintagmatic de nucleu X), sau mai simplu XG, conține un *nucleu*, reprezentat printr-o categorie lexicală (nevidă) sau printr-o categorie virtuală (vidă), înconjurat (prin relații de *coeziune*) de specificatori și/sau modificatori de tipul A (adjectival-adverbial). Este esențial să facem următoarea specificare: un XG din SCD nu include nici un complement (argument obligatoriu) sau adjunct. Complementele și adjunctii, împreună cu *nucleele* de nivel BAR = 1 formează nivelul BAR = 2 în FX-bar schema propusă în Fig. 3.2.1. Pentru un anumit nivel de specificare semantică, FX-bar schemele nu fac o distincție clară între complemente (argumente obligatorii) și adjuncti, considerând toate structurile subcategorizate ca fiind argumente sintactice; clasificări ulterioare (suplimentare) sunt făcute pe baza sabloanelor verbale și restricțiilor sintactice, semantice, și pragmatice asupra componentelor sablonului, la nivel de lexicon.

O problemă a cărei soluție poate influența în mod special și teoria X-bar este aceea a asignării corecte a complementelor și adjunctilor, în particular, a stabilirii corecte a dependentelor dintre grupurile nominale (NGs). Soluția acestei probleme nu se poate obține la nivel sintactic, iar o soluție completă nu se poate obține uneori nici chiar în contextul unui nivel semantic minimal (vezi [26], [27]). Chomsky remarcă realitatea că “... *the distinction between modifiers and arguments is notoriously difficult in certain cases*” [9, p. 44]. Exemple simple ilustrează această problemă: în TBarr [8], sintagmele “*the students of physics*” este

vazuta ca un NP cu un *argument* PP, în timp ce sintagma "the students in the yard" este considerata a fi un NP cu un *adjunct modifier* PP. De fapt, în numeroase LNs, inclusiv engleza, se pot aduce multiple argumente serioase pentru ca cele doua sintagme sa poata fi la fel de bine interpretate fie într-un fel, fie în celalalt.

Solutia SCD pentru acest exemplu foarte particular este urmatoarea (schitând si solutia problemei generale): substantivul "students" este obiectual, adica *nu are* o natura *predicationala* prin el însusi, astfel ca ambele sintagme nominale care îl succed sunt considerate de catre SCD ca fiind *modificatori* pentru NG "students". Natura acestor modificatori poate fi diferita deoarece "physics" este introdus de markerul de caz (genitiv) "of", în timp ce "the yard" este introdus de markerul *prepositional* "in". În general, când nucleul lui NG posedea o trasatura *predicationala*, atunci NG care urmeaza nucleului predicational asigura o distributie sintactica ce satisface un anumit sablon (verbal) al predicatului (verbului) corespunzator.

Clasele din PS-Ms si ierarhiile lor din SCD [7] sunt responsabile pentru delimitarea structurilor sintagmatiche propuse de schemele FX-bar, si pentru stabilirea dependentelor sintactico-semantice. Diferitele tipuri de markeri sunt adesea aplicate simultan (deci multiplu) asupra acelorasi categorii gramaticale, în cadrul anumitor nivele de structurare (proiectii pe BAR-nivel). Similar cu unele teorii lingvistice (LFG, FG, si partial HPSG) dar contrar altora (GB, GPSG etc.), SCD nu considera *prepozitia* ($X = P$) ca fiind o categorie lexicala majora. În SCD, P primeste rolul unui marker (functional), având atât proprietati de marker de caz cât si de complementizator. Categoriile HPSG PP[+PRD] sau PP[-PRD] (vezi [16]) sunt irelevante pentru SCD deoarece trasatura +PRD în HPSG este atribuita numai lui PP subcategorizat de un V, în timp ce trasatura (*predicationala*) PRED din SCD poate fi în mod egal atribuita lui V, N, sau A (la nivelul lexiconului, cel puțin) dar nu si lui P.

În S-C-D proprietatile de subcategorizare sunt exploatate *ab initio*, la nivelul de organizare al lexiconului, pe baza *trasaturii functionale* PRED de *predicationalitate*, asignata sau nu, unora din categoriile sintactice majore N, V, A. Observatii lingvistice empirice ne-au convins, înca de la începuturile cristalizarii SCD [22], ca o *taxonomie* functionala si predicativa adecvata ar trebui sa reprezinte punctul de plecare al oricarei teorii lingvistice, atât din motive teoretice cât si pragmatice, si ca multiple abordările actuale (cum ar fi [27]-[32]) aduc o sustinere puternica pentru multe din ideile esentiale din SCD, în special folosirea intensiva a predicativitatii si functionalitatii descrierilor lexical-semantice ale categoriilor lingvistice atât în *procesarea automata* a LN cât si în cadrul *bazelor de cunostinte lexicale*.

[19] propune urmatoarea specificare a *Principiului Proiectiei Maximale* (PMP) [Eng: *Principle of Maximal Projection*], ca un pas important catre folosirea intensiva a trasaturilor predicationale (functionale) ale categoriilor lexicale majore în SCD. Propunem aici

O specificare a PMP (forma actualizata):

Proprietatile de subcategorizare ale categoriilor sintactice majore N, V, A depind de trasatura lor lexical-semantica PRED(icativity), cu valorile ACT si EXIST, si de trasatura lor morfo-semantica TENS(e), cu valorile FINI(te) si INFI(nite).

Trasatura PRED, atribuita categoriilor majore N, V, A la nivel de lexicon, primește două valori: valoarea ACT, pentru acele categorii care au *proprietati predicationale* (în literatura este folosit adesea termenul “*deverbale*”), și valoarea EXIST, pentru acele categorii N, V, A cu caracter *existential*, obiectual, non-predicational. Trasatura TENS primește valorile FINI(te) pentru acele forme ale categoriei V care posedă un timp sau aspect finit, personal, și valoarea INFI(nite) pentru toate celelalte categorii. Exemple:

[Eng: *boy, pencil*; Rom: *baiat, pix*] PRED:= EXIST; și TENS:= INFI;

[Eng: *attempt, showing, proved*; Rom: *încercare, aratând, demonstrat*]

PRED:= ACT; și TENS:= INFI;

[Eng: *are*; Rom: *sunt*]

PRED:= EXIST; și TENS:= FINI;

[Eng: *gives*; Rom: *da*]

PRED:= ACT; și TENS:= FINI.

Într-un grup verbal VG reprezentând un compus la un timp finit, valorile “pozitive” de trasaturi, cum sunt ACT sau FINI sunt *mostenite* de la nucleul V al VG de către întreaga sintagma VG, sau pot fi obținute *cumulativ* prin *proiectia* morfo-sintactica.

Specificarea PMP de mai a funcției *proiecției maxime* este necesară în SCD deoarece în multe LNs, inclusiv în română, *calitatea deverbala* (predicatională, deci funcțională) a categoriilor lexicale tradiționale *non-verbale* cum ar fi N și A trebuie descoperită cât mai devreme posibil și asignată la nivel de lexicon. De exemplu, în engleză, deși pentru substantivele care ‘*verbalizează*’ în “-ing” valoarea trasaturii lor TENS este INFI, aceste substantive posedă, pentru trasatura PRED, aceeași valoare ACT sau EXIST pe care o au verbele din care provin substantivele (sau gerunziile) în “-ing”, și astfel posedă *aceleași* proprietăți de subcategorizare ca ale verbului de origine.

3.2. Ipoteze de lucru și aspecte caracteristice ale FX-bar schemei

Continuând ideile de bază ale schemelor AX-bar din [19], propunem, pentru SCD, FX-bar schema generală din Fig. 3.2.1. Muchiile din stânga conțin *noduri* cu rol *funcțional* sau *relational*: marcheri, cuantificatori, specificatori, modificatori (eventual adjuncti). Pentru a obține reprezentări sintactice și semantice corecte, nodurile funcționale se aplică (recursiv) *nucleelor* X_k și CL_k, k = 0, 1, 2, iar nucleele, cu rol funcțional (predicational, X₁) sau relational (eventual X₂), au ca argumente clauze infinite (complemente, X₁) sau finite (X₂). Precizăm că la acest nivel nu se poate face distincția dintre complementele COMPL_i (argumente obligatorii) și adjuncti ADJCT_i (argumente optionale). În mod normal, în Fig. 3.2.1., ADJCT_i sunt “amestecați” printre ARG_j, la nivel sintactic nefiind discernabili de complementele obligatorii ale unui nucleu predicational. Poziția funcțională (la stânga nodului X₁) a nodurilor ADJCT poate rezulta doar în urma unor calcule semantice și pragmatice suplimentare, din care se obține rolul tematic al argumentelor ARGs ale lui X₁.

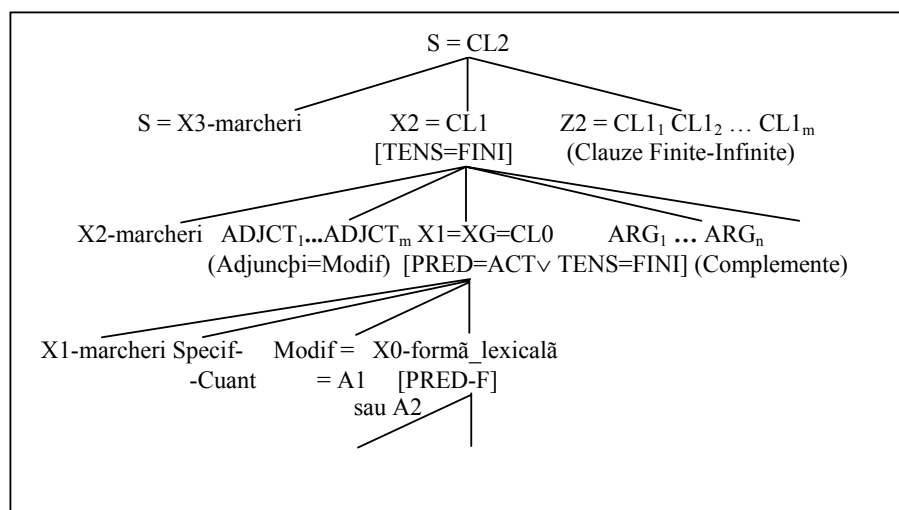


Figura 3.2.1. Schema (funcțională) FX-bar generală

(♣) *Aspecte specifice ale schemei FX-bar propuse:* (♣1) Sunt permise un număr arbitrar de *argumente* (sau *sateliti* în sensul [10], [31]), toate notate cu ARGs. În SCD, ARGs sunt formate din *complemente* obligatorii (COMPLs) și din *adjuncti* (ADJCTs), sau *complemente optionale*. ADJCTs pot fi reprezentați la nivel sintactic tot ca argumente ale nucleului, însă la nivel semantic ADJCTs au rol de *modificatori* ai nucleului. Notatia “A-pozitie” din teoriile Chomskyene, care înseamnă ARG-pozitie, nu trebuie confundată cu notatia noastră pentru categoria A = adjectiv-adverb. În teoriile și notatia lui Chomsky, COMPLs sunt în A-pozitie (ARG-pozitie), în timp ce ADJCTs nu. SCD se situează pe o poziție sintactică similară cu HPSG [16], care utilizează lista SUBCAT pentru a codifica toate sintagmele pe care le subcategorizează un nucleu semantic, adică atât COMPLs cât și ADJCTs (sau ARGs din SCD). (♣2) Sintagmele AG = A0 sau A1, sau AP = A2 sunt *postulate* de către SCD ca fiind de tipul categoriei *functionale* Modif, manifestate prin categoriile A (de nivel X0, și aplicabile la nivel X0), ADJCTs (de nivel X1, și aplicabile la nivel X0 și X1), și clauza relativă (de nivel X2, și aplicabilă la nivel X0 și X1). (♣3) Categoria generică Specif (sau Spec), în care intra cuvintele și sintagmele ce desemnează *cuantificatori* de toate tipurile (generalizati), determinatori (în particular), este postulată de către SCD ca fiind o *categorie funcțională* ce poartă trasaturi de natură *cuantificatională* la nivel lexical (în particular, *negatia* la nivelul X1), inclusiv articularea (hotărâta sau nu), suprapunându-se deci uneori peste X1-marcheri de trasaturi funcționale cum este *acordul*. Relațiile (funcționale) de *acord* sunt esențiale pentru *coeziunea locală și globală* în cadrul strategiei SCD: acordul dintre X0-Modif și X0-Specif cu nucleul X0 (la nivel X1), acordul Nucleu-Subj (sau chiar Nucleu-COMPL) și acordul COMPL-Pron_{Emfat} (Pronume emfatic) (la nivel X2), o anumită *corespondență* a timpurilor evenimentelor într-o clauză și între clauze. Aceste tipuri de relații de acord, referință și coreferință, coeziune, coerentă, etc. sunt responsabile pentru o largă clasă de dependente locale și globale, inclusiv dependente la distanță mare și în extra-pozitie. Accentul în componenta de *coeziune* a strategiei SCD (Segmentare-Coeziune-Dependentă) cade pe mijloacele *sintactice* și de “*suprafata*”, mai

curând decât pe cele semantice, încercând să găsim, să extragem, și să utilizăm într-o măsură maximală informații de ordin superior, cum ar fi informația de discurs [34], pragmatică, semantică etc. (♣4) Sintagma tradițională PP din teoriile lingvistice clasice, iar în SCD, grupul prepositional PG (format dintr-un grup nominal NG care este precedat de o prepoziție sau o locuțiune prepositională) este întotdeauna considerată un ARG (COMPL sau ADJCT) în FX-bar schemele al căror nucleu (lexical nevid sau vid) este N, V, A. Această ipoteză de bază asupra PG este justificată de SCD prin faptul că P este considerată o categorie majoră, adică o categorie de nivel X1 în schema FX-bar din Fig. 3.2.1. și doar o categorie de nivel X0. Proprietățile de subcategorizare ale N, V, A (dar nu și P) pot fi asignate *ab initio*, la nivel de lexicon, începând cu trasatura lexicală PRED a categoriilor predicative. Categoria P poate primi proprietăți funcționale, cel mai adesea ca *marcher de caz*, uneori proprietăți *relationale* (de exemplu [Eng: *on*; Rom: *asupra*]), dar nu și proprietăți de subcategorizare. (♣5) Subiectul (Subj) în SCD, lexical nevid sau vid (PRO), este considerat ca un argument special al proiecțiilor maxime ale categoriilor X = N, V, A într-o clauză finită (de nivel X2) sau infinită (de nivel X1). (♣6) În ipotezele (♣5) și (♣2) de mai sus, categoria lingvistică tradițională VP este dizolvată într-un grup verbal VG (finit sau infinit), înconjurat (de cele mai multe ori urmat) de nucleu de ARGs și formând o clauză finită, respectiv infinită. (♣7) *Teoria limitării* și multe probleme majore legate de TBarr [8], [9], [17] sunt explicitate și rezolvate în cadrul realizat de SCD și schemele FX-bar, în principal datorită delimitării clare a funcțiilor și relațiilor care se aplică cuvintelor și sintagmelor, a reprezentării lor lexicale prin clasele de PS-Ms, și a specificării domeniului lor de aplicare. Acest rol este realizat explicit în cadrul claselor și ierarhiilor de marcheri propuse și utilizate de SCD [2], [6], [7]. Trebuie să remarcăm că lucrările sale cele mai recente [34], [35], Chomsky adoptă o tehnică similară de “limitare” a operațiilor de *construire* [Eng: *merge*] și *transformare* [Eng: *move*] doar la “domeniul” sintactic al unei “faze” [Eng: *phase*], o unitate textuală (care în general coincide cu clauza!) în care Chomsky propune următorul *principiu de impenetrabilitate* “Într-o fază (clauză n.n.) F cu nucleul H, domeniul lui H nu este accesibil la operații în exteriorul lui F, ci este accesibil numai H și muchia sa (nodul sau ascendent)” [34]. Exact așa este construită și funcționează schema FX-bar! De asemenea, fenomene de *teoria legării* [9], [8], [3], [16], *legăturile* [Eng: *linking*] din [27], mecanisme de *coeziune* (locală și globală) și *discurs* întâlnite în [36], [31], [33], etc. sunt mai ușor de pus în evidență și de rezolvat în cadrul oferit de strategia lingvistică SCD și teoria FX-bar.

(♦) **Observații asupra ipotezelor de lucru pentru schema FX-bar din Fig. 3.2.1.:**

(♦1) Schema FX-bar este proiectată să lucreze în asociere cu un parser care este capabil să recunoască clasele de PS-Ms și structurile sintagmatice considerate de strategia lingvistică SCD. Schema FX-bar este organizată pe *patru nivele* de proiecție $BAR = 0 \div 3$ (deasupra nivelului de lexicon, notat conventional $BAR = -1$); *trei nivele X0-X1-X2* corespund proiecției dintre nivelul lexical ($BAR = 0$) și nivelul *clauzal*, al structurilor *uni-eveniment*; alte *trei nivele CL0(=X1)-CL1(=X2)-CL2* corespund proiecției dintre nivelul *clauzal minimal* $CL0 = X1$ și nivelul *frazel*, al structurilor *multi-eveniment*. Nivelele uni-eveniment X0-X1-X2 exprimă predicatia clauzei (propoziției) simple în care sunt distribuite categoriile lexicale de bază și sintagmele pe care le generează, în timp ce nivelele CL0-CL1-CL2 exprimă relațiile logice și predicative (de ordinul doi) dintre clauzele simple. Schema FX-bar lucrează într-o manieră recursivă (top-down sau bottom-up), atât în situațiile de analiză cât și în cele de

generare în care este antrenat parserul asociat, în strânsă cooperare cu strategia lingvistică SCD, cu clasele de PS-Ms și ierarhiile lor și, mai ales, pe baza *meta-algoritmilor* SCD de analiza-generare [1], [2], [6], [7]. Sa mai observăm că FX-bar schema din Fig. 3.2.1. poate fi utilizată independent de așa numita *ordine canonică* (sau *sistemică*) a cuvintelor și sintagmelor dintr-o clauză, specifică fiecărui LN [37], [38]. (♦2) Valoarea ACT de trasatură (funcțională) pentru categoriile N și A (și implicit V) este atribuită acestor categorii la nivelul lexicon atunci când ele corespund unor evenimente cu actanți și/sau stări multiple. Valoarea EXIST este implicit sau explicit introdusă de formele și înțelesurile verbelor existențiale (*a fi*), modale (*a trebui*), etc. (♦3) Trasatura (funcțională) TENS este similară cu categoriile virtuale I (INFL) și T (Tense) din teoriile GB și TBarr ale lui Chomsky și din schemele S-bar corespunzătoare (Fig. 1.3. și Fig. 2.2.3.). Pentru un VG finit (TENS = FINI), structura V2 corespunzătoare devine clauză finită clasică. Dacă sintagma XG (X1) este un grup a cărei categorie-nucleu X posedă valorile de trasaturi PRED = ACT și TENS = INFI, atunci XG devine noul nucleu al unei clauze infinite ce face parte dintr-o structură de nivel X2 (XP). (♦4) Poziția specială a *subiectului sintactic* (Subj) este considerată de către SCD atât o ARG-pozitie (asemanătoare, de fapt, cu o COMPL-pozitie) cât și o *Caz*-poziție. În concordanță cu TBarr [8] și cu HPSG [16], Subj primește poziția specială a *primului element* din lista SUBCAT [16]. Aceasta este în esență o poziție sintactică, iar Subj poate primi o funcție tematică (θ -poziție) autentică doar ca rezultat al unor calcule sintactice și semantice suplimentare. (♦5) Așa cum rezultă din schema FX-bar din Fig. 3.2.1., sintagmele AP și PP din teoriile lingvistice clasice sunt segmentate de către marcherii SCD [7] în sintagme mai mici XG, X = N, V, A. Așa cum am precizat deja, SCD atribuie noilor sintagme următoarele roluri: AG = Modif, cu rol funcțional la nivelul de proiecție X1, și PG = ARG (COMPL sau ADJCT), ADJCT purtând de asemenea rol de Modif al nucleului de nivel X2. PG devine deci un NG P-marcant, iar orice categorie A are de la început reprezentarea (nesaturată) A(X), unde X = N, V, A este nucleul (existent, viitor, sau lipsind pur și simplu) al sintagmei de nivel X1 în care Modif = A. În mod similar, orice categorie Specific (determinator, cuantificator, etc.) joacă un rol similar, schema FX-bar impunând reprezentarea funcțională Specific(X), unde X este nucleul sintagmei. (♦6) În ciuda anumitor asemănări (inerente) între schemele FX-bar și versiunea MinP a teoriei X-bar, există diferențe de bază în ce privește organizarea și funcționarea constructivă dintre schemele (F)X-bar din Fig. 3.2.1. și Fig. 2.2.1. De exemplu, în schema FX-bar, fiecare element lexical se proiectează într-o categorie obiectuală sau funcțională (relațională), aceasta este (coeziv și recursiv) înconjurată de către Specific și/sau Modif, iar dacă valoarea ACT a trasaturii PRED a nucleului este prezentă, atunci această valoare ACT este mostenită de către întreaga sintagma al cărei nucleu a fost specificat sau modificat. Aceasta sintagma cu nucleu predicțional își subcategorizează complementele (argumentele obligatorii COMPLs) și adjunctii ADJCTs (care modifică sintagma-nucleu). În schema X-bar din Fig. 2.2.1., se întâmplă tocmai invers deoarece “*The Head-Complement relation is the “most local” relation of an XP to a terminal Head Y, all other relations within YP being Head-Specifier (apart from adjunction, ...)*” [3: p. 53]. (♦7) Deși schema FX-bar generală a fost proiectată având în vedere în primul rând limba română, ea poate fi aplicată pentru a reprezenta, grafic și logic, structuri sintactico-semantice ale LNs cu valori ale parametrilor gramaticali foarte diferite, cum ar fi engleza-germana sau franceza-germana. Distribuția complementelor (argumentelor) în română (engleza, franceza) poate fi foarte

diferita de cea din germana; de exemplu, într-o clauza al carei verb principal din compusul sau verbal VG se afla în poziție finală, sau pentru o categorie A (adjectiv-adverb) având valoarea de trasatură PRED = ACT.

Ex. 3.2.2.R. /Paharul /spart / /de Ion/ cu mingea /de fotbal/

Ex. 3.2.2.E. /The glass /broken / /by Ion/ with / the football /

Ex. 3.2.2.G. /Das /von Ion /mit /dem Fußball / /zerbrochene / /Glass/

Dupa cum am remarcat în (♦1), schema FX-bar poate fi utilizată independent de regulile structurilor sintagmatice și ordinea lor (din româna sau germana), aceasta deoarece principiile rămân aceleași și diferă numai anumii parametri și valorile lor pentru LNs distincte: în româna (și engleza, franceza) argumentele succed o categorie A ce reprezintă un nucleu predicțional, în timp ce în germana ele îl pot (!) precede. Dacă un nucleu V al unei clauze are valorile de trasaturi PRED = ACT și TENS = FINI, atunci distribuția ARGs este similară cu cea din româna, cu posibile (și probabile) diferențe impuse de *ordinea sistemică*, strict dependentă de LN, a ARGs (a se vedea [37] dar și [27]).

Dacă se încearcă utilizarea formei FX-bar ca “schelet” pentru un automat (sau gramatica formală) de analiză și generare a LN, un asemenea automat ar trebui să mimeze atât forma generală a schemei FX-bar cât și regulile gramaticale de analiză-generare. Partea din automat care reflectă cele patru nivele de organizare a structurilor LN în schema FX-bar ar trebui să fie independentă de limbaj (cel puțin pentru o largă clasă de limbaje europene), în timp ce (sub)partea constituenta care recunoaște structurile lingvistice pe fiecare nivel individual X_k ($k = 1, 2, 3$) trebuie să fie dependentă de limbaj (acest fapt este binecunoscut și parametrizat). Reprezentarea schemei FX-bar pentru Ex.3.2.2.G. este aceeași cu reprezentările FX-bar pentru Ex.3.2.2.R.-E., și similară cu figura pentru Ex.4.2.R.-E.

4. Exemple de aplicare a schemelor FX-bar

Vom expune câteva exemple de aplicare a schemelor FX-bar la reprezentarea sintagmelor, clauzelor și frazelor. În exemplele prezentate, categoriile gramaticale pentru care PRED = ACT sau TENS = FINI vor fi *subliniate*, iar PS-Ms care se aplică sintagmelor X_k ($k = 0, 1, 2$) sunt reprezentați grafic în text prin apariția unui sau mai multe semne ‘slash’ /. Să notăm că schemele (augmentate) AX-bar din [19], deși oarecum asemănătoare în spirit sunt efectiv scufundate în schema FX-bar generală, diferențele substanțiale constând în forma unitară a FX-bar schemei și în criteriile sintactice și logico-semantice mai clare, pe baza cărora clasele de PS-Ms și ierarhiile lor sunt explicit propuse și aplicate în funcționarea schemei FX-bar.

Care este relația dintre exemplele de FX-bar scheme și formulele logice atasate după reprezentarea grafică? *Prima formulă* este o reprezentare uzuală a LN, care folosește limbajul logicii predicatelor, reprezentare mai apropiată de exprimarea în LN, conținând toate variabilele ce exprimă referințele-coreferințele, dar (pentru simplitate) fără cuantificatorii corespunzători. *A doua formulă* este traducerea mai completă a primei formule în limbajul de programare logică Prolog, folosind tehnici clasice de reprezentare a

cunostintelor de LN în Prolog. Pe o scala ascendenta a masurii în care sintagmele LN ar fi analizate, *schema* FX-bar poate fi vazuta ca un prim rezultat al procesului de parsare (analiza), *prima formula* ar urma procesului de parsare, incorporând fenomenele de referinta (si coreferinta, rezolutie a anaforei, etc.), iar *a doua formula* ar reprezenta o rafinare a primei formule. Formulele de tipul doi reprezinta de asemenea atât un stadiu final al procesului de analiza a frazei cât si punctul de pornire în procesul de generare a unei fraze (conform cu abordarea [39], [6] a generarii automate a LN, însa diferita de [31], de exemplu).

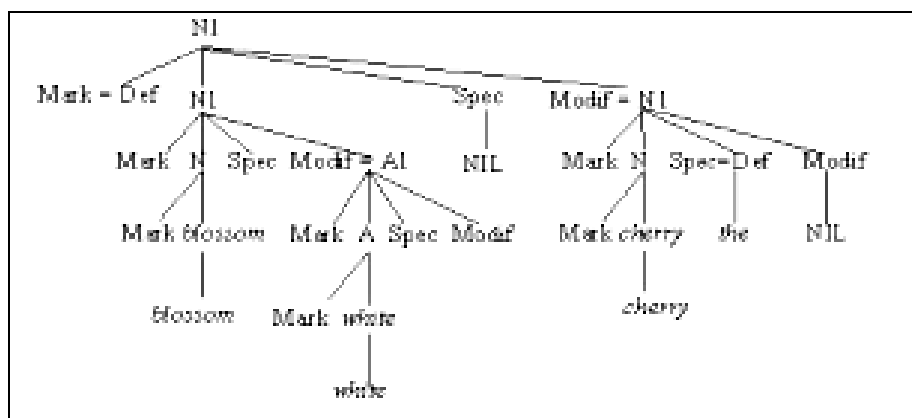
Este important sa remarcam ca schema FX-bar propusa reflecta, în principal, relatiile de *dependentă* dintre diferitele categorii, sintagme, si clauze dintr-o fraza, împreuna cu markerii corespunzatori care controleaza, în parte, si comportamentul lor distributional. Deoarece am vazut în ce masura schemele ordinea argumentelor este (parametric) dependentă de limbaj în schemele FX-bar, acestea pot codifica nu numai situatii în care argumentele succed (situatia obisnuita) sau în care ele preced nucleul lor semantic (Ex.3.2.2.), dar si în care argumentele aceluiasi nucleu sunt interschimbabile. Deci aceleiasi schema FX-bar i se pot atribui mai multe formule logice corespunzatoare "echivalente".

4.1. De la text la scheme FX-bar

Strategia SCD propune urmatoarele scheme FX-bar pentru exemplele de mai jos. Desi muchiile ale caror noduri sunt Modif sau Specif sunt situate în dreapta nucleului corespunzator (pentru conveniente grafice), ele trebuie înțelese ca având rol functional (situate la stânga si aplicându-se nucleului X1). La fel si cu unii adjuncti, la nivel X2. Diferentele dintre codificarea formei pentru engleza si cea pentru româna sunt nesemnificative (cu exceptia unor aspecte suplimentare de acord, care sunt puse în evidenta). Forma codificata a textului pentru limba engleza este un argument suplimentar pentru versatilitatea schemelor FX-bar propuse.

Ex. 4.1.1.R. /floare albă / de cireș /

Ex. 4.1.1.E. /the cherry / white blossom /



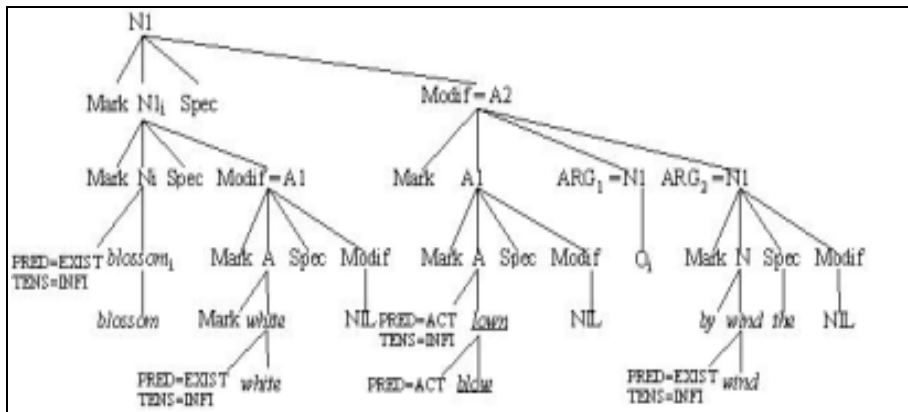
LR.4.1.1.R. de(cireș)(albă(floare(X)));

LR.4.1.1.E. quant(indef, X, white(blossom(X)), cherry(X)).

Ex. 4.1.2.R. /floare albă /, // bătută // de vânt /

Ex. 4.1.2.E. /the white blossom /, // blown // by the wind /

object_i = O_i; event_j = e_j



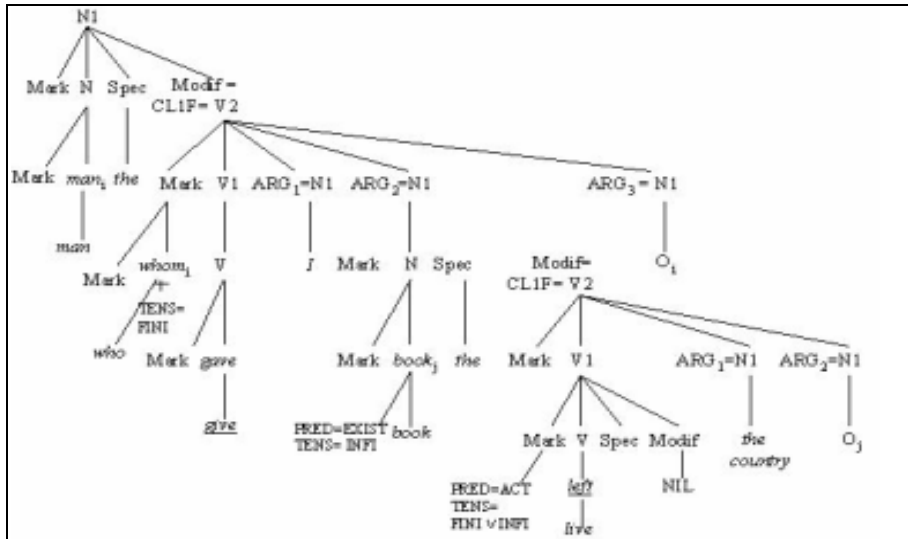
LR.4.1.2.R. albă(floare(X)) \wedge bătută(de(vânt(Y)), X);

LR.4.1.2.E. quant(indef, X, white(blossom(X)),

quant(indef, Y, by(the(wind(Y))), blown(Y, X)).

Ex. 4.1.3.R. // educat // [de tatăl său] // corespunzător // cu vechile principii /

Ex. 4.1.3.E. // educated // [by his father] // accordingly // with old-fashioned principles /

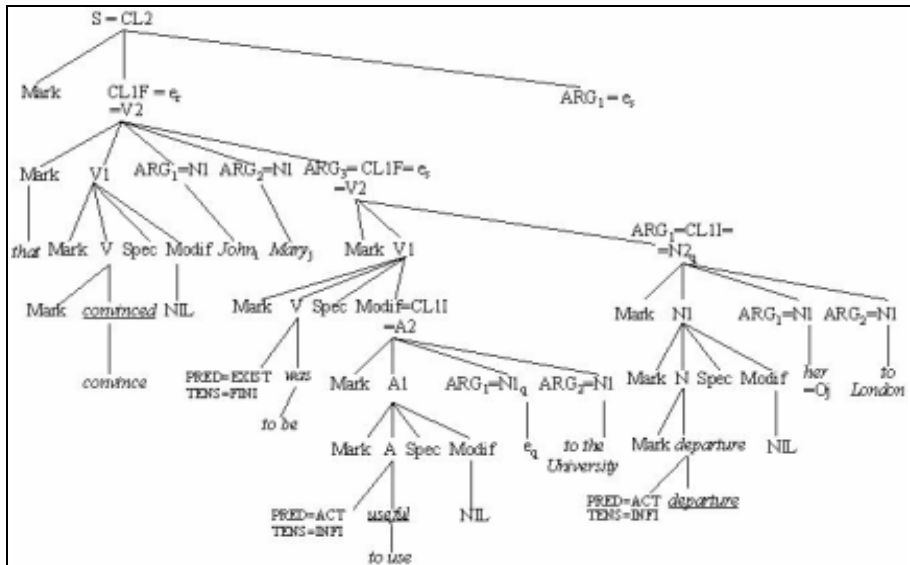
4.1.4.E. **Reading2** (*left* = past_participle(*leave*))4.1.4.E. **Reading3** ([Eng: *left*] = [Rom: *stânga*]) . .

LF.4.1.4.R. a-părăsit(omul(X) \wedge am-dat(Y, cartea(Z), X), țara(T);
 LF.4.1.4.E. quant(def, X, and(man(X), quant(def, Y, I(Y), quant(def, Z, book(Z), gave(Y, Z, X))), quant(def, T, country(T), left(X, T))).

Ex. 4.1.5.R. // Ion_i // a convins-o_j // pe Maria_j // că // deplasarea_k // ei_i // la Lodra // a fost utilă // e_k // Universității . //

Ex. 4.1.5.E. // John_i // convinced // Mary_j // that // her_j // departure_k // to London // was useful // e_k // to the University . //

object_i = O_i ; event_j = e_j



- LF.4.1.5.R. a-convins(ion, -o(pe(maria;)), că(a-fost-utilă(
deplasarea(ei(X;), la(londra), universității(Y))));
- LF.4.1.5.E. convinced(john = X, mary = Y, quant(def, X, her(X), departure(X, to(
london)) = E, quant(def, Z, university(Z), was-useful(E, Z))))).

4.2. Observatii generale

(♣1) Nu este scopul prezentei lucrari sa arate cum sunt obtinute reprezentarile FX-bar ale structurilor LN (într-o maniera mai mult sau mai putin algoritmica), ci doar sa propuna schema FX-bar generala ca un mecanism esential de reprezentare a informatiei lingvistice, sa sugereze cum lucreaza, si sa explice ratiunile introducerii acestui mecanism. Teoria FX-bar este integrata ca o componenta importanta a strategiei lingvistice SCD, însa ea poate utilizata si în alte contexte computationale, cu conditia de a include ingredientele necesare, si anume, clasele de PS-*Ms*, ierarhiile acestor clase, o taxonomie functionala (predicationala) si relationala a categoriilor majore si a markerilor, un algoritm (în particular, algoritmi SCD) de obtinere a structurilor de dependenta, etc. Aspecte mai detaliate ale SCD au fost prezentate în [1], [2], [6], [7]. (♣2) Functionarea corecta a schemelor FX-bar expuse arata clar cât de necesara este utilizarea (intensiva) a trasaturilor *predicative* si *functional-relationale* pentru fiecare categorie lexicala. Din experienta noastra în ce priveste analiza si generarea automata a limbii române [6], consideram ca accentul pus pe trasaturile functionale ale categoriilor gramaticale, cuplat cu punerea în evidenta a PS-*Ms*, reprezinta elemente-cheie în utilizarea cu succes a teoriilor X-bar curente în procesarea automata a LN si în cadrul unor teorii lingvistice moderne (UG, FG,

HPSG, etc.). (♣3) Punerea în valoare a trasaturilor *functionale* (în particular, predicationale) ale categoriilor majore N, V, A, și a celor *relationale* ale claselor de *marcheri* (marcheri numiți în literatura și “*cue phrases*” [Rom: *sintagme indicatoare*] [28], [31], sau *conective* [29], [30], etc., deși esențiale, nu poate rezolva toate problemele. De exemplu, asignarea dependentelor corecte în juxtapunerea de NGs este o problemă binecunoscut de dificilă, imposibil de rezolvat complet doar la nivel sintactic. Există însă în prezent un puternic curent către acest tip de abordări, aceasta deoarece ele reflectă mult mai adecvat structura reală a textului de LN (cel puțin pentru o clasă largă de LNs europene). Aceste abordări pot diferi substanțial în instrumentele și tehnicile de parsare, însă principiile rămân foarte similare (de exemplu, [19], [29], [31], [33], etc. (♣4) PS-Ms (marcherii de structuri sintagmatice) joacă un rol fundamental în delimitarea structurilor sintactice și semantice, și stabilirea dependentelor corecte între aceste structuri. SCD a pus accentul încă de la începuturi pe acest aspect [22]. Se remarcă în prezent o întreagă mișcare către reconsiderarea rolului esențial al marcherilor, în special la nivel de *discurs* și în analize complexe ale marilor unități textuale (regăsirea informației, rezumare automată, planificare și generare automată de text, etc.). Strategia SCD, cu componenta ei de teorie FX-bar, încearcă să pună la lucru întreaga paletă de PS-Ms, de la nivel *lexical* și de *coeziune* (locală), până la nivel de *discurs* (*coeziune* și *coerentă* globală), punând accentul pe sintaxă (nivelul de “*suprafata*”, [Eng: *shallow*]) și pe un nivel minimal de semantism. În funcție de problema de LN ce trebuie rezolvată, aceste niveluri pot fi amplificate în mod corespunzător. (♣5) Cuplarea schemelor FX-bar cu: (a) clasele de marcheri SCD și cu ierarhia lor ce corespunde celor patru nivele de proiecție lingvistică din FX-bar [7]; (b) o taxonomie bazată pe predicționalitate a categoriilor majore N, V, A; (c) exploatarea maximală a trasaturilor *functionale* (predicazionale) și *relationale* a tuturor categoriilor lexicale și nelexicale (deci a PS-Ms); (d) o schemă X-bar simplă și unică, apelată recursiv pe cele patru nivele ale sale, pornind de la *lexicon* (conventional, BAR = -1) și până la nivelul de *discurs* al frazei multi-eveniment (BAR = 3), aceste aspecte reprezintă principalele diferențe (și noutăți) dintre teoria FX-bar și teoriile X-bar precedente. (♣6) Schema FX-bar poate fi de asemenea asociată cu un automat dependent de limbaj (pentru o largă clasă de LNs), care începe să lucreze pentru fiecare frază, primește *on-line* cuvânt cu cuvânt, și se oprește odată cu semnul de punctuație final al frazei. Pentru valori adecvate ale parametrilor de LN cum sunt *ordinea cuvintelor* (*argumentelor*) și *direcția proiecție lingvistice* pentru categoriile majore și pentru marcheri, schema FX-bar poate reprezenta corect dependențele structurilor lingvistice (inclusiv pentru Ex.3.2.2.G).

5. Problema X-bar teoriei actuale

Mai este necesară X-bar teoria sau nu? Este teoria X-bar pe moarte sau nu? Care este valoarea teoretică și, mai ales, practică a (sub)teoriei X-bar în teoriile lingvistice și a tehnologiilor actuale ale LN? Cum trebuie să percepem în mod corect X-bar teoria atunci când, în aceeași carte a lui Chomsky, găsim următoarele două pasaje:

(Chomsky1): *“The concepts of X-bar theory are therefore fundamental. In a minimalist theory, the crucial properties and relations will be stated in the simple and elementary terms of X-bar theory.”* [3, p. 172],

(Chomsky2): *“Standard X-bar theory is thus largely eliminated in favor of bare essentials.”* [3; p. 246].

Subliniem ca aceste citate nu sunt extrase din text astfel încât sa nu aiba relevanta în context, intentia de a provoca confuzie. Dimpotriva! De asemenea, scopul nostru nu este de a cauta o posibila incoerenta ci de a pune în evidenta noua pozitie a lui Noam Chomsky, între 1992 si 1995. Încercam sa deschidem o discutie pe aceasta tema deoarece consideram ca exista o problema, si ca ea este de o reala importanta.

În aceasta sectiune urmarim patru obiective: **(A)** Sa enuntam problema X-bar teoriei. **(B)** Sa rezumam solutiile existente în momentul de fata. **(C)** Sa stabilim rolul X-bar teoriei în interiorul contextului teoriilor lingvistice si sa sugeram posibile dezvoltari. **(D)** Sa specificam pozitia FX-bar schemelor propuse privitor la dilema eliminarii complete a X-bar teoriei si, în special, relatia noii FX-bar teorii conturate în contextul strategiei lingvistice SCD. **(E)** Cateva concluzii si perspective.

(A) Sa consideram urmatoarea *problema*: reflecta teoria X-bar o realitate lingvistica a LNs, si daca da, prin ce mijloace aceasta realitate lingvistica ar putea fi cel mai bine reflectata? Proiectia categoriilor lingvistice este un fapt lingvistic de netagaduit. Chomsky si alti distinsi lingvisti nu au fost în completa eroare în ultimii 25-30 de ani? Credem ca nu. Problema este daca teoria X-bar poate înca sa mai fie un *bun model*, sau vehicul, care sa exprime acest *fapt*, si cu ce pret de utilitate. *Principiul Proiectiei Extinse* [3, p. 55] si *Principiul Proiectiei Maximale* (propus în [19] si sectiunea 3.1.) au ca scop sa stabileasca forma si marginile cele mai probabile ale unitatilor textuale obtinute în cadrul procesului de proiectie a categoriilor lingvistice.

(B) Ipoteza (Chomsky1) de mai sus da un raspuns afirmativ la aceasta întrebare în timp ce (Chomsky2) reprezinta, aparent, opusul acestui raspuns. Abordarea din [3, Cap. *Categories and Transformations*] pentru ipoteza (Chomsky2) este ca disolutia schemelor X-bar, deci a proiectiei categoriilor lingvistice, poate fi înlocuita cu succes prin folosirea proprietatilor de functionalitate, predicativitate, tipologie si transformare intrinseci acestor categorii, desi aceste proprietati sunt reprezentate în [3] cu acelasi aparat X-bar pe care îl combat! În cadrul unei teorii a "*structurii sintagmatice pure*", operatiile unui sistem computational al NL "*construiesc recursiv obiecte sintactice*", iar "*categoriile sunt constructii elementare rezultate din proprietatile elementelor lexicale*", cu conditia "*sa nu fie adaugate obiecte noi în cursul procesarii, înafara de rearanjari ale proprietatilor lexicale*" [3]. Rezultatul pare sa fie spectacular: dispar nivelele de proiectie (în sensul teoriei X-bar), astfel spus, nu se face nici o deosebire între elementele lexicale si nucleele proiectate din ele, în timp ce "*teoria structurilor sintagmatice poate fi eliminata în întregime, se pare, pe baza celor mai elementare ipoteze*" [3, p. 294].

Nu ar fi pentru întâia oară când teoria lingvistică încearcă să renunțe la (sub)teoria X-bar. Chomsky sugerează ca nivelele de proiecție lingvistică pot fi înlocuite de către "*proprietățile (funcționale n.n.) ale elementelor lexicale*". Acesta este chiar cazul *gramaticii funcționale* (FG) [25] în care, formal, lipsește teoria X-bar. Dar chiar și în gramatica funcțională a lui Dik, conținutul ascuns al teoriei X-bar este scufundat de fapt în cele *patru nivele* de structuri ierarhice ale *functorilor* și *operatorilor* ce se aplică pe categoriile și structurile cu care FG lucrează la fiecare nivel sintactic. O situație specială avem în SCD, unde nivelele de proiecție a categoriilor lingvistice sunt recuperate pe baza unei funcționalități ierarhice a elementelor lexicale, iar FX-bar schema propusă poate fi utilizată (recursiv) ca un invariant sintactic constructiv al structurilor sintagmatice în cadrul proceselor de analiză și generare automată a LN (limbii române).

Schema FX-bar propusă (Fig. 3.2.1.) poate fi considerată ca un compromis, o negociere, între (Chomsky1) și (Chomsky2), deoarece (Chomsky2) se prezintă fără mecanisme concrete pentru a-și susține ipoteza: în timp ce teoriile X-bar clasice nu mai pot fi utilizate ca instrumente operationale pentru a reflecta o viziune exclusiv funcțională (și relațională) asupra sintaxei, teoria FX-bar propusă poate face acest lucru.

(C) Poziția noastră privind problema (A) asupra teoriei X-bar poate fi rezumată astfel: (C1) Proiecția categoriilor gramaticale este un fapt lingvistic. (C2) Acest fapt poate fi corect reflectat prin "*nuclee*" și "*nivele (bar) de proiecție*" în interiorul schemelor X-bar, dar și prin proprietățile funcționale "*intrinseci*" ale categoriilor lexicale și gramaticale. (C3) Teoria X-bar include deci o componentă de adevărată construcție lingvistică, iar ingredientul sau de bază este confecționat din relațiile funcționale stabilite între elementele lexicale (și nelexicale) conținute în cadrul schemelor X-bar. (C4) Atunci când proprietățile funcționale ale categoriilor lexicale nu sunt evaluate și exploatate corespunzător, teoria X-bar este inconsistentă și produce dificultăți de calcul și rezultate incorecte. (C5) Acestea sunt consecințele unui aspect mult mai general, și anume că teoria X-bar nu trebuie să fie văzută ca o teorie gramaticală singulară, construită pentru sine, ci ca un dispozitiv component al unui *mecanism lingvistic teoretic și computațional mai general*, ale cărui principii să guverneze teoria X-bar. Axiomatica (bazele constructive ale) teoriei X-bar trebuie să fie un *rezultat* al bunei ei funcționări, pe fenomenele concrete de limbaj, și nu invers! (C6) *Ad limitum*, se poate concepe ca mecanismul lingvistic teoretic menționat mai înainte poate funcționa și fără includerea dispozitivului reprezentat de teoria X-bar, așa cum încearcă teoria MinP să propună în [3, Cap. *Categories and Transformations*] (dar folosindu-se în explicare tot de aparatul de reprezentare al teoriei X-bar), precum și în cazul FG [25].

(D) Considerăm ca schemele (funcționale) FX-bar propuse furnizează un (sub)sistem necesar și folositor în cadrul oricărei teorii sintactice asupra LN, inclusiv (și în special) pentru strategia lingvistică SCD. O condiție esențială pentru schemele FX-bar este că ele să reflecte corespunzător proprietățile *funcționale* și *relaționale* ale categoriilor tuturor *lexicale* și *gramaticale*. Exemplele 4.1.1.-4.1.5. arată cum sunt construite schemele FX-bar, cum se obțin (prin apel recursiv pe nivele) structurile sintagmatice complexe ale

LN, și cum acestea rămân închise la operatorul de compunere (adjuncție) pe baza principiilor și regulilor SCD.

Schimbând perspectiva, prin definirea teoriei FX-bar ca o componentă a strategiei lingvistice SCD, și parafrazând formalismul bine-cunoscut al gramaticilor TAG [Eng: *tree adjoining grammar*], strategia SCD poate fi văzută și ca o teorie a evaluării și adjuncției de FX-bar scheme. Este doar o mostră a rolului important pe care teoria X-bar îl poate juca în cadrul teoriei și tehnologiei LN.

(E) Un *element original* propus de schemele FX-bar în peisajul teoriilor X-bar cunoscute este rolul lor dublu ce îl pot juca în cadrul strategiei SCD (și nu numai): Schemele FX-bar pentru $X = N, V, A, CL$ ($CL =$ clauză) trebuie concepute ca un set de invariante sintactice (dinamici) ce pot fi folosiți (1) la *reprezentarea* informației lingvistice la nivel de lexicon (în mod similar cu structurile de trasaturi lingvistice [18], dar într-o manieră mai simplă și mai regulată), și (2) la *procesarea (analizarea și generarea) automată* de text în LN (inclusiv, și mai ales, limba română), de la structurile sintagmatice simple până la cele de discurs.

Derivarea de *automate* și *gramatici formale* bazate pe schema FX-bar, pentru analiza LN, ar fi o consecință normală și o provocare a prezentei propuneri. Modul *recursiv, ascendent și incremental* (prin apelul de funcții și relații cu rol lingvistic multiplu), dar și *descendent* (bazat pe sateliții nucleelor semantice), utilizarea la maximum a *contextualității* marcherilor de toate tipurile poate reprezenta o motivație naturală pentru cercetarea relației dintre strategia SCD (cu componenta ei de *teorie* FX-bar), și modelele generative generoase oferite de către *gramaticile contextuale* Marcus [41], [42], un formalism *context-dependent* puternic, destinat reprezentării, parsării, dar și analizei semantice și de discurs (articularea *topic-focus* [37]) a LN. *Gramaticile contextuale* Marcus aparțin unei serii de formalisme care includ *gramatici* TAG [43], *gramatici orientate-nucleu* [15], [16], *gramatici indexate*, *gramatici X-bar*, *gramatici context-free marcate* [44] etc., formalisme ce realizează o modelare mai realistă a comportamentului sintactic, semantic și discursiv al LN.

Referințe bibliografice

- [1] N. Curteanu (1990). *A Marker-Hierarchy-based Approach Supporting the SCD Parsing Strategy*. Research Report no. 18, Institute of Technical Cybernetics, Bratislava.
- [2] N. Curteanu (1994). *From Morphology to Discourse Through Marker Structures in the SCD Parsing Strategy. A Marker-Hierarchy Based Approach*. Language and Cybernetics, Akademia Libroservo, Prague, 61-73.
- [3] Noam Chomsky (1995). *The Minimalist Program*. The MIT Press, Cambridge, Massachusetts.
- [4] N. Chomsky (1970). *Remarks on Nominalizations*. In R. Jacobs and P. Rosenbaum (eds.), *Readings in Transformational Grammar*, Ginn and Co., Boston, 184-221.

-
- [5] T. Stowell (1981). *Origins of Phrase Structure*. Ph.D. Dissertation, Dept. of Linguistics and Philosophy, MIT, Cambridge.
- [6] N. Curteanu, G. Holban (1996). *Strategia lingvistica SCD aplicata la analiza si generarea limbii române*. Limbaj si Tehnologie (Dan Tufis, Ed.), Academia Româna, Bucuresti, p. 169-176.
- [7] N. Curteanu, C. Lintes (2002). *Segmentation Algorithms for Clause-Type Textual Units*, Research Report, Institute of Theoretical Informatics, Romanian Academy.
- [8] Noam Chomsky (1986). *Barriers*. The MIT Press, Cambridge.
- [9] Noam Chomsky (1981). *Lectures on Government and Binding*. Foris, Dordrecht.
- [10] Simon C. Dik (1989). *The Theory of Functional Grammar*. Foris Publishers, Dordrecht.
- [11] Carl Pollard, Ivan Sag (1994). *Head-Driven Phrase Structure Grammar*. The University of Chicago Press, Chicago & London.
- [12] Gerald Gazdar, E. Klein, G. Pullum, I. Sag (1985). *Generalized Phrase Structure Grammar*. Harvard University Press, Cambridge, Massachusetts.
- [13] Peter Sells (1985). *Lectures on Contemporary Syntactic Theories*. CSLI, Stanford, California.
- [14] Stuart Shieber (1986). *An Introduction to Unification-Based Approaches to Grammar*. CSLI, Stanford, California.
- [15] Carl Pollard, Ivan Sag (1987). *Information-based Syntax and Semantics*. CSLI, Stanford, California.
- [16] Carl Pollard, Ivan Sag (1994). *Head-Driven Phrase Structure Grammar*. The University of Chicago Press, Chicago & London.
- [17] E.P. Stabler Jr. (1992). *The Logical Approach to Syntax: Foundations, Specifications and Implementations of Theories of Government and Binding*. The MIT Press, Cambridge, Massachusetts.
- [18] N. Curteanu, G. Holban (2000). *A Set-Theoretic Approach to Linguistic Feature Structures and Unification Algorithms* (I, II). Computer Science Journal of Moldova, 8(2): 116-149, 8(3): 223-246.
- [19] Neculai Curteanu (1988). *Augmented X-bar Schemes*. COLING'88 Proceedings, Budapest, 130-132.
- [20] Neculai Curteanu, A. Todirascu, G. Holban (1997). *Teorii sintactice ale limbajului natural*. Raport de cercetare, Institutul de Informatica Teoretica, Academia Româna, Iasi, 66 p.
- [21] Alain Lecomte (1998). *Multimodal Logic for Syntax*. Logica Trianguli, 2: 49-72.
- [22] Neculai Curteanu (1983). *Algoritmi de analiza sintactica a frazei si propozitiei*

-
- românești. INFO-IASI'83*, p. 553-548.
- [23] M. Moortgat (1997). *Categorial Type Logics*. Handbook of Logic and Language, Elsevier.
- [24] E.P. Stabler Jr. (1997). *Derivational Minimalism*. Logical Aspects of Computational Linguistics, LNCS no.1328, Springer-Verlag, Berlin.
- [25] Simon Dik (1989). *The Theory of Functional Grammar*. Foris Publishers, Dordrecht.
- [26] Robert Kasper (1993). *Adjuncts in the Mittelfeld*. In "German Grammar in HPSG" (J. Nerbonne *et al.*, Eds.), CSLI, Stanford, California.
- [27] Denis Bouchard (1995). *The Semantics of Syntax. A Minimalist Approach to Grammar*. The Univ. of Chicago Press, Chicago & London.
- [28] Julia Hirschberg, D. Litman (1993). *Empirical Studies on the Disambiguation of Cue Phrases*. Computational Linguistics 19(3): 501-530.
- [29] Jacques Jayez, C. Rossari (1999). *Pragmatic Connectives as Predicates. The Case of Inferential Connectives*. In "Predicative Forms in Natural Language and in Lexical Knowledge Bases" (P. Saint-Dizier, Ed.), Kluwer Academic Publishers, Dordrecht.
- [30] Patrick Saint-Dizier (Ed.) (1999). *Predicative Forms in Natural Language and in Lexical Knowledge Bases*. Kluwer Academic Publishers, Dordrecht.
- [31] Daniel Marcu (2000). *The Theory and Practice of Discourse Parsing and Summarization*. The MIT Press, Cambridge.
- [32] V. Raskin, S. Nirenburg (1999). *Lexical Rules for Deverbal Adjectives*. In "Breadth and Depth of Semantic Lexicons", Kluwer Academic Publishers, Dordrecht.
- [33] O. Popârda, N. Curteanu (2002). *L'évolution du discours juridique français analysé par la stratégie linguistique SCD*. In "Représentation du Sens Linguistique" (D. Bouchard, Ed.), ELCOM Studies in Theoretical Linguistics, ELCOM EUROPA.
- [34] Noam Chomsky (2000). *Minimalist inquiries: the framework*. In R. Martin *et al.* (Eds) "Step by step. Essays on Minimalist Syntax in Honor of Howard Lasnik", MIT Press, Cambridge, p. 89-155.
- [35] Noam Chomsky (2001). *Derivation by phase*. In M. Kenstowicz (Ed.) "Ken Hale: a life in language", MIT Press, Cambridge, p. 1-52.
- [36] Jane Morris, G. Hirst (1991). *Lexical Cohesion Computed by Thesaural Relations as an Indicator of the Structure of Text*. Computational Linguistics 17(1): 21-48.
- [37] Eva Hajicova, H. Skoumalova, P. Sgall (1995). *An Automatic Procedure for Topic-Focus Identification*. Computational Linguistics, 21(1): 81-94.
- [38] P. Sgall, E. Hajicova, J. Panevova (1986). *The Meaning of the Sentence in Its Semantic and Pragmatic Aspects*. Kluwer Academic Publishers, Dordrecht.
- [39] S. Shieber, F. Pereira, G. Van Noord, R. Moore (1990). *Semantic Head-Driven Generation*. Computational Linguistics 16(1): 30-41.

- [40] Steven Abney (1996). *Part-Of-Speech Tagging and Partial Parsing*. In “Corpus-Based Methods in Language and Speech”, (K. Church *et al.*, Eds.), Kluwer Acad. Publishers, Dordrecht.
- [41] Solomon Marcus (1997). *Contextual Grammars and Natural Language*. In Cap. 5 (Vol. 2) din “The Handbook of Formal Languages”, G. Rozenberg, A. Salomaa, Eds., Springer-Verlag, Berlin, 215-235.
- [42] Gheorghe Paun (1997). *Marcus Contextual Grammars*. Kluwer Academic Publishers, Dordrecht.
- [43] Michele Abrusci, Christophe Fouqueré, Jacqueline Vauzeille (1999). *Tree Adjoining Grammars in a Fragment of the Lambek Calculus*. *Computational Linguistics*, 25(2): 209-236.
- [44] Philip Miller (1999). *Strong Generative Capacity. The Semantics of Linguistic Formalism*. CSLI Publications, Stanford, California.

Teoria HPSG. Studiu de caz: acordul încrucișat

Ana-Maria BARBU
RACAI, Calea 13 Septembrie nr.13, Bucuresti
abarbu@racai.ro

1. Introducere

Oricât ar fi de mare entuziasmul creat de performanțele realizate cu calculatorul, care cuprinde deopotrivă și domeniul prelucrării limbajului natural, rezultate temeinice nu se pot obține dacă acestea nu sunt fundamentate pe îndelungi și profunde analize teoretice. Nu putem aspira la obiective majore în ingineria lingvistică, precum analizarea și generarea de texte, construirea de verificatoare ortografice și gramaticale sau chiar de traducătoare automate, dacă se ignoră particularitățile inerente ale obiectului în studiu, anume ale limbajului natural în general, și a limbii de aplicație, în special. Or aceste particularități sunt oferite, sub un aspect sau altul, tocmai de teoriile gramaticale. Experiența a dovedit că eșecurile din ingineria lingvistică au avut ca posibile surse eșecurile în descrierea corespunzătoare a fenomenelor de limbă, dar și succesele, la rândul lor, s-au datorat în parte acurateții, exactității, și nu în ultimul rând caracteristicilor computaționale ale unui model gramatical teoretic.

Iată de ce alegerea unei teorii lingvistice adecvate, cu scopul de a scrie pe baza acesteia o gramatică computațională a unei limbi particulare, în speța a limbii române, este un act de primă însemnătate.

După anii primelor dezvoltări ale gramaticii generative, sintaxa formală este, de aproape două decenii, repusă în discuție ca obiect de studiu autonom distinct în același timp de cel al lexicului și cel al sensului. Mai multe curente teoretice, cunoscute sub numele generice de „gramatici de unificare” sau „gramatici bazate pe constrângeri”, s-au născut din această reconsiderare a sintaxei. Este vorba de modele recente (cele mai vechi datând de la începutul anilor ‘80), dezvoltate în cea mai mare parte în Statele Unite, și în general aproape necunoscute publicului român. Aceste modele se pretează scrierii de gramatici pentru calculator, dar ambiția lor este mai întâi de a constitui teorii lingvistice de sine statatoare. Autorii lor se înscriu pe linia programului gramaticii generative chomskyene din 1957, de la care preiau grija pentru o formalizare operatorie a sintaxei, dar se disting suficient de modelul actual al Școlii de la Cambridge (numit Government and Binding) pentru a prezenta teorii alternative. Printre punctele comune ale gramaticilor de unificare, se află pe de o parte atenția acordată unei articulare mai explicite a lexicului, sintaxei și semanticii, pe de altă parte accentul pus pe descrierile lingvistice și recurgerea la un stil de

analiza sintactica mai "concret", care limiteaza recurgerea la elemente "vide" (nerealizate concret) si care restrânge numarul etapelor intermediare în producerea unui enunt.

În acest articol vom prezenta pe scurt una dintre teoriile lingvistice amintite, anume „Gramatica sintagmatica ghidata de centru”, denumita abreviat HPSG dupa numele sau din engleza „Head-driven Phrase Structure Grammar”. Apoi vom ilustra modul în care poate fi aplicata aceasta teorie în reprezentarea unui fenomen mai special de limba româna prin aceea ca presupune dependente încrucisate de acord. Este vorba de structuri relative de tipul *baiatul a carui sora cânta* unde articolul genitival *a* se acorda cu substantivul *sora*, iar pronumele relativ *carui* se acorda cu substantivul *baiatul*.

2. Teoria lingvistică HPSG

2.1. Scurt istoric

Modelul gramaticii sintagmatice ghidate de centru (engl. *Head-driven Phrase Structure Grammar*, sau HPSG) a fost conceput la începutul anilor '80 de Carl Pollard si Ivan Sag cu scopul de a permite o integrare mai explicita a diferitelor nivele de analiza lingvistica: fonetic, sintactic si semantic. El a luat nastere în principal din Gramatica Sintagmatica Generalizata (GPSG) si din lucrarile lui C. Pollard despre *Head Grammar* [1], dar autorii lor s-au inspirat deopotriva din numeroase alte teorii. Ei au preluat de la modelul chomskyian al Guvernarii si Anaforicitatii (GB) notiunea de modularitate si recurgerea la principii foarte generale (Principiul anaforicitatii, al controlului etc.). De la gramatica functionala de unificare FUG [2] au împrumutat reprezentarea uniforma a elementelor lexicale, a sintagmelor si regulilor gramaticale sub forma de structuri de trasaturi. S-au inspirat de la gramatica lexical functionala LFG pentru îmbogătirea cadrelor de subcategorizare si a notiunii de regula lexicala. Au luat de la gramaticile categoriale ideea de saturare progresiva a predicatelor si recurgerea la o ierarhie de functii gramaticale (cf. [3]). S-au inspirat, în sfârșit, dintr-un punct de vedere mai formal, din lucrari de logica si informatica asupra tipurilor si mostenirii.

Teoria este prezentata în cele doua lucrari ale lui C. Pollard si Ivan Sag: [4] si [5]. Majoritatea exemplelor privesc limba engleza si trateaza fenomene variate: fenomene de acord, constructii infinitivale, anafore, constructii relative si comparative. Fenomenele de control sunt totodata dezvoltate în [6], iar o analiza a anaforelor este propusa în [2]. Primele lucrari au conferit de la bun început o dimensiune multilinguala acestei teorii prin abordari privind germana ([8], [9]), catalana ([10]), japoneza ([11]), dar si coreana ([12]), franceza ([13]) si italiana ([14]).

C. Pollard si I. Sag preiau din modelul GPSG notiunea de gramatica sintagmatica, cu distinctia între o componenta ierarhica (schema DI –de dominanta imediata) si o componenta liniara (principii de precedenta liniara), precum si recurgerea la principii foarte generale de partaj si de propagare a trasaturilor. Totusi ei se separa de modelul original în câteva puncte. Structurile sintagmatice sunt în întregime exprimate în termeni de de structuri de trasaturi, cu introducerea unui atribut Ramuri. Structurile de trasaturi sunt la

rândul lor organizate în ierarhii de tipuri, comportând fiecare trasaturi predefinite. Modelul HPSG ofera astfel anumite simplificari în raport cu GPSG: întregul arsenal de reguli DI este redus la șase scheme de baza; metaregulile sunt eliminate în favoarea regulilor lexicale. S-a urmarit deosebirea clara între ceea ce tine de domeniul constrângerilor universale și ceea ce tine de descrierea unei limbi particulare. Principiile de coocurența a trasaturilor din GPSG, care amesteca constrângerile universale și cele specifice unei limbi date, au fost suprimate.

2.2 Organizarea generala a HPSG

2.2.1 Caracteristici specifice gramaticilor de unificare

Se poate considera ca gramaticile de unificare, sau gramticile bazate pe constrângeri, reprezinta noile teorii sintactice ale anilor '80. Este vorba de modele care urmaresc o articulare explicita între lexic, sintaxa și semantica. Proprietatile lingvistice corespunzatoare sunt concepute ca "informatii" asociate morfemelor, sintagmelor sau constructiilor, combinate prin operatii variate, dintre care unificarea ocupa un rol central. Aceasta concepie "integratoare" este unul dintre atuurile lor pentru tratarea automata a limbajelor naturale. Un alt avantaj este ca ele se bazeaza pe modele logice sau matematice (gramatici de constituenți, structuri de trasaturi), pentru care au fost definite metode de programare. Ele sunt în general rezultatul unui compromis între expresivitatea lingvistica (grija de a facilita exprimarea diferitor principii lingvistice adaugând u-se variante notationale sau operatori) și eficacitate (notatii concentrate, putine operatii).

Aici, ne vom rezuma sa punctam trasaturile lor comune cele mai pregnante, dintre care:

- reabilitarea descrierilor de suprafata;
- reînnoirea descrierilor sinatctice prin definirea de trasaturi complexe;
- definirea de principii generale de buna formare a enunturilor;
- integrarea lexicului, sintaxei și semanticii.

Gramaticile de unificare îmbogatesc aparatul formal al gramaticilor de constituenți cu un numar de notiuni importante. În acest capitol ne vom limita la prezentarea principalelor notiuni utilizate pe parcursul lucrării, pentru detalii putând fi consultate S. Shieber 1986a sau H. Uszkoreit 1989.

2.2.1.1 Structuri de trasaturi

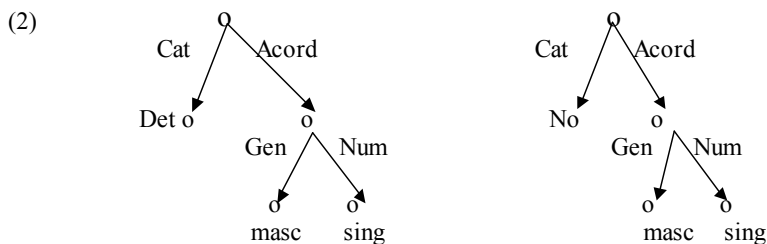
Structurile de trasaturi (engl. *feature structure*) sunt primitive ale teoriilor sintactice bazate pe unificare și reprezinta ansambluri de trasaturi, numite și complexe de trasaturi (engl. *feature complexes* sau *feature bundles*), care pot fi reprezentate sub forma de matrice. O **trasatura** este o pereche atribut-valoare, valorile putând fi simboluri atomice sau trasaturi. Trasaturile cu valoare non atomica conduc la structuri de trasaturi care prezinta îmbricari.

Spre exemplu, cuvintelor *acest* si *câine* li se asociaza o trasatura Cat cu valoare atomica (pentru categorie) si o trasatura complexa Acord care ia ca valoare conjunctia a doua trasaturi Num (pentru numar) si Gen:

$$(1) \quad \begin{array}{cc} \text{acest} & \text{câine} \\ \left[\begin{array}{l} \text{Cat} = \text{Det} \\ \text{Acord} = \left[\begin{array}{l} \text{Gen} = \text{masc} \\ \text{Num} = \text{sing} \end{array} \right] \end{array} \right] & \left[\begin{array}{l} \text{Cat} = \text{N} \\ \text{Acord} = \left[\begin{array}{l} \text{Gen} = \text{masc} \\ \text{Num} = \text{sing} \end{array} \right] \end{array} \right] \end{array}$$

O structura este rau formata când contine de doua ori acelasi atribut (la acelasi nivel de îmbricare) cu o valoare diferita.

Si alte reprezentari de structuri de trasaturi (sau structuri atribut-valoare) sunt posibile, fiind echivalente formal. Cele mai utile, pentru implementarea informatica, sunt cele care utilizeaza grafuri orientate: arcuri care poarta nume de trasaturi si puncteaza spre noduri care sunt etichetate cu valoarea trasaturii (daca e vorba de trasaturi cu valoare atomica) sau sunt puncte de plecare pentru alte arce (pentru trasaturi cu valoare non atomica). De pilda, pentru exemplele de mai sus vom avea urmatoarele reprezentari:



În termeni de grafuri, echivalentul interdictiei ca un acelasi atribut sa apara de doua ori la acelasi nivel cu valori diferite este interdictia ca doua arcuri sa puncteze, plecând din acelasi nod, catre doua noduri diferite care poarta aceeasi eticheta (ceea ce e o restrictie generala asupra grafurilor ce corespund automatelor deterministe).

Structurile de grafuri pot fi ciclice sau non ciclice. Acestea din urma se numesc **grafuri aciclice orientate** (engl. *Directed Acyclic Graph* sau DAG), denumire adesea folosita pentru a desemna structurile de trasaturi.

În lucrul cu structuri de trasaturi complexe se impun unele distinctii, de pilda, între structurile identice si **structurile cu valori partajate** (sau **reentrante**). Cele din urma sunt identice si vor ramâne astfel indiferent de modificarile suferite ulterior, ceea ce nu se întâmpla cu primele. În exemplul ce urmeaza structura de trasaturi A comporta doua

atribute cu valori identice Acord si Num. În structura B, cele doua atribute Acord sunt coindexate (prin indicele 1), ceea ce face ca ele sa partajeze în mod egal trasatura [Num = sing].

$$(3) \quad \begin{array}{cc} \text{A:} & \text{B:} \\ \left[\begin{array}{l} \text{Det} = [\text{Acord} = [\text{Num} = \text{sing}]] \\ \text{Nume} = [\text{Acord} = [\text{Num} = \text{sing}]] \end{array} \right] & \left[\begin{array}{l} \text{Det} = [\text{Acord} = | 1 | [\text{Num} = \text{sing}]] \\ \text{Nume} = [\text{Acord} = | 1 |] \end{array} \right] \end{array}$$

Daca se unifica fiecare din aceste structuri cu structura C de mai jos, rezultatul nu va fi acelasi:

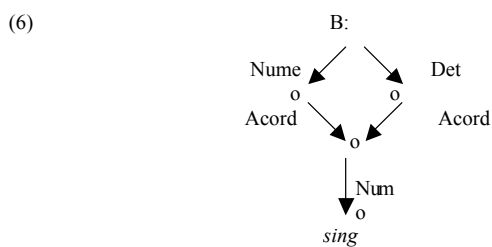
$$C: [\text{Det} = [\text{Acord} = [\text{Gen} = \text{masc}]]]$$

$$(4) \quad \begin{array}{c} C \cup A: \\ \left[\begin{array}{l} \text{Det} = [\text{Acord} = [\text{Num} = \text{sing}, \text{Gen} = \text{masc}]] \\ \text{Nume} = [\text{Acord} = [\text{Num} = \text{sing}]] \end{array} \right] \end{array}$$

$$(5) \quad \begin{array}{c} C \cup B: \\ \left[\begin{array}{l} \text{Det} = [\text{Acord} = | 1 | [\text{Num} = \text{sing}, \text{Gen} = \text{masc}]] \\ \text{Nume} = [\text{Acord} = | 1 |] \end{array} \right] \end{array}$$

Dupa unificare, trasatura Acord îmbricata sub atributul Nume va avea si el o trasatura Gen specificata în cazul lui $C \cup B$, dar nu si în cazul $C \cup A$.

În termeni de grafuri, reprezentarea unei structuri reentrante ca B este urmatoarea:



2.2.1.2 Extensiune și unificare

Se definește o relație de **extensiune** între trasaturi după cum urmează:

O structura de trasaturi A este o **extensiune** a unei structuri de trasaturi B (notându-se $A \supset B$) daca si numai daca:

-- toate trasaturile cu valoare atomica prezente în B sunt prezente si în A cu aceeasi valoare,

-- pentru orice trasatura $\langle f \rangle$ cu valoare non atomica, valoarea lui $\langle f \rangle$ în A este o extensiune a valorii lui $\langle f \rangle$ în B.

De exemplu, structura de trasaturi asociata cuvântului *câine* în (1) este o extensiune a structurii din (7), dreapta, dar reciproca nu este adevarata pentru ca structura de mai jos nu are trasatura [Num = sing] prezenta în cea a cuvântului *câine*:

$$(7) \quad \left[\begin{array}{l} \text{Cat} = \text{N} \\ \text{Acord} = \left[\begin{array}{l} \text{Gen} = \text{masc} \\ \text{Num} = \text{sing} \end{array} \right] \end{array} \right] \supset \left[\begin{array}{l} \text{Cat} = \text{N} \\ \text{Acord} = [\text{Gen} = \text{masc}] \end{array} \right]$$

Daca numarul de atribute nu este limitat se pot obtine o infinitate de structuri care sunt extensii ale unei structuri date. Relatia inversa a extensiei se numeste **subsumare**, A subsuma B daca si numai daca B este o extensie a lui A.

Pe baza acestei relatii de ordine putem defini o structura de stiva, cu o limita superioara si o limita inferioara. Este de notat ca aici nu exista o relatie de ordine stricta pentru ca orice structura este o extensie a ei însesi ($A \supset A$). Structura care le subsumeaza pe toate celelalte (pentru care toate celelalte sunt extensii) este structura vida (notata \mathcal{T}), pe care o putem interpreta ca disjunctia tuturor cuplurilor atribut-valoare ale gramaticii. Daca dorim sa plasam o limita superioara, structura care va fi o extensie a tuturor celorlalte (care este subsumata de toate celelalte) va fi cea care contine conjunctia tuturor cuplurilor atribut-valoare posibile (notata \perp) adica o structura "falsa" sau rau formata.

Aceasta relatie de ordine e folosita pentru a defini unificarea. Aceasta operatie a luat nastere din cercetarile în logica si informatica (limbajul Prolog). Definita la început ca procedura de rezolvare pentru logica predicatelor de ordinul întâi, cf. [15], ea a fost introdusa în lingvistica de A. Colmerauer, [16], apoi de M. Kay, [17], pentru a testa, fuziona si propaga trasaturi sintactice. Ea este definita în felul urmator:

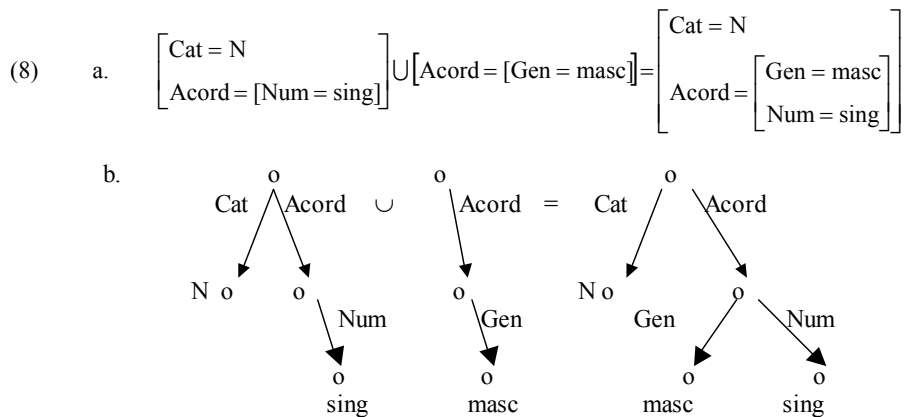
Unificarea a doua structuri de trasaturi A si B (notata $A \cup B$) este structura minimala care este în acelasi timp o extensiune a lui A si a lui B. Daca o astfel de structura nu exista, unificarea "esueaza" (ceea ce e notat cu \perp).

Altfel spus, unificarea verifica compatibilitatea dintre doua structuri de trasaturi si produce o structura rezultanta care este cea mai mica structura ce contine toata informatia din prima structura si toata informatia din a doua structura.

Unificarea este o operatie idempotentă ($A \cup A = A$), comutativă ($A \cup B = B \cup A$) si asociativă ($A \cup (B \cup C) = (A \cup B) \cup C$), spunem de asemenea ca este declarativă (daca $A = A'$ si $B = B'$ atunci $A \cup B = A' \cup B'$) si monotona ($A \cup B \supset A$ si $A \cup B \supset B$; daca A

$\supset B$ atunci $\forall C A \cup C \supset B \cup C$), ceea ce vrea să spună că relațiile de extensiune sunt conservate prin unificare. Colocvial spus, unificarea adaugă informație, fără să o scadă.

În termeni de grafuri, echivalentul operației de unificare este fuziunea definită pentru automatele cu număr finit de stări. Pentru exemplul din (8a) se obține reprezentarea grafică din (8b):



Anumii operatori pot fi adăugați structurilor de trasaturi (cf. L. Karttunen 1984), cei mai utili fiind negația (notată \sim sau \neq pentru trasaturi cu valoare atomică) și disjuncția (notată prin acolade sau semnul $/$). Folosirea negației permite să se renunțe la anumite disjuncții. Există de exemplu echivalența între următoarele două ecuații, dacă considerăm că atributul *Mod* are 8 valori posibile în română (indicativ, conjunctiv, imperativ, prezumtiv, infinitiv, gerunziu, supin, participiu):

$$[\text{Mod} \neq \text{inf}] \Leftrightarrow [\text{Mod} = \text{ind}/\text{conj}/\text{prez}/\text{imp}/\text{ger}/\text{sup}/\text{part}].$$

În secțiunea următoare vom trece la descrierea caracteristicilor specifice ale teoriei HPSG care o fac distinctă de toate celelalte teorii bazate pe unificare. Trebuie spus de la bun început că autorii modelului HPSG au preluat o multime de caracteristici ale teoriilor aparute anterior, inclusiv de la gramatica generativă, tocmai din dorința de a aduna într-un singur formalism tot ce e mai adecvat pentru reprezentarea lingvistică în general. Pentru o paralelă detaliată între HPSG și alte teorii bazate pe constrângeri a se vedea [18].

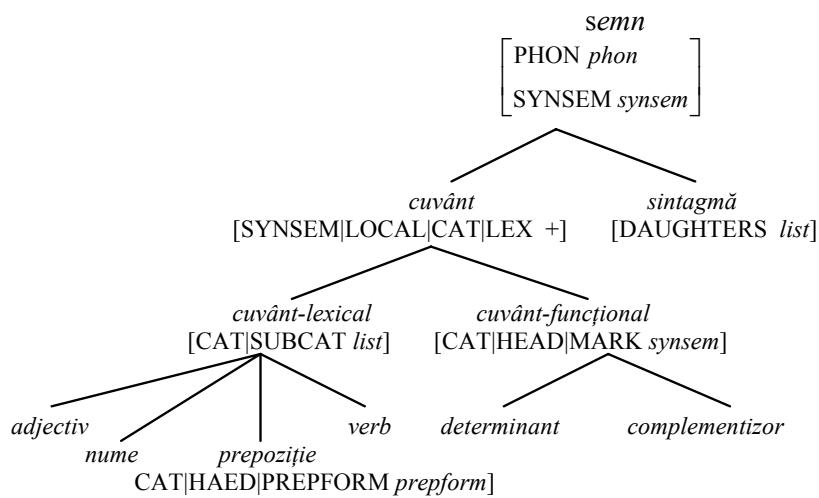
2.2 Caracteristici specifice HPSG

În HPSG, structurile de trasaturi, utilizate în LFG pentru reprezentarea funcțiilor gramaticale, iar în GPSG pentru reprezentarea categoriilor, sunt sistematizate pentru a include atât structurile de constituenți cât și regulile gramaticale. Ele corespund la ceea ce se numește un **semn** lingvistic, adică un cuvânt, o sintagma sau o regulă, conținând

informatii fonetice, sintactice, semantice si discursive. Structurile de trasaturi sunt cât se poate de adecvate pentru organizarea într-o notatie comuna a informatiilor lingvistice eterogene.

Spre deosebire de celelalte teorii lingvistice bazate pe unificare, HPSG utilizeaza ierarhizarea tipologica. Fiecare structura de trasaturi este încadrata într-un anumit tip pentru care sunt predefinite anumite constrângeri si care își are locul într-o ierarhie de tipuri. În cadrul ierarhiei functioneaza relatia de mostenire a constrângerilor tipurilor superioare asupra descendentilor lor. Un exemplu de ierarhie de tipuri este data în (9).

(9)



Pentru fiecare tip sunt definite anumite trasaturi specifice (sau anumite constrângeri) care se adauga constrângerilor mostenite de la tipurile din care descind. Trebuie adaugat ca într-o ierarhie de tipuri sunt permise mosteniri multiple, adica sunt permise tipuri care au mai multi parinti.

Cel mai general tip în HPSG este „semnul” (în engleza *sign*). El contine informatie fonologica (prin trasatura PHON) si informatie sintactico-semnatica (prin trasatura SYNSEM). Semnul, la rândul lui, poate fi un cuvânt sau o sintagma, dupa cum se vede în (9), mai sus. Sintagma are spre deosebire de cuvânt o trasatura în plus, numita DAUGHTERS (adica ramuri-surori) care are ca valoare o lista cu semnele combinate în sintagma. Un exemplu de semn lexical împreuna cu descrierea trasaturilor specifice acestuia este data în (10) pentru verbul *a vrea*.

$$(11) \left[\begin{array}{l} \text{SYNSEM | CAT [1]} \\ \text{DAUGHTERS | HEAD - DTR | SYNSEM | CAT | HEAD [1]} \end{array} \right]$$

Semnul HEAD-DTR poate fi sintagmatic sau lexical.

b. Principiul de Subcategorizare

Atributul SUBCAT are ca valoare o lista care este actualizata progresiv, pe masura ce sintagma se “satureaza”, în sensul ca atunci când complementele sunt realizate, ele sunt eliminate din lista SUBCAT a sintagmei respective. O sintagma se numeste saturata (sau completa) când valoarea listei SUBCAT este vida. Principiul de Subcategorizare poate fi enuntat astfel:

Valoarea listei SUBCAT a ramurii HEAD-DTR a unei sintagme trebuie sa corespunda concatenarii listei L1 ca valoare a atributului SUBCAT al sintagmei si a listei L2 a semnelor ce apartin ramurii de componente COMPS-DTR (sau, mai precis, nu lista semnelor, ci a trasaturilor SYNSEM a acestor semne).

Acesta poate fi reprezentat prin structura de trasaturi urmatoare (notând prin simbolul \oplus concatenarea listelor):

$$(12) \left[\begin{array}{l} \text{SYNSEM | CATEGORY | SUBCAT L1} \\ \text{DAUGHTERS} \left[\begin{array}{l} \text{HEAD - DTR | SYNSEM | CAT | SUBCAT L1} \oplus \text{L2} \\ \text{COMPS - DTR L2} \end{array} \right] \end{array} \right]$$

Tinând seama de Principiul de Subcategorizare pot fi descrise urmatoarele doua scheme DI:

1. Schema DI pentru o sintagma saturata cu ramura Componente: *head-compl* sau *head-subject*

$$(13) \left[\begin{array}{l} \text{SYNSEM | CATEGORY | SUBCAT } \langle \rangle \\ \text{DAUGHTERS} \left[\begin{array}{l} \text{HEAD - DTR | SYNSEM | CAT | SUBCAT } \langle X \rangle \\ \text{COMPS - DTR } \langle X \rangle \end{array} \right] \end{array} \right]$$

2. Schema DI pentru o sintagma non saturata cu ramura Componente: *head-compl*

$$(14) \left[\begin{array}{l} \text{SYNSEM | CATEGORY | SUBCAT } \langle X \rangle \\ \text{DAUGHTERS} \left[\begin{array}{l} \text{HEAD - DTR | SYNSEM | CAT | SUBCAT } \langle X, Y1, Y2...Yn \rangle \\ \text{COMPS - DTR } \langle Y1, Y2...Yn \rangle \end{array} \right] \end{array} \right]$$

3. Schema DI pentru o sintagma cu ramura Adjunct: *head-adjunct*

Modificatorii (adjective atributive, adverbe, complemente circumstantiale) sunt introdusi într-o ramura speciala numita ramura Adjunct (sau ADJCT-DTR). Modificatorii selectioneaza categoria pe care o modifica (N' pentru adjective, V sau GV pentru adverbe). Aceasta selectie se face printr-un atribut MODIF, care are ca valoare o structura de trasaturi SYNSEM. Pentru o sintagma centru-adjunct bine formata trebuie sa aiba loc unificarea valorii trasaturii MODIF a adjunctului cu valoarea trasaturii SYNSEM a centrului. Astfel adjectivele pot selectiona numele pentru care sunt atribuite, iar adverbele pot selectiona verbele respective, adica se poate preciza în intrarea lor lexicala trasaturile Categorie, Continut, Index etc. ale numelui sau verbului asteptat. Descrierea unei sintagme cu Adjunct este urmatoarea:

$$(15) \left[\text{DAUGHTERS} \left[\begin{array}{l} \text{HEAD - DTR} \mid \text{SYNSEM} \mid 1 \mid \\ \text{ADJCT - DTR} \mid \text{SYNSEM} \mid \text{CAT} \mid \text{HEAD} \mid \text{MODIF} \mid 1 \mid \end{array} \right] \right]$$

c. Principiul de Semantic

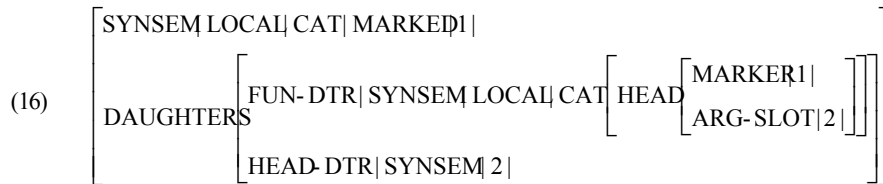
Principiul semantic reglementeaza propagarea trasaturilor semantice, adica cele doua trasaturi CONTENT si CONTEXT. Se urmareste pe de o parte ca sintagmele sa partajeze valoarea trasaturii CONTENT din ramura centrului cu trasatura proprie CONTENT, iar pe de alta parte sa determine "ridicarea" la nivelul sintagmelor superioare a eventualilor cuantificatori si a variabilelor care le pot corespunde.

HPSG face apel la notiunea de centru semantic, acesta fiind identic cu centrul sintactic, în afara cazului sintagmelor cu adjunct. În acest caz, centrul sintactic este categoria modificata, dar centrul semantic este modificatorul (care joaca rolul de predicat semantic). Principiul Semantic poate fi exprimat astfel:

Valoarea atributului CONTENT a categoriei dominante este identica cu valoarea atributului CONTENT a categoriei care este centru semantic (ramura Adjunct sau, implicit, ramura HEAD).

O alta schema DI, *head-functor*, propusa de Allegranza în [19], reprezinta o modificare a schemei *head-adjunct* cu scopul de a satisface exigentele de reprezentare a determinantilor într-un grup nominal. Determinatorii sunt tratati ca functori aplicati centrului. Ei selecteaza centrul prin atributul ARG-SLOT si marcheaza sintagma rezultata cu anumite trasaturi specifice determinantului respectiv prin partajarea valorii atributului MARKER între ramura Functor si nodul mama. Descrierea acestei scheme este data mai jos.

4. Schema DI pentru o sintagma cu ramura Functor: *head-functor*



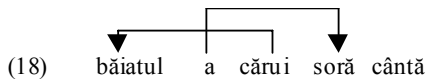
Cu aparatul formal oferit de HPSG, în secțiunea care urmează, dam spre exemplificare analiza unei structuri concrete din limba română. Structura propusă conține un centru nominal modificat de o propoziție relativă al cărei element de relație este în cazul genitiv precedat de articolul genitiv. Aceasta structură este interesantă prin faptul că prezintă un fenomen, acela de acord încrucișat, care pare să scape reprezentărilor gramaticilor independente de context. Avantajul teoriei lingvistice discutate aici, însă, oferă o soluție pe cât de unitară, pe atât de elegantă, după cum sperăm să reiasă din cele ce urmează.

3. Structuri relative cu acord încrucișat

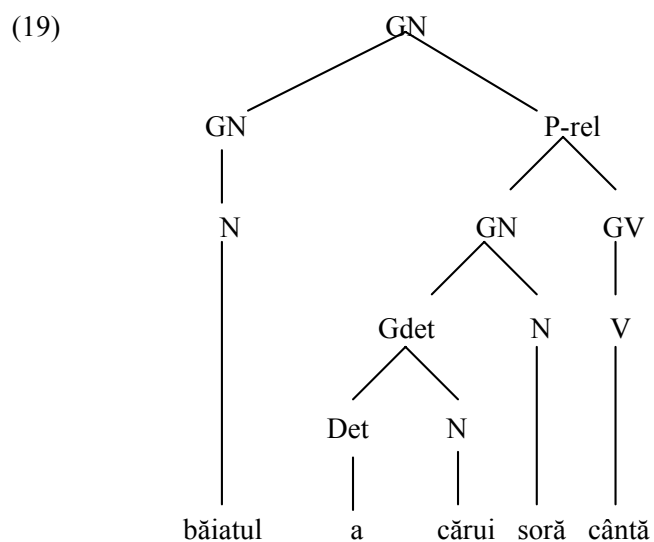
În limba română structurile care prezintă acord încrucișat sunt propozițiile relative în care pronumele relativ este precedat de articolul genitiv, ca în exemplul de mai jos.

(17) baiatul a cărui soră cântă

Acordul este încrucișat prin aceea că pronumele relativ propriu-zis se acordă cu substantivul determinat de propoziția relativă, *baiatul*, iar articolul genitiv *al* se acordă cu subiectul relativei, *sora*, după următoarea schemă:



Structura internă a acestui grup nominal este reprezentată în arborele de mai jos.

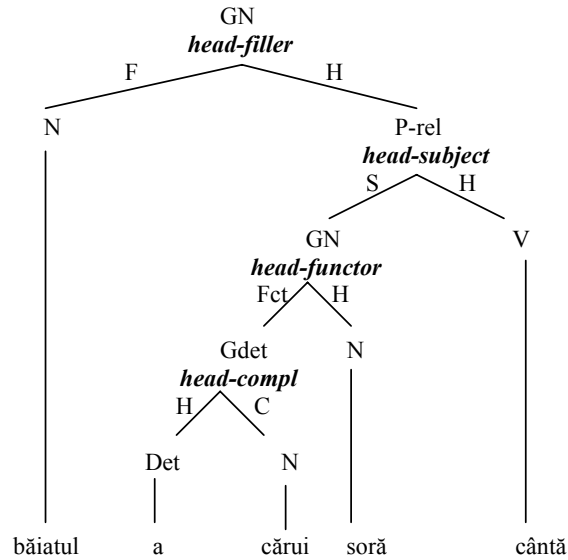


Dupa cum se vede în acest arbore, exemplul din (17) este format dintr-un substantiv centru, *băiatul*, modificat de o propoziție relativă al cărei subiect, *a cărui soră*, cuprinde elementul de relație care face legătura dintre numele amintit și propoziția relativă.

Dacă ne-am limita descrierea la regulile independente de context sugerate în arbore, nu am putea da seama de fenomenul de acord încrucișat pe care-l discutăm aici. Acest lucru este însă posibil dacă folosim o gramatică HPSG, beneficiind de avantajele oferite de mecanismul unificării și de reprezentările prin structuri de trasaturi.

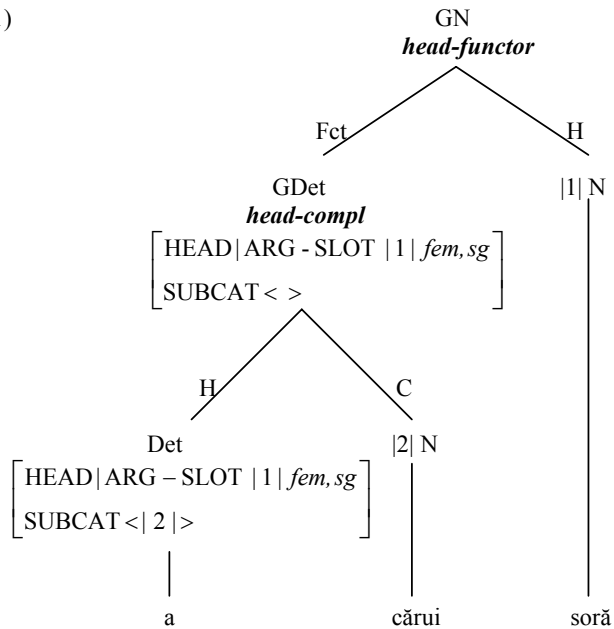
Aplicând schemele de dominanță imediată și principiile specifice teoriei HPSG, arborele de mai sus poate fi adnotat cu regulile HPSG aplicate, în felul următor (unde am folosit ca notații funcționale H=centrul sintagmei, C=complement, Fct=functor, F=filler).

(20)



Fenomenul de acord încrucișat presupune pe de o parte acordul determinatorului *a* cu substantivul *sora*, iar pe de alta parte acordul pronumei relative *cărui* cu substantivul *băiatul*. Primul acord amintit se face relativ banal. Intrarea lexicală a determinatorului *a*, în calitate sa de functor, specifică în valoarea atributului sau central ARG-SLOT ce trasaturi de acord trebuie să aibă substantivul pe care urmează să-l modifice. Când determinatorul *a* se combină cu complementul sau *cărui*, principiul trasaturilor centrale face ca această informație să fie percolată la nodul mama GDet. Mai departe, schema DI head-functor verifică dacă trasaturile de acord ale GDet unifică cu cele ale centrului sau nominal. Acest mecanism este ilustrat în arborele de mai jos.

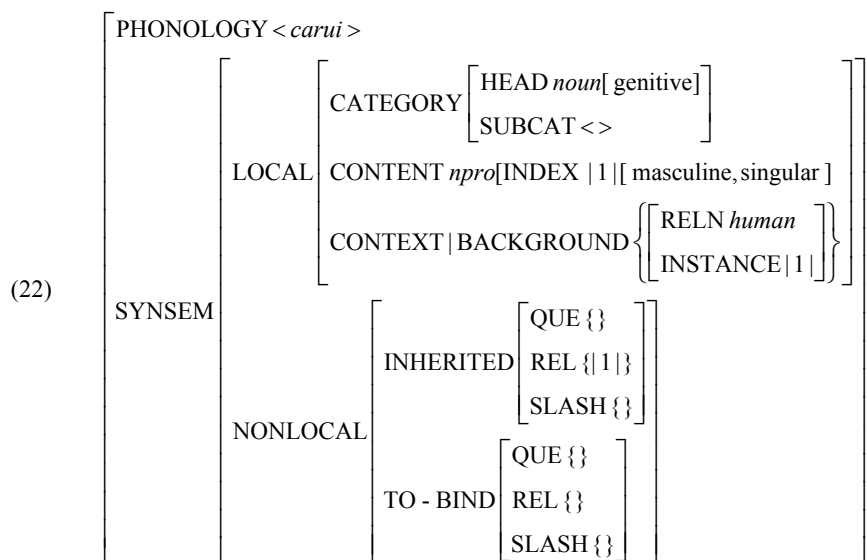
(21)



A doilea tip de acord, în schimb, ridică anumite dificultăți prin faptul că nu se realizează într-un arbore local, adică nu se realizează între ramurile surori ale aceluiași nod. Prin urmare, trasaturile de acord ale pronumelui relativ trebuie să percolate până la nivelul nodului P-rel (din (19)) pentru a putea fi controlate prin unificare de regula head-filler cu trasaturile de acord corespunzătoare substantivului determinat.

Mecanismul din teoria HPSG care dă seama de propagarea la distanță a anumitor trasături se numește mecanismul dependentelor la distanță și se aplică fenomenelor de limbă precum interogațiile, topicalizarile și, cum este cazul nostru, construcțiile relative. Aici ne vom ocupa numai de tratarea relativelor, pentru celelalte fenomene a se vedea [5].

Ideea principală a acestui mecanism este că pronumele relative poartă în intrările lor lexicale informații despre numele la care se referă. Intrarea lexicală a pronumelui relativ din exemplul nostru va conține, prin urmare, informațiile date în (22).



Valoarea trasaturii NONLOCAL | INHERITED indica acele trasaturi care vor fi supuse Principiului Trasaturilor Nonlocale. Aceste trasaturi pot fi specifice elementelor interogative, definite prin atributul QUE, elementelor dislocate, date de atributul SLASH sau pot fi specifice elementelor relative indicate prin atributul REL. Dupa cum se observa în (22), acest ultim atribut are în cazul de fata valoare non-vida, coindexata cu continutul semnificativ de masculin-singular al pronumelui.

Potrivit Principiului Trasaturilor Nonlocale, formulat în (23), valoarea atributului nonlocal INHERITED („mostenit”) este trecuta din nod în nod spre vârful arborelui pâna va întâlni o ramura sora ale carei trasaturi locale unifica cu cele mostenite.

(23) **Principiului Trasaturilor Nonlocale**

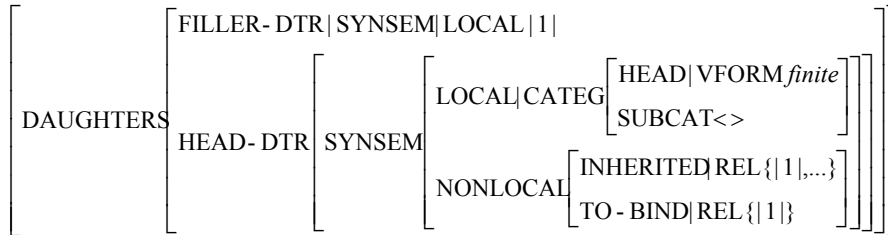
Pentru fiecare trasatura nonlocala, valoarea atributului INHERITED a nodului mama este egala cu reuniunea valorilor atributului INHERITED ale ramurilor fiice mai putin valoarea atributului TO-BIND a ramurii centru.

Atributul TO-BIND, practic, opreste propagarea trasaturilor mostenite în momentul în care se realizeaza elementul cautat, adica elementul care a facut necesara aceasta propagare. De exemplu, trasaturile de acord ale pronumelui relativ, în exemplul nostru *carui*, se propaga la nivelul propozitiei relative pâna când este realizat substantivul la care se refera acest pronume, adica *baiatul*.

Regula care asigneaza o valoare atributului TO-BIND în momentul în care are loc unificarea trasaturilor locale ale unui element cu trasaturile mostenite pe ramura centru este

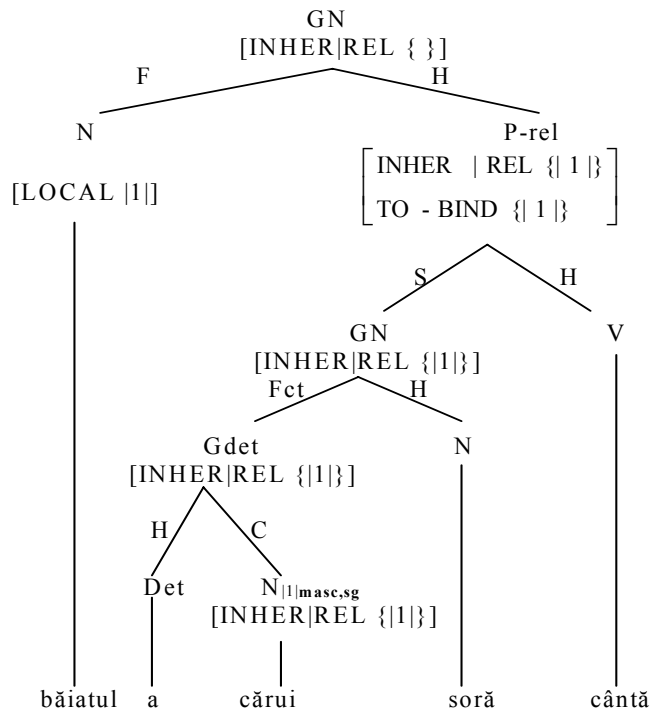
o schema de dominanta imediata numita *head-filler* (*filler* ar putea fi parafrazat drept „ceea ce vine sa completeze o lipsa”) si este descrisa în (24).

(24) Schema DI *head-filler*



În sfârșit, daca aplicam Principiul Trasaturilor Nonlocale si schema DI *head-filler*, acordul la distanta pe care îl avem în vedere se realizeaza în maniera ilustrata în arborele de mai jos.

(25)



În concluzie, acordul încrucișat avut în vedere presupune, pe de o parte, un acord local, cel dintre articolul genitival și substantivul determinat, în cazul nostru subiectul propoziției relative, iar pe de altă parte un acord la distanță, cel dintre pronumele relativ și substantivul determinat, exterior propoziției relative. Primul tip de acord se face pe baza Principiului Trasaturilor Centrale și a acordului banal dintre functor și centrul său, pe când cel de al doilea tip de acord face uz de Principiul Trasaturilor Nonlocale și de schema de Dominanță Imediată *head-filler*.

4. Concluzii

Analiza oferită aici pune în lumină faptul că un fenomen dificil precum acordul încrucișat poate fi tratat într-o manieră relativ simplă și elegantă cu ajutorul unei teorii lingvistice adecvate, cum este teoria Head-driven Phrase Structure Grammar.

Prin aparatul formal și adecvarea lingvistică pe care le oferă această teorie, descrierea fenomenelor limbii române devine incontestabil mai unitară, mai explicită și mult mai riguroasă. Acestor avantaje li se adaugă încă unul, extrem de important, acela al adecvării teoriei pentru implementarea informatică. Este deschis astfel drumul pentru construirea de gramatici computaționale ale limbii române și dezvoltarea componentei informatizate a acesteia.

Aplicațiile informatice ale teoriei HPSG sunt, de altfel, în plină dezvoltare și nu am dori să încheiem înainte de a aminti câteva aspecte în acest sens.

Modelul HPSG a făcut parte încă de la origine dintr-un sistem de tratare automată a englezei dezvoltat în laboratoarele de cercetare Hewlett Packard din Palo Alto ([20]). Apoi, au fost propuse diferite implementări, unele bazate pe sistemul PATR ([21]), altele realizate direct în Prolog ([22], [23]). Dintre implementările de sisteme de gestiune a structurilor de trasaturi tipologizate și cu moștenire, se poate cita sistemul *Typed Feature Structure* (TFS) al lui M. Emele și R. Zajac [24] și sistemul ALE al lui B. Carpenter [25].

Teoria HPSG a inspirat deopotrivă noul formalism european ALEP, a cărui implementare (în Prolog) presupune un mecanism de gestionare de gramatici și lexicoane, un analizor, un generator și un modul de transfer pentru traduceri automate. Este de altfel utilizat în mai multe centre de cercetare universitară (precum DFKI la Saarbrücken, *Center for Cognitive Science* în statul Ohio, CSLI la Stanford) sau industriale, în special la ATR în Japonia (pentru traducerea automată englezo-japoneză pentru stabilirea de întâlniri prin telefon).

O altă aplicație informatică a acestei teorii, pe cât de recentă, pe atât de importantă este cea cuprinsă în proiectul Verbmobil, [26], care s-a ocupat cu traducerea bidirecțională, în timp real, a textelor vorbite în trei limbi (germană, engleză și japoneză).

Head-driven Phrase Structure Grammar este o teorie care s-a impus incontestabil în lingvistica modernă atât prin numeroasele sale aplicații informatice, cât și prin

„generalitatea” aparatului sau care o face adecvata pentru numeroase limbi ale lumii, asa cum se poate vedea din impresionanta bibliografie electronica HPSG oferita de pagina www.dfki.de/lt/HPSG. Nu trebuie trecute cu vederea lucrarile de limba româna dezvoltate în acest cadru, dintre care le amintim pe cele ale lui Ionescu ([27]-[33]), Monachesi ([34]-[36]) si Barbu ([37]) la care s-ar cuveni sa se adauge multe altele spre afirmarea limbii române în lingvistica internationala.

Referinte bibliografice

- [1] Pollard, C. (1984) *Generalized Context-Free Grammars, Head Grammars and Natural Language*. Teză de doctorat. Universitatea din Stanford.
- [2] Kay, Martin (1979) ‘Functional Grammars’, *Actes 5° annual meeting of the Berkeley Linguistics Society*, Berkeley, pp. 142-158.
- [3] Oehrle, Richard; Bach, Emmon; Wheeler, Deirdre (eds.) (1988) ‘Categorial Grammars and Natural Language Structures’, Dordrecht: Reidel.
- [4] Pollard, C.; Sag, I. (1987) *Information-based Syntax and Semantics*, CSLI, University of Chicago Press.
- [5] Pollard, C.; Sag, I. (1994) *Head-driven Phrase Structure Grammar*, CSLI, University of Chicago Press.
- [6] Sag, I.; Pollard, C. (1991) ‘An integrated theory of complement control’, *Language*, 67:1, pp. 63-113.
- [7] Pollard, C.; Sag, I. (1992) ‘Anaphors in English and the scope of binding theory’, *Linguistic Inquiry*, 23:2, pp. 261-303.
- [8] Pollard, C. (1990) ‘On head non-movement’, *Actele Colocviului Discontinuous constituency*, Tilburg.
- [9] Nerbonne, J.; Netter, K.; Pollard, C. (eds.) (1993) ‘German grammar in HPSG’, CSLI, University of Chicago Press.
- [10] Balari, S. (1993) ‘Feature structures, linguistic information and grammatical theory’, Teza de doctorat, Universitatea Autonoma din Barcelona.
- [11] Gunji, T. (1987) *Japanese Phrase Structure Grammar*, Reidel.
- [12] Chung, C. (1993) ‘Korean auxiliary verb constructions without VP modes’, *Harvard Workshop on Korean Linguistics*, V; în C. Pollard, I. Sag (eds.), *Readings in HPSG*.
- [13] Miller, P.; Sag, I. (1993) *French clitic movement without clitics or movement*, LSA Meeting, Los Angeles
- [14] Monachesi, P. (1993) ‘Object clitics and clitic climbing in Italian HPSG grammar’ *Actes 6° European ACL*, Utrecht, pp. 431-437

-
- [15] Robinson, J. (1965) 'A machine-oriented logic based on the resolution principle', *Journal of the ACM*, 12, pp.23-44
- [16] Colmerauer, A. (1975) 'Les grammaires de métamorphose', Université d'Aix Marseille, reluat în L. Bolc (ed.) *Natural Language Communication with computers*, Springer, Verlag, 1978.
- [17] Kay, M. (1979) 'Functional grammars', *Actes 5° annual meeting of the Berkeley Linguistics Society*, Berkeley, pp. 142-158.
- [18] Abeillé, A. (1993) *Les nouvelles syntaxes. Grammaires d'unification et analyse du français*, Armand Colin, Paris.
- [19] Allegranza, V. (1998) 'Determiners as Functors: NP Structure in Italian' în S. Balari & L. Dini (eds.) *Romance in HPSG*, CSLI, Stanford.
- [20] Proudian, D.; Pollard, C. (1985) 'Parsing Head-driven Phrase Structure Grammar', *Actes 23° ACL*, Chicago, pp. 167-171.
- [21] Shieber, S. (1986) *An Introduction to unification-based theories of grammar*, CSLI, University of Chicago Press.
- [22] Oliva, K. (1990) 'Simple parser for an HPSG-style grammar implemented in Prolog', *Actes 13° COLING*, Helsinki, vol.3, pp.434-436.
- [23] Carpenter, B. (1991) 'The generative power of Categorical grammars and Head-driven Phrase Structure grammar with lexical rules', *Computational Linguistics*, 17:3, pp.301-314.
- [24] Emele, M.; Zajac, R. (1990) 'Typed-unification grammars', *Actes 13° COLING*, Helsinki, vol.3, pp.293-298.
- [25] Carpenter, B. (1992) 'The Logic of typed Feature Structures with Applications to Unification Grammars, Logic Programs and Constraint Resolution', Cambridge University Press [Implementarea sistemului ALE]
- [26] Wahlster, W. (ed.) (2000) *Verbmobil: Foundations of Speech-to-Speech Translation*, Springer, Berlin.
- [27] Ionescu, E. (1995-1996) 'A Type of SOV Construction in Romanian', "*Cahiers de Linguistique Théorique et Appliquée*", tomes XXXII-XXXIII, 19-39
- [28] Ionescu, E. (1995-1996), 'Accusative Weak Pronouns in Romanian', *Cahiers de Linguistique Théorique et Appliquée*, tomes XXXII-XXXIII, 1995-1996, 40-52
- [29] Ionescu, E. (1995-1996), 'Accusative Clitic Doubling in Romanian', *Cahiers de Linguistique Théorique et Appliquée* tomes XXXII-XXXIII, 1995-1996, 53-73
- [30] Ionescu, E. (1995-1996), 'Accusative Clitic Climbing in Romanian', *Cahiers de Linguistique Théorique et Appliquée*, tomes XXXII-XXXIII, 1995-1996, 74-87

-
- [31] Ionescu, E. (1995-1996), 'A Quantification-based Approach to Negative Concord in Romanian' in Geert-Jan M. Kruijff and Richard T. Oehrle (editori), *Proceedings of Formal Grammar Conference Utrecht*, 1999, p. 25-36
- [32] Ionescu, E. (1995-1996), *pro-Drop: An HPSG Account without Lexical Rules*, "Bucharest Working Papers in Linguistics", vol. I, nr.1, 1999, 117-124
- [33] Ionescu, E. (1995-1996), *On the Status of PE in the Direct Object Construction in Romanian*, Romanian Journal of Information Science and Technology, volume 4, numbers 3-4, 2001, p. 293-310
- [34] Monachesi, P. (1998) 'The morphosyntax of Romanian cliticization' în P.-A. Coppen, H. van Halteren, & L. Teunissen, eds., *Proceedingd of Computational Linguistics in The Netherlands 1997*, pp. 99-118, Amsterdam-Atlanta:Rodopi.
- [35] Monachesi, P. (1999) 'Linearization properties of the Romanian verbal complex' în *Proceedings of WECOL 98*, Tempe.
- [36] Monachesi, P. (2000) 'Clitic Placement in the Romanian verbal complex', în B. Gerlach and J. Grijzenhout (eds.) *Clitics in Phonology, Morphology and Syntax*, LA 36, Amsterdam: John Benjamins Publishing Company
- [37] Barbu, A.M. (1998) 'Romanian determiners:order and classification' în *Revue Roumaine de Linguistique*, XLIII, nr.5-6, pp.299-315, Bucuresti

Dupa 10 ani de experienta terminografica: noul model de date terminologice al TermRom

Dan MATEI

A. Preambul

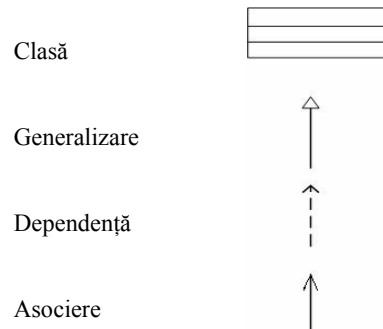
Din 1991 — când a fost înființată — Asociația Română de Terminologie (TermRom) a desfășurat o activitate terminografică materializată într-o bază de date proprie (accesibilă, în parte, pe web la www.cimec.ro/tr/) și într-o serie de publicații specifice. Formatul terminografic utilizat — descris în [Matei.1996] —, derivat din formatul standard MicroMATER (ISO 6156), se bazează pe un model de date (relativ) complex, serializat pe două nivele: nivelul conceptului și nivelul termenului. Practica terminografică (ce se traduce prin prelucrarea unei mari diversități de date terminologice) ne-a revelat o tensiune între complexitatea datelor reale și insuficienta complexitate a modelului folosit. În plus, necesitatea transferului de date între aplicații diverse a scos la iveală utilitatea consemnării cu o granularitate sporită a elementelor înregistrării terminologice. Mai mult, „entuziasmul” cu care ISO revizuieste standardele terminologice în ultimii ani⁶⁴, cu alte cuvinte, relativa instabilitate a standardelor din acest domeniu, îndeamnă la o și mai fină granularitate, pentru a spori șansele de compatibilitate cu normele de transfer viitoare. Pe de altă parte, pe măsura acumulării experienței, era din ce în ce mai limpede că modelul de date folosit ar trebui să acomodeze o mai mare diversitate și complexitate de metadate bibliografice, ca și o fină și flexibilă tratare a metadatelor „administrative”, de gestionare a colecției terminologice (vezi și [ISO 16642]).

Aceste considerente au dus la elaborarea unui model de date obiectual, care, pe lângă cerințele expuse mai sus, să fie și suficient de abstract ca să permită o serializare convenabilă (pentru transfer de date), — probabil bazată pe XML, de exemplu în formatul MARTIF [ISO 12200] — și să nu ceară elaborarea de aplicații informatice de o complexitate excesivă.

⁶⁴ Atât [ISO 12200] cât și [ISO 12620] sunt în revizie (deși ambele datează doar din 1999), iar [ISO 16642], este încă nedefinitivă. Desigur, această stare a lucrurilor probează și faptul că domeniul nu este încă bine „asezat”.

B. Modelul

Modelul este prezentat în continuare, într-un formalism UML [Unified Modelling Language] (mult simplificat), folosind următoarele notatii:



Conventional, modelul este împartit în secțiuni („pachete” [packages], în terminologia UML). La nivelul cel mai de sus, se disting secțiunea (asa zis) functională și secțiunea administrativă.

B.1. Secțiunea functională

În fig. 1 se prezintă clasele functionale esențiale și asocierile lor. Practic, orice element al modelului este o 'înregistrare'. Cu alte cuvinte, 'înregistrare' este clasa generică. Existența unei clase generice oferă — pe lângă gruparea proprietăților comune tuturor elementelor — și posibilitatea de a avea un identificator unic pentru fiecare înregistrare din baza de date ce implementează acest model.

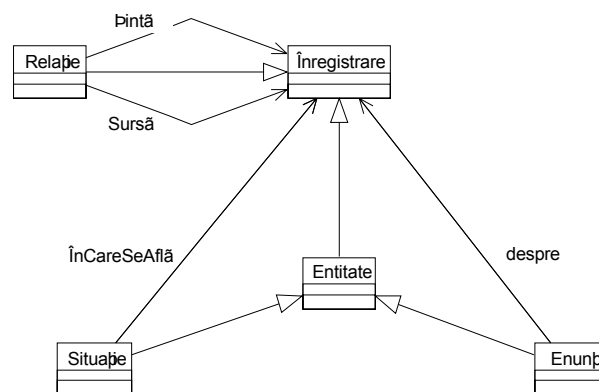


Figura 1 – Secțiunea functională (generică)

Clasa 'înregistrare' are doua subclase: 'entitate' (care grupeaza elementele ce au o existenta autonoma) si 'relatie' (care grupeaza asocierile binare între înregistrari). Se observa ca sunt acceptabile chiar si relatiile binare între relatii, lucru folositor si în practica.

Reificarea relatiilor binare între înregistrari simplifica mult modelul si constituie o maniera flexibila de a consemna o mare varietate de asocieri între elementele modelului. O relatie R poate avea doua caracteristici fundamentale, utile în cadrul modelului:

- a. simetria: daca x este în relatia R cu y, y este în relatia R cu x;
- b. tranzitivitatea: daca x este în relatia R cu y si y este în relatia R cu z, x este în relatia R cu z.

Pentru fiecare instanta a clasei 'relatie', aceste caracteristici (sau lipsa lor) se consemneaza ca un atribut al tipului respectiv de relatii (nereprezentat grafic în model)⁶⁵. Consemnarea acestor proprietati ale relatiilor poate fi foarte folositoare pentru programele care ar exploata baza de date.

Pentru a se rezolva (relativ) simplu si flexibil asocierile multiple între înregistrari, s-a introdus subclasa 'situatie' a clasei 'entitate'. Dupa cum se vede în figura, o instanta (sau mai multe) a clasei 'situatie' se asociaza cu o instanta a clasei 'înregistrare', iar obiectul 'situatie' este conectat cu oricâte alte elemente prin instante banale ale clasei 'relatie'. În practica, cele mai frecvente utilizari ale acestui tip de obiect sunt ca încarnari de contexte si evenimente. În fine, cea de-a doua subclasa a clasei 'entitate' este 'enunt'. Acest tip de obiect este destinat a consemna atribute ale unei înregistrari care n-au fost aprioric prevazute în model, cu alte cuvinte el gazduieste mentiuni pentru care se doreste un statut superior simplelor note, si anume care se doresc a fi colocabile si/sau indexabile.

În continuare se prezinta doar subsectiunile sectiunii functionale care sunt de interes în contextul acestui volum.

B.1.1. Sec?iunea terminologic?

Fig. 2 prezinta entitatile (i.e. subclasele clasei 'entitate') de natura terminologica.

⁶⁵ O categorie de relatii — importanta în terminologie — este cea a relatiilor ierarhice, i.e. cele tranzitive si asimetrice.

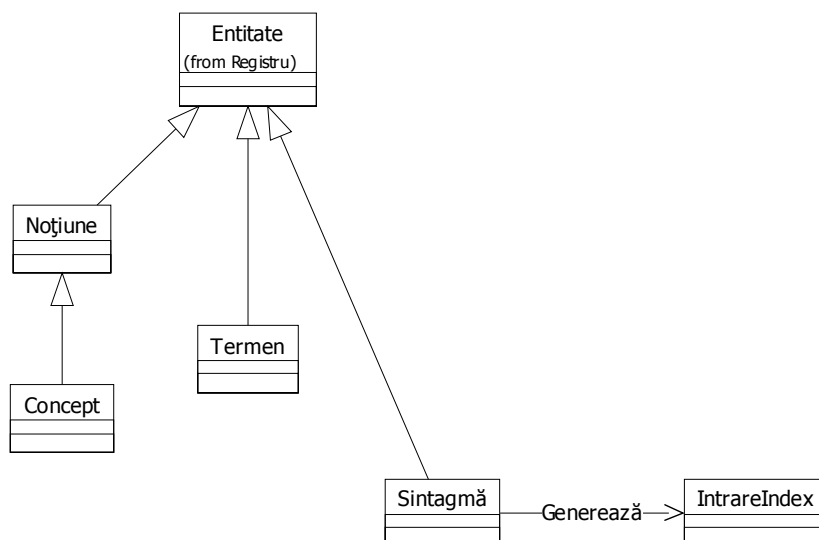


Figura 2 – Secțiunea terminologica

Principala clasa a acestei subsecțiuni este 'notiune'. Instancele ei consemnează notiunile vehiculate în baza de date terminologica, independent de limba. Din rațiuni practice, și anume din necesitatea de a cuprinde în baze de date terminologice și materialul organizat de obicei în tezaure terminologice, s-a decis să se cuprindă în modelul de date nu doar conceptele pure, ci și unități semantice mai largi, precum cele desemnate de termenii compusi într-un tezaur (sau ceea ce ISO 12620 numește 'unități frazeologice' [A.2.1.18]). Clasa acestor unități conceptuale care cuprinde conceptele și unitățile semantice mai largi este clasa 'notiune'. Distincție fină între 'notiune' și 'concept' este formulată în logica astfel [Chetan&Sommer.1978]:

Notiune: forma logică fundamentală care reflectă însușirile caracteristice necesare și generale ale unei clase de obiecte.

Concept: notiune care reflectă însușirile esențiale ale unei clase de obiecte⁶⁶.

Asadar, o notiune care nu e concept cuprinde mai mulți factori semantici, deci poate fi factorizată⁶⁷.

⁶⁶ Exemple de noțiuni care nu sunt concepte: "bărbat blond", "scriitor important".

⁶⁷ O regulă simplă, pragmatică de a distinge o notiune care este concept de una care nu este, ni se pare: notiunea care e concept și-ar găsi locul într-un dicționar, pe când cea care nu e, nu.

A doua subclasa a acestei secțiuni este 'termen'. Instanțele ei consemnează doar „denumirile” conceptelor (A.1. în ISO 12620). Cu alte cuvinte, consemnează ceea ce au în comun o familie de expresii lingvistice ce desemnează un concept⁶⁸. Expresiile lingvistice propriu-zise sunt consemnate în instanțele clasei 'sintagma'⁶⁹. Din pricina faptului că un termen poate fi exprimat printr-un set de expresii lingvistice (flexiuni, variante ortografice etc.), s-a preferat separarea „termenului” de expresiile sale lingvistice, în felul acesta nu ne conformăm strict definiției pentru 'termen', din ISO 12620 (A.1): "a designation of a defined concept in a special language by a linguistic expression".

Se poate observa în figura faptul că sintagmele generează intrări de index. În fapt, o sintagma poate genera — prin inversare/permutare — mai multe intrări de index, dacă terminograful decide că asta ar fi în folosul utilizatorilor, prin colocarea sintagmei la fiecare „factor” semnificativ. Exemple:

Sintagma	Intrări de index
efect Doppler	efect Doppler Doppler, efect
pseudofonetism	pseudofonetism fonetism, pseudo-
completivă indirectă anticipată	completivă indirectă anticipată indirectă anticipată, completivă anticipată, completivă indirectă

Clasa 'relatie' este vitală pentru consemnarea asocierilor între entitățile modelului. Pentru a ilustra modul în care se consemnează informația terminologică esențială, în fig. 3 s-au reprezentat tipurile de relații esențiale care asociază, pe de o parte, conceptele cu termenii care le desemnează, iar pe de altă, termenii cu sintagmele care-i exprimă. De asemenea, se vede cum o „situație” (care — în această ilustrare — implică (cel puțin) un loc, o perioadă și un agent) caracterizează designarea.

⁶⁸ Exemple de "familie de expresii lingvistice" sunt: a) *cladire, cladiri*; b) *expresiv, expresiva, expresivi, expresive*.

⁶⁹ în acest context, 'sintagma' desemnează — printr-un abuz de limbaj — atât sintagme cât și cuvinte.

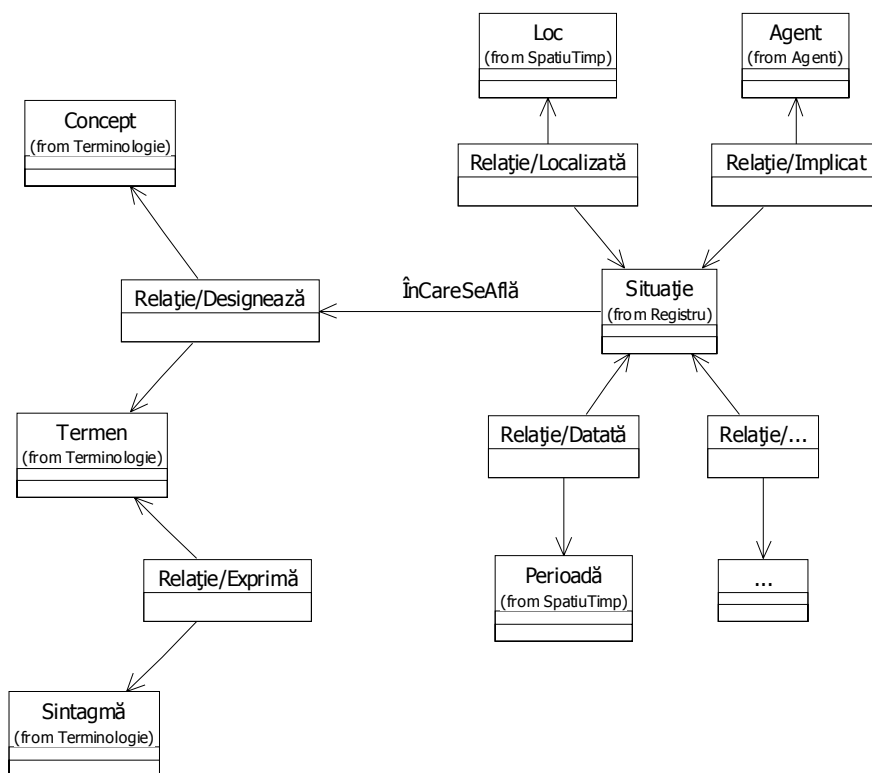


Figura 3 – Ilustrare a reprezentării informației terminologice

Într-o astfel de schema, se pot reprezenta cu acuratețe cazuri precum:

a) Concept: *mic arbust cu flori roșietice din familia ericaceae ...*

- Relație/designează:
 - Termen (științific) [latina]:
 - Relație/exprima:
 - Sintagma: *Kalmia latifolia*
- Relație/designează:
 - Situatie/context:
 - Relație/localizează:
 - Loc: *nordul Statelor Unite*
 - Termen [engleza]:

-
- Relatie/exprima:
 - Sintagma: *mountain laurel*
 - Relatie/designeaza:
 - Situatie/context:
 - Relatie/localizeaza:
 - Loc: *sudul Statelor Unite*
 - Termen [engleza]:
 - Relatie/exprima:
 - Sintagma: *calico bush*
 - Relatie/designeaza:
 - Situatie/context:
 - Relatie/localizeaza:
 - Loc: *sudul Statelor Unite*
 - Termen [engleza]:
 - Relatie/exprima:
 - Sintagma: *sheep's bane*
 - Relatie/designeaza:
 - Termen [româna]:
 - Relatie/exprima:
 - Sintagma [s.m.sg.]: *laur de munte*
 - Relatie/exprima:
 - Sintagma [s.m.pl.]: *lauri de munte*
- b) Concept: comandant de calarime
- Relatie/designeaza:
 - Situatie/context:
 - Relatie/localizeaza:
 - Loc: Moldova
 - Relatie/localizeaza:
 - Loc: Tara Româneasca
 - Relatie/dateaza:
 - Perioada: sec. XVII-XVIII
 - Termen [româna]:
 - Relatie/exprima:
 - Sintagma [s.m.sg.]: *serdar*
 - Relatie/exprima:
 - Sintagma [s.m.pl.]: *serdari*
- c) Concept: boier de rang mijlociu

- Relatie/designeaza
Situatie/context:
Relatie/dateaza:
Perioada: *sec. XVIII-XIX*
Termen [româna]:
Relatie/exprima:
Sintagma: *serdar* [s.m.sg.]
Relatie/exprima:
Sintagma: *serdari* [s.m.pl.]

Tot ca o ilustrare, în fig. 4 se prezintă modul cum se consențează etimologia unui termen, cu ajutorul clasei 'situație': o situație de tip 'etimologie' se asociază cu termenul de baza, iar termenii din care acesta provine sunt asociați cu situația prin intermediul unor relații de tip 'provineDin'.

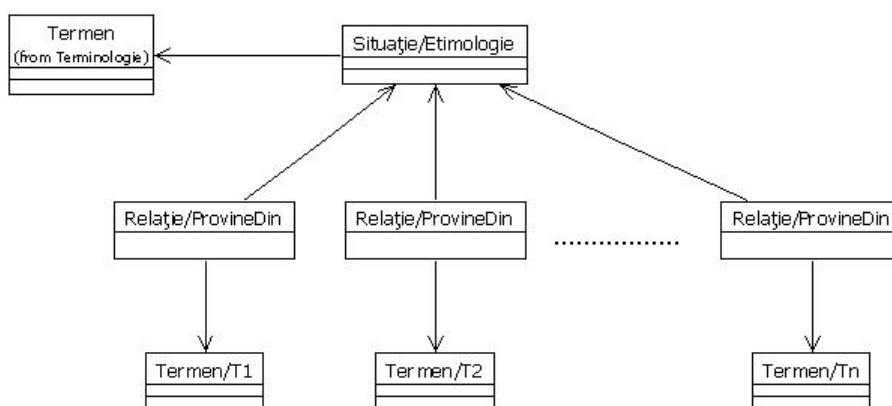


Figura 4 – Ilustrare a reprezentării etimologiei

De pilda:

- Concept: fixat la vârf
- Relatie/designeaza:
Termen [româna]:
Relatie/exprima:
Sintagma: acrofix
Situatie/etimologie:

Relatie/provine din:
 Termen [greaca]:
 Relatie/exprima:
 Sintagma: acro

Relatie/provinedin:
 Termen [latina]:
 Relatie/exprima:
 Sintagma: fixus

B.1.2. Sectiunea bibliografica

Fig. 5 prezinta entitatile (i.e. subclasele clasei 'entitate') de natura bibliografica, cu alte cuvinte este o sectiune de metadate. Sectiunea pare simpla, deoarece o buna parte din multitudinea de date bibliografice sunt consemnate cu ajutorul relatiilor. Clasa esentiala este 'editie'; cea care consemneaza fisă bibliografică a unei editii citate.

Entitatea 'lucrare' consemneaza metadatele specifice unei creatii (mai ales textuale, în cazul nostru), i.e. „abstractizeaza” ceea ce au în comun toate editiile unei lucrari. Utilitatea ei imediata este colocarea tuturor manifestarilor unei lucrari, indiferent de limba sau editie. O subclasa importanta a clasei 'lucrare' este entitatea 'serial'. Aici se consemneaza si periodicele, adica entitatile ce grupeaza instantele clasei 'NumarPeriodic', cu alte cuvinte publicatiile-gazda ale articolelor. Discutia asupra acestor clase si a relatiilor între ele depaseste cadrul acestui articol.

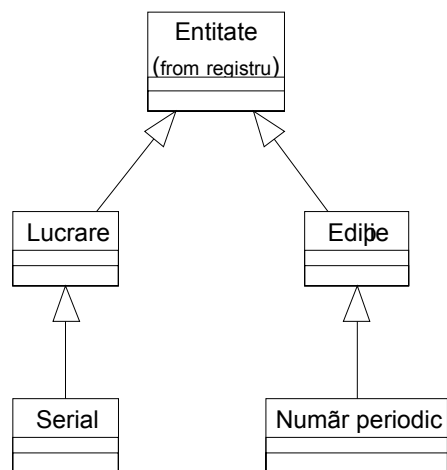


Figura 5 – Sectiunea bibliografica

B.2. Sectiunea administrativa

În fig. 6 se prezinta clasele de natura administrativa si relatiile esentiale între ele. Rolul acestor clase este de a consemna modificarile survenite în baza de date, în succesiunea lor. În acest fel se poate urmări geneza înregistrărilor si se pot identifica responsabilitățile. În plus, deoarece se prevede și stocarea datelor modificate, se creează premisele revenirii la stări anterioare ale bazei de date. În instanțele clasei 'interventie' se consemnează fiecare modificare operată asupra unei înregistrări. Fiecare asemenea instanță este asociată — prin intermediul instanțelor clasei 'contributie' — cu agentul (i.e. operatorul) care a produs-o. În plus o intervenție este asociată și cu sursele ei documentare. Se observă cum clasa 'referință' poate avea ca instanțe atât referințe bibliografice (citând o editie), cât și referințe personale (citând o comunicare personală).

Clasa 'ÎnregistrareArhivă' este foarte importantă, instanțele ei fiind chiar versiunile „desuete” (i.e. cele dinaintea de modificări) ale atributelor înregistrărilor.

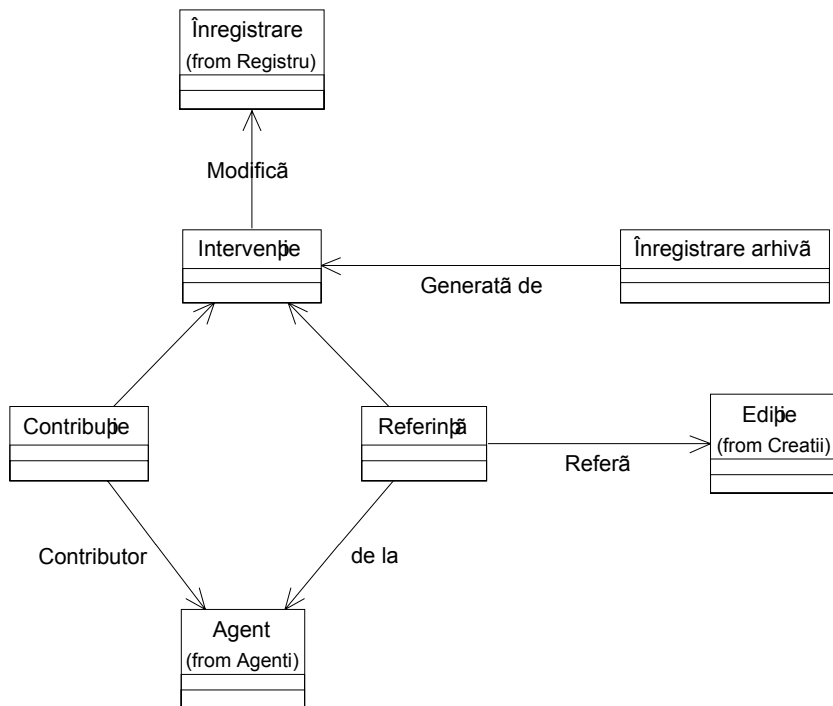


Figura 6 – Sectiunea administrativa

C. Remarci finale

Modelul prezentat pare suficient de flexibil pentru a satisface cerințele funcționale atât ale unei baze de date terminologice, cât și a unei lexicografice (mai ales datorită distincției între termeni și expresiile lor lingvistice). El este și suficient de abstract pentru ca schema unei baze de date ce l-ar folosi ca fundament să fie relativ comodă la implementare.

TermRom are în curs un proiect de elaborare a unei astfel de baze de date terminologice. După finalizarea acesteia, este de așteptat un proces traumatic de convertire a bazei de date curente. Sporul de funcționalitate obținut va compensa însă efortul.

D. Referințe

- [Chețan&Sommer.1978] Chețan, Octavian, Radu Sommer. *Dicționar de filozofie / Coordonare științifică Octavian Chețan, Radu Sommer*. — București: Editura Politică, 1978
- [ISO 12200] *ISO 12200:1999, Computer applications in terminology – Machine-readable terminology interchange format (MARTIF) – Negotiate interchange*
- [ISO 12620] *ISO 12620:1999, Computer applications in terminology – Data categories*
- [ISO 16642] *ISO/CD 16642:1999, Computer applications in terminology – Metamodel for representing terminological data collections*
- [Matei.1996] Matei, Dan. *Banca de date terminologice a TermRom și problemele ei neologice*, în *Limbaj și Tehnologie / Dan Tufiș – editor*. — București: Editura Academiei Române, 1996

Probleme de reprezentare a datelor terminografice într-o baza de date relationala

Dr. Sorin GHETARU

Oriunde și oricând se creează, comunică, înregistrează, prelucrează, stochează, transformă sau refolosește informație sau cunoștințe de specialitate este implicată într-un fel sau altul și terminologia. Comunicarea într-un anumit domeniu a devenit un discurs specializat cu texte de specialitate diferențiate în nenumărate forme. Atunci când se definește terminologia ca o mulțime structurată de concepte și denumirile lor într-un anumit domeniu, ea poate fi socotită ca fiind infrastructura cunoașterii de specialitate. Scrierea textelor tehnice și documentarea tehnică devin astfel imposibile fără o utilizare corectă a unor resurse terminologice. Deoarece producerea textelor tehnice implică frecvent mai multe limbi, terminologiile multilingve de înaltă calitate au devenit bunuri mult dorite greu de găsit pe înfloritoarea piață a industriilor limbajelor și cunoașterii.

Există numeroase baze de date terminologice disponibile pentru interogare on-line sau pe CD-ROM (TERMIUM, EURODICAUTOM), pe dischete sub forma unor dicționare electronice sau ca baze de date personale realizate și întreținute de ingineri, specialiști în calculatoare, chimiști care lucrează ca terminologi, traducătorii, autori de texte tehnice. Aceste baze de date sunt utilizate pentru:

- traducere asistată de calculator;
- scrierea de texte tehnice și științifice asistată de calculator;
- sisteme informatice (administrarea componentelor etc.);
- cercetări terminologice în lingvistică, filozofia științei, sociologia tehnologiei etc.

Pentru asemenea obiective au fost dezvoltate aplicații specializate (programe de management al bazelor de date terminologice), unele disponibile pe piața terminologică internațională, altele ca prototipuri în cadrul unor proiecte de cercetare academice.

MARTIF este formatul standardizat pentru managementul informației terminologice. Posibilitatea organizării terminologiei în baze de date având formate diferite face nerealistă presupunerea că s-ar putea cădea de acord asupra unui anumit format de bază de date relatională, așa cum este SQL, care să fie folosit pentru schimburile terminologice. De aceea s-a mers pe linia producerii unui format la dispoziția publică fără obligații materiale și care să fie independent de platforma de lucru. Rezultatul este **MARTIF (Machine-Readable Terminology Interchange Format)** cunoscut și ca ISO 12200.

În ISO 12620 sunt descrise 150 de categorii de date, un număr imens care nu urmărește decât să le arate pe cele posibile și modul în care acestea pot fi structurate. Categoriile MARTIF sunt împărțite în 10 secțiuni grupate în 4 clase. Acestea sunt:

- termen: cuprinde categoria de date termen (1);
- informație în legătură cu termenii: conține informația legată de termeni (2) și informația privind gradul de echivalență;
- informație descriptivă: relație cu domeniul (4), descrierea conceptului (5), relații între concepte (6), categorii de date care leagă un concept de poziția sa în sistemul de concepte (7), note (8);
- informație administrativă: categorii de date care leagă un concept de un element al unui tezaur sau de o altă formă de documentare (9), categorii de date care cuprind informații administrative.

Un avantaj major al faptului că MARTIF este scris folosind cod SGML este acela că, deși se poate aprecia că lectura codului nu este facilă, ea este totuși posibilă ca urmare a faptului că nu face apel decât la caracterele ASCII. Un alt avantaj al sistemului MARTIF este acela că el acceptă referințe către alte documente chiar din interiorul documentului. Inițial MARTIF presupune că înainte de implementarea produselor software pentru importul sau exportul datelor programatorii sunt obligați să examineze sursele implicate. Pentru a asigura un acces așa numit “orb” care să permită oricui să transfere baze de date terminologice din orice sistem spre sau dinspre MARTIF este necesară o standardizare suplimentară a categoriilor de date, domeniilor specifice etc.

Tabela ce urmează enumeră acea parte a “elementelor” MARTIF care sunt de cea mai mare importanță pentru realizarea unei resurse terminologice Multilingve.

<termEntry>	Set complet unic de date terminologice pentru un concept exprimat într-o singură limbă, și cuprinzând unul sau mai mulți termeni și datele descriptive și administrative asociate lor, sau, în cazul unei abordări bilingve sau multilingve, două sau mai multe concepte foarte apropiate, exprimate în fiecare limbă, precum și datele descriptive și administrative asociate lor. Atributele includ: type, care clasifică setul de date terminologice conform categoriile de date specificate de ISO 12620.
<langSet>	Limba; în cadrul unui element <termEntry> va fi folosit pentru a grupa mai multe <tig> și <ntig> asociate unei singure limbi. Prezența atributului lang este obligatorie, în afara cazului în care el este mostenit.

<tig>	Grup de informatii terminologice; în cadrul unui element <termEntry>, va contine elemente de informatii asociate cu un singur termen, fiecare dintre acestea functionând la acelasi nivel; cu alte cuvinte nu este permisa imbricarea între elementele subordonate unui <tig>. Prezenta atributului lang este obligatorie, în afara cazului în care el este mostenit.
<ntig>	Grup încuibat de informatii terminologice; va fi folosit în cadrul unui element <termEntry> daca anumite elemente informationale sunt asociate mai curând cu elemente interne, decât cu întregul <tig>. Urmatoarele elemente vor fi folosite în cadrul <ntig> pentru a gazdui alte date terminologice: <termGrp>, <termNoteGrp>, <descripGrp> si <adminGrp>. Prezenta atributului lang este obligatorie, în afara cazului în care el este mostenit.
<term>	Va contine un termen format dintr-un singur cuvânt sau din mai multe cuvinte, sau o desemnare simbolica privita ca un termen tehnic.
<termGrp>	Va contine un element <term> si, posibil, cel putin înca un element încuibat în plus fata de termen.
<termNote>	Va contine informatii legate de termen. Atributele includ: type, care clasifica <termNote> conform categoriilor de date specificate în ISO 12200.
<termNoteGrp>	Va contine un element <termNote> si posibil cel putin un element încuibat în plus fata de informatia legata de termen. Va fi folosit pentru a gazdui un nivel suplimentar de imbricare în cadrul elementului <termGrp>
<descrip>	Va contine informatii descriptive precum definitia, contextul sau explicatii descriind concepte si termeni. Atributele includ: type, care clasifica <descrip> potrivit categoriilor de date specificate în ISO 12200.
<descripGrp>	Va contine un element <descrip> si, posibil, cel putin un element imbricat în plus fata de informatia descriptiva.
<admin>	Va contine date administrative. Atributele includ: type, care clasifica <admin> în functie de categoriile de date specificate în ISO 12200.
<adminGrp>	Va contine un element <admin> si, posibil, cel putin un element imbricat în plus fata de informatiile administrative.
<date>	Va contine o singura data de formatul YYYY-MM-DD, cu optiunea notarii data-timp YYYY-MM-DD hh:mm:ss. Atributele includ: type, care clasifica <date> dupa categoriile specificate în ISO 12200.
<note>	Va contine o nota sau o adnotare drept comentariu legat fie de un întreg <termEntry>, un întreg <tig> sau <ntig> ori de unul din elementele <...Grp>.

<descripNote>	Va fi folosit în cazul informațiilor de tipul <note> folosite în cadrul <descripGrp> când conținutul notei este legat de o lista de opțiuni.
<adminNote>	Va fi folosit în cazul informațiilor de tipul <note> folosite în cadrul <adminGrp> când conținutul notei este legat de o lista de opțiuni.
<ptr> ⁷⁰	Va consta dintr-un indicator către o altă locație din documentul curent. Atributele includ: type, care clasifică <ptr> conform Anexei A, A.12 target, care precizează destinația referirii, ca unul sau mai mulți identificatori SGML.
<ref>	Va defini o referire către o altă locație din documentul curent, în termeni de unul sau mai multe elemente identificabile. <ref>GI este asociat cu text suplimentar drept conținut al elementului, deci consta dintr-o eticheta-start cu o tinta integrată, urmată de textul asociat și închisă de o eticheta-sfârșit. Atributele includ: type, care clasifică <ref> conform Anexei A. target, care precizează destinația referirii ca unul sau mai mulți identificatori SGML.
<xref> ⁷¹	Va defini o referință la un grafic, ilustrație, figură, tabel sau alt document extern sau fișier folosind o notație indicativă extinsă ca valoare a atributului tinta a <xref>, de ex. <xref target='documentIdentifier'>, unde valoarea 'documentIdentifier' este un cod de identificare pentru documentul tinta. Utilizatorul va documenta notația indicativă extinsă care este folosită incluzând un comentariu adecvat în elementul <encodingDesc> ale header DTD. Atributele includ: type, care clasifică <xref> conform Anexei A. target, care precizează destinația referirii ca unul sau mai mulți identificatori SGML.
<hi> ⁷²	Va fi folosit pentru a marca un cuvânt sau o frază ca evidențiat grafic în contrast cu textul înconjurător. Atributele includ: type, care clasifică <ref> conform Anexei A. target, care precizează destinația referirii ca unul sau mai mulți identificatori SGML.

⁷⁰ Nota – <ptr> GI nu poate fi asociat cu text suplimentar drept conținut al elementului, întrucât consta doar dintr-o eticheta-start cu o tinta integrată. Elementele <ptr>, <ref> și <xref> sunt toate considerate link-uri pentru că ele conectează locația lor curentă cu o altă locație tintită în cadrul unui document sau cu o locație externă documentului.

⁷¹ Nota – Elementele externe tintite de <xref> trebuie să fi accesibile sistemului-tinta pentru scopuri de importare.

⁷²Nota – În managementul terminologiei o utilizare frecventă a <hi> se face pentru a sublinia termeni necesari, adică termeni folosiți într-o definiție, nota sau alt material textual care sunt definiți altundeva în resursa terminologică. Vezi de asemenea Anexa A, A.2.2.2.

<foreign>	Va identifica un cuvânt sau o fraza ca aparținând altei limbi decât cea a textului înconjurător. Atributele includ: lang, care identifica limba cuvântului sau frazei marcate.
<refObjectList>	Va fi folosit în back-matter și va conține unul sau mai multe obiecte back-matter, mai ales resurse comune ca: date bibliografice, date de responsabilitate, identificatori de namespace (URL-uri și FPI-uri), material textual la care se fac referiri dese, liste de locații geografice, fișiere externe și altele asemenea. Atributele includ: type, care clasifică <refObjectList> după categoriile de date specificate în ISO 12620 Anexa A, A.11.4.1.
<refObject> ⁷³	Va conține o dată constând în general dintr-o resursă comună ca: date bibliografice, date de responsabilitate, identificatori de namespace (URL-uri și FPI-uri), material textual la care se fac referiri dese, liste de locații geografice, fișiere externe și altele asemenea. Datele bibliografice ar trebui să rezide în back matter sau într-un document extern (caz în care se va face referire la datele bibliografice din back matter folosind elementul <xref>). Atributele includ: type, care clasifică <refObject> după categoriile de date specificate în ISO 12620 Anexa A, A.11.4.2. Dacă se specifică altfel, tipul <refObject> este moștenit de la <refObjectList> respectiv.
<itemSet>	Va fi folosit în back matter și va conține unul sau mai multe obiecte individuale care în mod tradițional sunt grupate împreună, de ex. obiectele numele autorului și prenumele autorului vor fi grupate împreună într-un <itemSet> de tip=autor Atributele includ: type, care clasifică <itemSet> în principal conform categoriilor de date listate în ISO 12620 Anexa B. Totuși acest Standard Internațional nu specifică întregul spectru al categoriilor de date care pot fi folosite cu <itemSet>
<item>	Va conține un exemplu individual de informație back matter. Atributele includ: type, care clasifică <itemSet> în principal conform categoriilor de date listate în ISO 12620 Anexa B pentru informații bibliografice. Totuși acest Standard Internațional nu specifică întregul spectru al categoriilor de date care pot fi folosite cu <item>

⁷³ Nota – Unele documente terminologice cuprind date bibliografice complete în format nediferențiat drept conținut al categoriei de date sursă (vezi ISO 12620:1999, A.10.19). Această practică încurajează redundanța și efortul mare pentru îngrijirea datelor. Aceste informații ar trebui convertite în obiecte back matter (informații bibliografice) dacă este posibil.

<itemGrp>	Va contine unul sau mai multe <item> împreuna cu <ptr>, <ref> sau <note>. Atributele includ: type, care clasifica <item> în principal conform categoriilor de date listate în ISO 12620 Anexa B pentru informatii bibliografice. Totusi acest Standard International nu specifica întregul spectru al categoriilor de date care pot fi folosite cu <itemSet>
------------------------	--

Din acest tabel au mai fost eliminate elementele (aproape la fel de numeroase) specifice informatiilor bibliografice. Instantierea elementelor enumerat mai sus se face prin intermediul „categoriilor de date” standardizate de ISO 12620. Numarul acestora este de aproximativ 200. În cea mai ampla resursa terminologica (EURODICAUTOM) sunt în prezent prezente mai puțin de 20 astfel de categorii de date.

Uniunea Europeana în activitatea sa este unul dintre utilizatorii majori ai procedurilor de tradulare a textelor si terminologiei. Aceasta se datoreste partial faptului ca legislatia sa este direct aplicabila în statele membre si de aceea ea trebuie sa fie disponibila în toate limbile de lucru oficiale. Ca rezultat, traducatorii Comisiei Europene produc mai mult de 1 milion de pagini pe an si au de-a face cu cel puțin 6-7 milioane de termeni (în medie sunt 8 sau 9 termeni care ridica probleme pe fiecare pagina).

Unitatea pentru Terminologie a Comisiei Europene este destinata asigurarii suportului lingvistic pentru toate limbile oficiale ale Uniunii Europene. Au fost elaborate glosare de specialitate, multe dintre le în noua limbi. Domeniile acoperite sunt tratatele importante cum ar fi cele de la Maastricht si Roma, cele economice si administrative (Taxa pe Valoarea Adaugata, buget) dar si unele legate de subiectele centrale sau puternic inovatoare ale stiintei si tehnologiei (fizica plasmei, biotehnologie, minerit). Deosebit de rolul lor de resurse terminologice si de surse terminologice pentru domeniile de inovare, aceste glosare documenteaza ceea ce se numeste “Eurolect”, adica frazele si cuvintele care își au origina în cadrul Uniunii Europene si pentru care nu exista echivalente nationale.

Monitorizând toate modificarile aparute ca urmare a unei evolutii permanente a bazei de date EURODICAUTOM am constatat ca, recent, a avut loc schimbarea suportului hardware si odata cu aceasta pot fi observate urmatoarele:

- Indicarea mult mai frecventa a referintei la documentul sursa a termenului;
- Indicarea frecventa a referintei la documentul sursa al definitiei acestuia;
- Indicarea documentului sursa si pentru sinonime si abrevieri;
- Utilizarea mai frecventa a notelor pentru adaugarea unor informatii suplimentare asupra termenilor, acestea putând fi grupate astfel:
 - o {NTE} explicatii si informatii generale asupra termenilor;
 - o {TXT} contextul (de cele mai multe ori un exemplu de utilizare a termenului respectiv);
 - o {GRM} informatii gramaticale (gen, numar);

-
- o {USG} indicarea mediului în care este utilizat termenul: “*technical jargon*”;
 - o {REG} nota asupra unor utilizari locale speciale sau asupra regionalismelor;
 - o {DOM} indicarea unui domeniu sau subdomeniu care completează clasificarea obisnuită folosită anterior și care a rămas încă prezentă.

De asemenea se prevede ca în cel mai scurt timp să fie implementate următoarele:

- afișarea tuturor caracterelor și diacriticelor (ca și a informației nelingvistice, dacă se cere);
- îmbunătățirea sistemului de clasificare a domeniilor;
- introducerea link-urilor interne și externe.

Modelele de date terminologice orientate în exclusivitate către terminologie au avantajul de a fi relativ intuitive pentru terminolog. Transcrierea directă a elementelor și relațiilor dintre acestea într-o bază de date este din ce în ce mai dificilă și mai riscantă.

Există încercări meritorii de realizare a unor interfețe “cuprinzătoare” pentru consultarea resurselor terminologice. Exemplele următoare sunt edificatoare în acest sens.

Primul exemplu ar putea provoca comentarii legate de complexitatea reală a înregistrării referințelor bibliografice cele mai obișnuite.

K - Club Cycom
 Copyright Cycom Limited 2002 (<http://www.cycom.co.uk/>)
 These details identify the source of some text appearing within one of the term entries.

Identifier	iso1087-1.2	Generate unique identifier
Author given name		
Author family name	TC 37/SC 1	
Article title		
Page numbers		
Book title	087-1.2 Terminology work - Vocabulary - Part 1: Theory and application	
ISBN		
Book edition	Draft	
Publication date (YYYY-MM-DD)	1999-04-22	
Publisher		

Commit changes Commit changes and close Rollback changes and close

Al doilea, ne determina sa luam în considerare urmatoarele:

La nivelul Uniunii Europene numarul limbilor pentru care este necesar suport terminologic este atât de mare (si speram înca în crestere) încât nu mai este posibila multiplicarea tabelor bazelor de date potrivit numarului de limbi de lucru. Din fericire, “balizarea” documentelor permite identificarea si prelucrarea corect dependenta de limba în care au fost concepute acestea. Se vine astfel în sprijinul “globalizarii” aplicatiilor informatice care sunt suport al resurselor terminologice multilingve dând posibilitatea acceptarii, prelucrarii si prezentarii numeroaselor scrisuri, formate de date si limbi existente. În acelasi timp trebuie adaptata si interfata utilizator potrivit locului si culturii careia îi apartine acesta printr-un proces nu mai putin important de “localizare”

Multa vreme, prelucrarea automata a datelor a fost considerata satisfactor realizabila prin utilizarea setului ASCII de caractere. În prezent este însa absolut necesar ca:

- Utilizatorul calculatorului sa poata tasta caractere si simboluri (vest-europene, est-europene, grecesti si cirilice, cel putin) folosind o claviatura standard.
- Aplicatia sa prelucreze si sa afiseze sau imprime siruri de caractere formate corect folosind seturi de caractere specifice fiecarei limbi.

Aceste cerinte pot fi realizate prin valorificarea calitatilor standardului Unicode de codificare prin utilizarea unor coduri de 16 biti pentru reprezentarea tuturor caracterelor pentru calculatoarele moderne care includ simbolurile tehnice si semnele speciale necesare imprimarii textelor.

Cu alte cuvinte la nivelul seturilor de semne necesare unei resurse terminologice multilingve se poate conta pe serviciile standardului Unicode si pe cele ale oricarei baze de date relationale care accepta Unicode.

Pentru indicarea formatelor de prezentare (fonte, punere în pagina, seturi de caractere) si a limbii utilizate se face apel la balizare astfel încât la nivelul câmpului vom gasi siruri de caractere Unicode balizate.

Înscrierea datelor terminologice este facilitata de înscrierea lor în „categorii de date” bine definite (vezi ISO 12620). Dar numarul mare al acestor categorii si mai ales

incidenta ridicata a aparitiilor neprevazute dinainte a unora noi face imposibila alocarea unui câmp de date fiecărei categorii de date. Aceeasi observatie poate fi facuta si asupra relatiilor dintre diferitele categorii de date care reflecta direct relatiile dintre elementele MARTIF. O solutie este o abstractizare suplimentara a datelor terminologice dupa încadrarea lor succesiva în siruri de caractere balizate, categorii de date, elemente MARTIF.

În centrul modelului de date se afla un set de 13 entitati (atomi):

Entitate	Descriere
data category	o anumita clasa de informatii terminologice (de exemplu: term, part of speech)
data category name	un nume agreat de utilizator (user-friendly), dependent de limba, al unei anumite categorii de date (de exemplu, în româna, "termen" pentru term)
data category index type	o strategie de indexare corespunzatoare unei anumite categorii de date (ISO 12620) (de exemplu: nu se indexeaza, se indexeaza ca valoare unica, se indexeaza cuvânt cu cuvânt)
lang	o anumita limba, care dispune de o schema de codare uniforma care utilizeaza un singur set de caractere (de exemplu: French, German, Italian)
charset	o combinatie unica de caractere care poate fi utilizata pentru reprezentarea unei singure sau mai multor limbi (de exemplu: ISO 8879-1. ISO 8859-2)
picklist	o multime de valori posibile ale unor date terminologice aparținând unei anumite categorii de date (ISO 12620) (de exemplu, pentru categoria "parte de vorbire": noun, verb, adjective)
element	o data terminologica unica
date value	o data (time stamp) care constituie valoarea unui element
number value	un numar care constituie valoarea unui element
picklist value	un membru al unei liste care reprezinta valoarea unui element
text value	sir de caractere care constituie valoarea unui element
index value	un sir de caractere care reprezinta forma normalizata indexata a unui element particular sau a unei parti a acesteia
link	legatura între doua elemente

Primele 6 "articole" sunt "meta-entitati"; ele sunt create si tabellele corespunzatoare sunt completate cu informatii înainte de încarcarea oricarei date terminologice în baza de date. Prin completarea acestor table se contureaza si se activeaza chiar modelul de date al bazei de date terminologice. Cu alte cuvinte, ansamblul "meta-tabelelor" defineste structura care impune conditii si unifica datele terminologice de nivel molecular. Ele pot fi considerate atomi catalizatori ai reactiilor necesare combinarii altor atomi în interactiuni moleculare.

Celelalte 7 entitati se încarca direct prin proceduri de introducere a datelor sau prin import si cuprind datele terminologice vizibile pentru utilizatorul bazei de date. Informatiile continute de aceste entitati pot fi validate la nivel molecular folosind interogari SQL

standard. Majoritatea interogarilor formulate de utilizatorii bazei de date se concentreaza aproape în întregime asupra informatiilor încarcate în aceste entitati.

Elementul central al aplicatiei pentru întretinerea unei astfel de baze de date este componenta de tip *parser* pentru crearea, validarea si prelucrarea documentelor MARTIF în particular (fara a ignora documentele SGML, HTML, XML). În mod obisnuit un *parser* este un modul software care examineaza un document SGML prin confruntarea acestuia cu DTD-ul corespunzator. Rezultatul acestei examinari este de cele mai multe ori simplu: 'da' în situatia în care documentul reprezinta o instantiere valida a DTD-ului si 'nu' în cazul contrar. De cele mai multe ori *parser*-ul este capabil sa 'normalizeze' documentul validat (aducându-l la o 'forma canonica') astfel încât faciliteaza formatarea, editarea si încarcarea documentului în baza de date.

Alaturi de *parser* si legat de acesta se afla un *editor structurat*. Pornind de la DTD acesta propune utilizatorului pas cu pas optiunile de compunere, sau modificare a unui document în conformitate cu definitia tipului corespunzator documentului. În cazul în care obiectivul este compunerea unui document SGML el poate asigura completarea *tag*-urilor necesare.

De cele mai multe ori sistemele de management al bazelor de date orientate spre text folosesc *fisiere inversate de indexare* a continutului acestora pentru regasirea informatiilor. Cautarea poate urmari aparitia unui anume cuvânt, sau a unui model oarecare într-un document sau în o parte a acestuia. Identificarea subdiviziunilor documentului se poate face folosind tocmai *tag*-urile cu acesta este marcat, respectiv modul în care acestea au fost transcrise în relatiile dintre tabelele bazei de date.

În fine, o componenta deosebit de importanta este aceea care realizeaza functiile de *import-export* ale datelor terminologice spre si dinspre baza de date.

Terminologia calitatii

Realizarea unor resurse terminologice multilingve este de mai multa vreme în centrul preocuparilor Asociatiei Române pentru Terminologie (TERMROM). Începând de anul trecut pe lista temelor având aceeasi orientare se înscrie proiectul "Terminologie armonizata cu prevederile EURODICAUTOM în domeniul calitate si standardizare". Proiectul a fost initiat de Ministerul Educatiei si Cercetarii si este finantat în cadrul Programului CALIST.

Obiectivele principale ale acestui subprogram sunt:

- Asigurarea flexibilitatii necesare pentru a raspunde operativ la cerintele concrete de rezolvare a unor teme de cercetare care decurg din prioritatile stabilite prin strategiile guvernamentale adoptate pe domenii specifice, în procesul integrarii României în U.E.

- Asigurarea condițiilor de dezvoltare și armonizare a sistemului de standarde naționale în conformitate cu cerințele organismelor de standardizare europene și internaționale;
- Asigurarea unei baze terminologice științifice pentru elaborarea standardelor de calitate românești, precum și în ceea ce privește condițiile de aplicabilitate a prevederilor standardelor internaționale și europene adaptate ca standarde românești;
- Clarificarea condițiilor pe care trebuie să le îndeplinească produsele românești în vederea patrunderii lor pe piața unică a Uniunii Europene și produsele introduse în România.

Pentru realizarea obiectivelor proiectului au fost prevăzute următoarele activități:

- Întocmirea unui Proiect Terminologic pentru definirea și înregistrarea terminologiei domeniilor calitate și standardizare utilizate în documentele oficiale ale Uniunii Europene, conform prevederilor EURODICAUTOM și standardelor internaționale;
- Extragerea, traducerea și structurarea terminologiei domeniilor calitate și standardizare;
- Proiectarea, programarea și implementarea unei Baze de date conform Proiectului Terminologic capabilă să gestioneze toate domeniile EURODICAUTOM;
- Înregistrarea în baza de date a terminologiei domeniilor calitate și standardizare;
- Elaborarea unei aplicații informatice de administrare a bazei de date terminologice și de transfer de date terminologice conform formatului standard ISO pentru lucrul în rețea;
- Realizarea unui site web pentru promovarea Bazei de date terminologice și punerea acesteia la dispoziția publicului.

A fost avizat Proiectul Terminologic, au fost stabilite cerințele pe care să le satisfacă suportul informatic, s-a constituit un fond de termeni specifici extrasi din EURODICAUTOM și din Tezaurul rațional al CEI și au fost demarate activitățile pentru realizarea unei baze de date relaționale EUROCAST pentru înregistrarea acestora.

Bibliografie

1. **ISO 639:1988**
Code for the representation of names of languages
2. **ISO 639-2:1998**
Code for the representation of names of languages - Part 2: Alpha-3 code

-
3. **ISO 704:2000**
Terminology work - Principles and methods
 4. **ISO 860:1996**
Terminology work - Harmonization of concepts and terms
 5. **ISO 1087-1:2000**
Terminology work - Vocabulary - Part 1: Theory and application
 6. **ISO 1087-2:2000**
Terminology work - Vocabulary - Part 2: Computer applications
 7. **ISO 1951:1997**
Lexicographical symbols particularly for use in classified defining vocabularies
 8. **ISO 6156:1987**
Magnetic tape exchange format for terminological/lexicographical records (MATER)
 9. **ISO 10241:1992**
Preparation and layout of international terminology standards
 10. **ISO 12199:2000(E)**
Alphabetical ordering of multilingual terminological and lexicographical data represented in the Latin alphabet
 11. **ISO 12200:1999**
Computer applications in terminology - Machine-readable terminology interchange format (MARTIF) - Negotiated interchange
 12. **ISO/TR 12618:1994**
Computer aids in terminology - Creation and use of terminological databases and text corpora
 13. **ISO 12620:1999**
Computer applications in terminology - Data categories
 14. **ISO 15188:2001**
Project management guidelines for terminology standardization

SECTIUNEA II

TEHNOLOGII ALE LIMBAJULUI SCRIS

Ro-balkanet - ontologie lexicalizata, în context multilingv, pentru limba româna

Dan TUFIS, Institutul de Inteligența Artificială, Academia Română, București,
tufis@racai.ro

Dan CRISTEA, Facultatea de Informatică, Universitatea A.I.Cuza, Iași
dcristea@infoiasi.ro

Rezumat

Cerintele creării unei ontologii multilingve de tipul EuroWordNet sunt frecvent contradictorii și dacă problemele de compatibilitate nu sunt considerate în etapele timpurii ale construcției o armonizare tardivă se poate dovedi dificilă sau imposibilă. Mai exact, există două probleme majore de compatibilitate care trebuie avute în vedere și anume: acoperirea conceptuală – în sensul că fiecare lexicon monolingv ar trebui să conțină lexicalizări ale aceluiași fond conceptual și coeziunea interpretativă – în sensul că interpretarea relațiilor folosite în fiecare din ontologiile cuprinse în ontologia multilingvă trebuie să fie identică. În lucrare sunt discutate ambele aspecte și prezentate soluțiile adoptate în vederea satisfacerii criteriilor de consistență și coerență multilinguală a wordnet-ului pentru limba română.

1. Limba, resurse lingvistice și comunicare electronică

Cercetarea în domeniul tehnologiilor limbajului este un domeniu ce are deja istorie în știința calculatoarelor, dar, actualmente, motivațiile sale depășesc sfera interesului pur științific sau comercial. Pastrarea identității limbilor și culturilor naționale în cadrul globalizant al societății informaționale și a cunoașterii readuce în actualitate avertismentul lui Alain Danzin (1992): „**În era electronică, este esențial pentru supraviețuirea unei limbi ca ea să fie folosită în sistemele de informare electronică.**” Avansul științific și tehnologic obținut în cei 10 ani scurși de la raportul prezentat de Danzin Comisiei Europene, a condus la maturizarea unor teorii, tehnologii, metode și la dezvoltarea altora noi, dar mai ales a permis definirea unor standarde pentru realizarea unitară a ceea ce generic se numește *resurse lingvistice fundamentale* ale unei limbi. Caracterul multilingual al societății cunoașterii, în care conceptul de „unitate prin diversitate” se referă în primul rând la prezervarea limbilor și culturilor actuale, a generat o deosebită efervescență,

puternic stimulata de organisme internationale – în primul rând Comisia Europeană – asupra cercetării în domeniul resurselor multilingve. Metodologic, tehnologia limbajului natural creează o distincție netă între prelucrări și date, între „masinaria software de prelucrare a limbajului” numită și *lingware* și cunostintele lingvistice, numite cum aratăm *resurse lingvistice*, necesare funcționării acestei mașinării. Dihotomia *lingware* - *resurse lingvistice*, susținută de standardele de reprezentare și codificare a cunostintelor lingvistice permite dezvoltarea independentă a celor două componente ale unui sistem de prelucrare a limbajului. *Lingware*-ul este independent de limba și intra tot mai pregnant în zona ingineriei software. El poate fi dezvoltat de specialiști de oriunde fără ca aceștia să fie preocupați de limba pentru care va fi folosit. Resursele lingvistice însă sunt de competența specialiștilor vorbitori nativi ai limbii respective. În condițiile în care aceste resurse lingvistice sunt realizate în conformitate cu standardele sau practicile internaționale, ele pot fi integrate în sistemele de comunicare electronică, nu doar pentru prelucrare monolingvă ci mai ales pentru prelucrări multilingve. Beneficiile alinierii la standardele internaționale în realizarea resurselor lingvistice sunt enorme, și putem considera un exemplu foarte simplu. Să presupunem că suntem interesați de un anumit subiect și, folosind imensul ocean informațional ce este Internet-ul, apelăm la un așa numit „motor de căutare”, un program a cărui funcționalitate asigură identificarea documentelor electronice ce conțin informații potențial relevante pentru subiectul nostru de interes. Acest gen de serviciu informațional este asigurat de „motoare de căutare” precum Google, Altavista, Excite și multe altele. Documentele interesante din punctul nostru de vedere ar putea să fie scrise în limba engleză, franceză, germană, română sau orice altă limbă. Dar pentru a le regăsi pe toate, indiferent în ce limbă am formulat cererea noastră de regăsire, motorului general de căutare îi sunt necesare resursele lingvistice specifice limbilor în care documentele ar putea exista. Dacă aceste resurse lingvistice există pentru engleză, franceză, germană, italiană etc. și ele sunt reprezentate în același format standardizat, rezultatul cercetării noastre documentare va fi o colecție de documente tratând subiectul de interes în oricare dintre aceste limbi. Un astfel de serviciu, numit regăsire documentară multilingvă este o realitate pentru toate limbile „mari”, o calificare ce nu are acoperire în substratul cultural ci doar în ceea ce se numește „nivelul de informatizare al limbii”. Procesul de informatizare a unei limbi naturale permite potențarea și diseminarea ei prin mijloacele tehnologice ale societății informaționale.

2. Lexicalizarea abordărilor în tehnologia limbajului și conceptul „wordnet”

Lexicul este fără îndoială cea mai importantă resursă lingvistică a unei limbi. Marea majoritate a cercetării actuale, atât în lingvistica formală cât mai ales în tehnologia limbajului, plasează componenta lexicală în centrul modelelor de limbă, sub influența a ceea ce a fost numită abordarea *lexicalizată* sau *lexicalista* a studiului limbii. Nu este de mirare, deci, enormul interes pentru dezvoltarea de resurse lexicale multilingve. Studiul computațional al dicționarelor electronice, natura informației ce trebuie inclusă în ele și tipul de prelucrări pe care le poate facilita o anumită structurare a unui mare volum lexical

a fost, fara îndoiala, fundamental influentat de proiectul WordNet, lansat în urma cu mai mult de 25 de ani la Universitatea din Princeton sub conducerea reputatului psiholingvist George Miller. WordNet, resursa publica, este o uriasa retea semantica lexicala în care peste 100.000 de *întelesuri* lexicalizate în limba engleza prin mai mult de 130.000 de cuvinte sunt asociate între ele prin relatii semantice si/sau lexicale (Fellbaum, 1998). Fondul lexical este distribuit în 4 retele semantice corespunzând categoriilor gramaticale deschise: substantive, verbe, adjective si adverbe. Notiunea de *înteles* (*meaning*) este în WordNet echivalata cu cea de concept si este reprezentata printr-o serie sinonimica în care fiecare cuvânt al seriei are asociat un numar ce identifica sensul în care cuvântul respectiv are întelesul asociat conceptului. Seria sinonimica ce identifica un înteles se numeste *sinset*. Relatiile existente între sinseturi sunt de diferite tipuri, depinzând de categoria gramaticala a cuvintelor ce alcatuiesc un anumit sinset (antonimie/sinonimie, hiponimie/hiperonimie, holonimie/meronimie, troponimie etc.). Influenta proiectului WordNet a fost enorma în domeniul tehnologiei limbajului (exprimata poate si prin faptul ca acum, în limbajul tehnic cel putin, cuvintele „wordnet” si „synset” au devenit substantive comune, importate prin calchiere în mai toate limbile) iar beneficiile acestui concept sunt atât de evidente încât Comisia Europeana, între 1996 si 1998, a finantat un proiect similar de mare anvergura numit EuroWordNet (Blokma et al., 1996). Acest proiect, extrem de ambitios si-a propus nu numai realizarea concertata de wordneturi monolingve pentru limbile europene de circulatie internationala (engleza, franceza, germana, italiana, olandeza, spaniola) dar a introdus o cerinta fundamental noua, anume corelarea multilinguala a celor 6 retele semantice lexicales, astfel încât într-un sinset al unei limbi sa se poata ajunge în echivalentul de traducere al oricaror celorlalte 5 limbi. Fata de relatiile originale din WordNet, EuroWordNet propune un inventar mult mai bogat (90) de relatii cum ar fi cele tematice de tip casual (Agent, Patient, Instrument, Location, Direction) sau cele corelând sensurile derivatilor lexicali (XPOS-SYNONYMY: a adora - adoratie).

Solutia tehnica pentru corelarea multilinguala a retelelor semantice monolingve a fost definirea unui index interlingual (ILI), independent de limba, continând reprezentari conceptuale ale întelesurilor lexicalizabile în limbile proiectului. Fiecare înteles din oricare din limbile reprezentate în reteaua semantica multilingva este pus în corespondenta, în general, cu un singur concept al indexului interlingual. Aceste corespondente se realizeaza prin intermediul a 20 de tipuri distincte de relatii binare. Sinseturile (seriile sinonimice) din doua sau mai multe limbi care sunt puse în corespondenta cu acelasi concept din ILI sunt considerate echivalenti de traducere, natura echivalentei de traducere fiind definita de tipul relatiilor ce definesc corespondenta dintre sinseturile respective si conceptul comun.

Initial, indexul multilingual a fost constituit ca o multime nestructurata a tuturor întelesurilor lexicalizate în WordNet (cu alte cuvinte în engleza). Ulterior, prin dezvoltarea wordneturilor monolingve, ILI a fost îmbogatit si cu reprezentari conceptuale cu lexicalizare ce nu se regasesc în engleza.

O alta inovatie a proiectului EuroWordNet a fost adoptarea unei multimi de primitive semantice, independente de limbaj, în termenii carora asa-numitele *concepte de*

baza din ILI au fost asociate cu descrieri *ontologice*. Prin importul acestor descrieri la nivelul lexicalizarilor prin echivalenți de traducere (și, prin mostenire, la hiponimii acestora) în fiecare dintre wordneturile monolingve, în EuroWordNet se poate vorbi de o ontologie lexicală multilingvă. O prezentare în detaliu a proiectului EuroWordNet se poate găsi în (Vossen, 1998).

Dupa 3 ani, proiectul EuroWordNet inițial a fost extins pentru o perioadă de încă doi ani (EuroWordNet II) și a încorporat încă 4 limbi: basca, catalana, ceha și estoniană. Proiectul EuroWordNet II s-a încheiat în anul 2000 cu realizarea unor nuclee a caror extensie a ramas în exercitiul financiar al autoritatilor nationale.

3. Limba română în contextul proiectului BALKANET, extensie a EuroWordNet

În septembrie 2001 a fost lansat proiectul european BALKANET (IST – 2000 – 29388), o continuare firească a proiectului EuroWordNet II care aduce alături de cele 10 limbi europene alte 5 limbi din zona balcanică: bulgăra, greacă, română, sârbo-croata, turca (Stamou et al., 2002). Ca și în EuroWordNet, ontologiile lexicale monolingve sunt corelate printr-o multitudine de concepte interlinguale, corespondențele fiind stabilite cu ajutorul unor relații de echivalență complexe (*eq-synonymy*, *eq-near-synonymy*, *eq-has-hyperonym*, *eq-has-hypernym*, etc.).

Reprezentanții din România în acest proiect, care va dura trei ani, sunt Institutul Academiei Române de Cercetări pentru Inteligența Artificială din București (coordonator Dan Tufis) și Facultatea de Informatică a Universității A.I. Cuza din Iași (coordonator Dan Cristea) și în realizarea obiectivelor proiectului sunt implicați numeroși specialiști, atât informaticieni cât și lingviști. Desigur, participarea românească în acest proiect și angajarea față de obiectivele proiectului nu s-a bazat numai pe entuziasm ci pe activități și rezultate anterioare importante, pe *surse lingvistice* primare (Tufis, 2001) de referință ale limbii române, implementate ca *resurse lingvistice* (ibid.) în format standardizat și pe o multitudine de programe de prelucrare dezvoltate de-a lungul a mulți ani de cercetare, în cea mai mare parte prin finanțare internațională.

3.1. Corpusuri

În cadrul proiectelor europene Multext-East și TELRI (Erjavec et al., 1997), (Dimitrova et al., 1998), (Tufis, Bruda, 1997), (Tufis et al., 1997, 1999) a fost creat un corpus paralel în 7 limbi, foarte detaliat adnotat, bazat pe romanul “1984” al lui Orwell și un alt corpus paralel în 25 de limbi, bazat pe “Republica” lui Platon. Adnotarea folosită inițial a fost conformă cu standardul TEI (<http://www.tei-c.org/>), dar ulterior, odată cu cristalizarea standardului CES (Ide, 1998), corpusurile au fost re-adnotate (automat) în conformitate cu CES. Acestea sunt două corpusuri relativ mici (câte aproximativ 110.000 cuvinte în fiecare limbă) dar, datorită acurateții proceselor de etichetare și de aliniere

(validate manual), au fost extrem de folosite pentru diverse aplicatii, de la construirea modelelor lingvistice pentru etichetare morfo-sintactica (Tufis, 1999), clasificare a documentelor (Tufis et al., 2000), extragere de echivalenți de traducere (Tufis, 2002), până la discriminarea automată a sensurilor (Ide et al., 2002). Pe lângă corpusurile multilingve s-au construit alte două corpusuri monolingve mult mai mari: un corpus literar bazat pe diverse romane (continând aproximativ 1.500.000 cuvinte) și un corpus jurnalistic (continând peste 100.000.000 cuvinte). Ambele corpusuri au fost segmentate, etichetate și lematizate automat⁷⁴.

3.2. Dictionare explicative: WEB-LEX și XML-LEX

Principalul dictionar pe care l-am folosit în analiza noastră este Dictionarul Explicativ al Limbii Române (DEX, 1996), referința lexicografică pentru limba română contemporană, dictionar realizat de Institutul de Lingvistică „Iorgu Iordan”⁷⁵ al Academiei Române. În urma analizelor statistice de frecvență în corpusurile menționate, au fost selectate și introduse în format electronic cele mai frecvente 23.000 de cuvinte titlu din DEX. Acest nucleu DEX a fost convertit într-o bază de date lexicale în cadrul proiectului european CONCEDE (*CONortium for Central European Dictionary Encoding*) (Tufis et al., 1999) și al proiectului prioritar al Academiei WEB-LEX (Tufis, 2000). Ulterior, îmbogățit continuu prin culegere manuală din alte câteva dictionare explicative (DEX'84, DOOM, DLRM), la inițiativa unor tineri entuziaști atât din țară cât și din diaspora (vezi de pildă: <http://dex.francu.com>), WEB-LEX a fost corectat sub aspect sintactic-structural și codificat într-un format standardizat, respectând convențiile lexicografice utilizate de DEX și, în măsura posibilului, conținutul său textual. Uneori, din considerente legate de consistența structurală, s-au operat o serie de modificări asupra conținutului. De asemenea, o serie de erori evidente în sursa primară au fost corectate de specialiști avizați. Deși mai bogat (în prezent WEB-LEX conține aproape 70.000 de intrări, față de cele circa 56.000 de intrări din DEX'96), influența DEX a fost fundamentală în dezvoltarea WEB-LEX. Pe de altă parte, eventualele critici asupra conținutului, acolo unde ne-am despartit de DEX, în nici un caz nu trebuie puse în seama Institutului de Lingvistică „Iorgu Iordan-Al. Rosetti” ci a noastră. Din acest motiv, preferăm să ne referim la WEB-LEX ca la un dictionar *de tip* DEX și nu ca varianta computațională a DEX-ului.

Codificarea conținutului WEB-LEX, s-a realizat folosind limbajul de adnotare XML. Implementarea, ce explicitează toate convențiile tipografice precum și informațiile implicite, a condus la un volum textual de date de circa 8-10 ori mai mare față de conținutul textual echivalent al DEX-ului. Adnotarea XML a fost realizată automat, cu ajutorul compilatorului **DIC** (Tufis, 2000). Compilatorul a fost generat automat folosind JavaCC[®], pe baza unei gramatici LL(7) ce descrie structura formală a intrărilor în DEX. **DIC** poate fi folosit pentru a genera documente XML (conform cu DTD-ul CONCEDE) pentru orice

⁷⁴ Toate aceste resurse pot fi găsite pe situl Consorțiului de Informatizare pentru Limba Română (ConsILR) la adresa <http://consilr.info.uaic.ro>

⁷⁵ Noua sa denumire este Institutul de Lingvistică "Iorgu Iordan-Al. Rosetti"

dictionar ce foloseste conventiile tipografice adoptate în DEX. În (Vintila-Radulescu, 2002) sunt prezentate o multitudine de dictionare realizate sau aflate în curs de realizare la Institutul de Lingvistica „Iorgu Iordan-Al. Rosetti” si presupunând ca ele urmaresc conventiile tipografice si lexicografice adoptate în DEX, toate aceste surse lingvistice de referinta pentru limba româna ar putea fi transformate, cu efort minim, în resurse computationale fundamentale pentru prelucrarea automata.

Varianta codificata a dictionarului nostru este numita XML-LEX iar structura sa este descrisa de DTD-ul (*Document Type Definition*) pe care îl reproducem în figura 1, dezvoltat în cadrul proiectului CONCEDE.

```
<!-- CONCEDE project - Deliverable DR2.1: concede.dtd -->
<!-- copyright CONCEDE project consortium, 1999 -->
<!-- ENTITY DECLARATIONS -->
<!ENTITY % a.global '
  id ID #IMPLIED
  n CDATA #IMPLIED
  lang IDREF #IMPLIED' >
<!ENTITY % a.text '
  %a.global;
  rend CDATA #IMPLIED
  wsd CDATA #IMPLIED' >
<!ENTITY % basetags '
  (orth|pron|hyph|syll|stress|pos|gen|case|number|gram|tns|
  mood|q|source|gloss|usg|def|per|aspect|degree|voice|eg|
  etym|xr|trans|itype|subc)' >
<!ENTITY % dictbase.seq '#PCDATA | na' >
<!-- STRUCTURAL ELEMENTS -->
<!ELEMENT dictionary (body) >
<!ATTLIST dictionary %a.global;
  type CDATA #IMPLIED
  version CDATA #REQUIRED
  xml:space (default | preserve) 'preserve' >
<!ELEMENT body (entry+) >
<!ATTLIST body %a.global; type CDATA #IMPLIED >
<!ELEMENT entry
```

```

    (hw, (%basetags;|struc|alt|brack)*)    >
<!ATTLIST entry %a.global; type CDATA #IMPLIED >
<!ELEMENT struc (%basetags; | struc | alt | brack)* >
<!ATTLIST struc %a.global; type CDATA #IMPLIED >
<!ELEMENT trans (%basetags; | struc | alt | brack)* >
<!ATTLIST trans %a.global; type CDATA #IMPLIED >
<!ELEMENT alt (%basetags; | brack )* >
<!ATTLIST alt %a.global; type CDATA #IMPLIED >
<!ELEMENT brack (%basetags;)* >
<!ATTLIST brack %a.global; type CDATA #IMPLIED >
<!-- CONTENT ELEMENTS -->
<!ELEMENT voice (%dictbase.seq;)* >
<!ATTLIST voice %a.text; >
<!ELEMENT tns (%dictbase.seq;)* >
<!ATTLIST tns %a.text; >
<!ELEMENT syll (%dictbase.seq;)* >
<!ATTLIST syll %a.text; >
<!ELEMENT subc (%dictbase.seq;)* >
<!ATTLIST subc %a.text; >
<!ELEMENT stress (%dictbase.seq;)* >
<!ATTLIST stress %a.text; >
<!ELEMENT source (%dictbase.seq;)* >
<!ATTLIST source %a.text; >
<!ELEMENT pos (%dictbase.seq;)* >
<!ATTLIST pos %a.text; >
<!ELEMENT per (%dictbase.seq;)* >
<!ATTLIST per %a.text; >
<!ELEMENT number (%dictbase.seq;)* >
<!ATTLIST number %a.text; >
<!ELEMENT na (#PCDATA) >
<!ATTLIST na %a.text; >
<!ELEMENT mood (%dictbase.seq;)* >

```

```
<!ATTLIST mood %a.text; >
<!ELEMENT m (%dictbase.seq;)* >
<!ATTLIST m %a.text; >
<!ELEMENT lang (%dictbase.seq;)* >
<!ATTLIST lang %a.text; >
<!ELEMENT itype (%dictbase.seq;)* >
<!ATTLIST itype %a.text; >
<!ELEMENT hw (%dictbase.seq;)* >
<!ATTLIST hw %a.text; >
<!ELEMENT gram (%dictbase.seq;)* >
<!ATTLIST gram %a.text; >
<!ELEMENT gen (%dictbase.seq;)* >
<!ATTLIST gen %a.text; >
<!ELEMENT degree (%dictbase.seq;)* >
<!ATTLIST degree %a.text; >
<!ELEMENT case (%dictbase.seq;)* >
<!ATTLIST case %a.text; >
<!ELEMENT aspect (%dictbase.seq;)* >
<!ATTLIST aspect %a.text; >
<!ELEMENT hyph (%dictbase.seq;)* >
<!ATTLIST hyph %a.text; >
<!ELEMENT eg (source | q | gloss)* >
<!ATTLIST eg %a.global; >
<!ELEMENT pron (%dictbase.seq;)* >
<!ATTLIST pron %a.text; type CDATA #IMPLIED >
<!ELEMENT q
  (%dictbase.seq; | gloss | ptr | xptr |oref)* >
<!ATTLIST q %a.text; type CDATA #IMPLIED >
<!ELEMENT etym
  (%dictbase.seq; | gloss | lang | m | ptr | xptr |oref)* >
<!ATTLIST etym %a.text; type CDATA #IMPLIED >
<!ELEMENT xr (%dictbase.seq; | ptr | xptr )* >
```

```

<!ATTLIST xr %a.text; type CDATA #IMPLIED >
<!ELEMENT def (%dictbase.seq; | ptr |xptr |oref |usg)* >
<!ATTLIST def %a.text; type CDATA #IMPLIED >
<!ELEMENT gloss (%dictbase.seq; | ptr |xptr |oref)* >
<!ATTLIST gloss %a.text; type CDATA #IMPLIED >
<!ELEMENT orth (%dictbase.seq; | ptr |xptr |oref |usg)* >
<!ATTLIST orth %a.text;
  expansion NMTOKEN #IMPLIED
  extent (full | pref | suff | part) "full"
  type CDATA #IMPLIED >
<!ELEMENT usg (%dictbase.seq;)* >
<!ATTLIST usg %a.text;
  type (syn|hyper|colloc|comp|plev|acc|lang|gram|obj|
  subj|verb|hint|geo|dom|register|time|style|
  hyponym | antonym | other) "other" >
<!ELEMENT oref EMPTY >
<!ATTLIST oref %a.text;
  target IDREF #IMPLIED
  fullform NMTOKEN #IMPLIED >
<!ELEMENT ptr EMPTY >
<!ATTLIST ptr %a.text;
  corresp IDREFS #IMPLIED
  next IDREF #IMPLIED
  prev IDREF #IMPLIED
  type CDATA #IMPLIED
  resp CDATA #IMPLIED
  crdate CDATA #IMPLIED
  targType NMTOKEN #IMPLIED
  targOrder (y | n | u) "u"
  evaluate (all | one | none) #IMPLIED
  target IDREFS #REQUIRED >
<!ELEMENT xptr EMPTY >

```

```
<!ATTLIST xptr %a.text;
  corresp IDREFS      #IMPLIED
  next IDREF          #IMPLIED
  prev IDREF          #IMPLIED
  type CDATA         #IMPLIED
  resp CDATA         #IMPLIED
  crdate CDATA       #IMPLIED
  targType NMTOKEN   #IMPLIED
  targOrder (y | n | u)  "u"
  evaluate (all | one | none) #IMPLIED
  target NMTOKEN      #REQUIRED >
```

Figura 1: DTD-ul Concede, utilizat la implementarea XML-LEX

Aceasta structura de codificare a fost adoptata în implementarea unui numar mare de dictionare, reprezentând un standard „de facto” în lexicografia computationala actuala (Erjavec *et al.*, 2000). Detalii suplimentare privind semantica entitatilor folosite în codificare si a atributelor acestora, pot fi gasite în documentatia tehnica a proiectului la adresa www.itri.bton.ac.uk/projects/concede/. În tabelul de mai jos, sunt exemplificate reprezentarea tipografica (de tip DEX) si reprezentarea codificata în XML.

DEX	XML-LEX
<p>ZA², zale, s.f. 1. Fiecare dintre ochiurile unui lanț; p. gener. (la pl.) lanț. ♦ Lănțișor de metal întrebuințat uneori ca podoabă. ♦ Cusătură în formă de lănțișor, executată de obicei la broderii. 2. (La pl.) Împletitură executată din inele mici de fier legate unul de altul; p. ext. armură făcută din această împletitură, cu care se îmbrăcau oștenii în antichitate și în evul mediu, spre a se apăra de loviturile dușmanilor. [Var.: (reg.) zălă, zea s.f.] - Cf. ngr. záva.</p>	<pre> <entry type="homonym" id="ZA.2"> <hw>ZA</hw> <alt> <brack> <gram>nominativ_feminin_singular_indefinit</gram> <orth>ZA</orth> </brack> <brack> <gram>nominativ_feminin_plural_indefinit</gram> <orth>zale</orth> </brack> </alt> <pos>substantiv</pos> <gen>feminin</gen> <struc n="1"> <alt> <def>Fiecare dintre ochiurile unui lanț</def> <brack> <usg type="hyper">prin generalizare </usg> <usg>la pl.</usg> <def>lanț.</def> </brack> </alt> <struc type="Sec"> <def>Lănțișor de metal întrebuințat uneori ca podoabă. </def> </struc> <struc type="Sec"> <def>Cusătură în formă de lănțișor, executată de obicei la broderii.</def> </struc> </struc> <struc n="2"> <usg>La pl.</usg> <alt> <def>Împletitură executată din inele mici de fier legate unul de altul</def> <brack> <usg type="hyper">prin extensiune</usg> <def>armură făcută din această împletitură, cu care se îmbrăcau oștenii în antichitate și în evul mediu, spre a se apăra de loviturile dușmanilor.</def> </brack> </alt> </pre>

	<pre> </struc> <struc type="Varianta"> <alt> <brack> <orth> zală</orth> <stress> zălă</stress> <usg>reg.</usg> </brack> <orth> zea</orth> </alt> <pos>substantiv</pos> <gen>feminin</gen> </struc> <etym> Cf. <lang>ngr.</lang> záva. </etym> </entry> </pre>
--	---

Figura 2: Continut primar si codificarea echivalenta în XML (cf. CONCEDE.dtd)

În tabelul din Figura 2, sunt exemplificate reprezentarea tipografică (de tip DEX) și reprezentarea codificată în XML. Menționăm că reprezentarea tipografică din coloana stânga a Figurii 2 s-a obținut automat, folosind un convertor XML de format, proiectat astfel încât rezultatul generării (interpretarea marcajului XML) să fie cât mai apropiat de aspectul dicționarului tipărit. Structura de dicționar, definită mai jos, este suficient de generală pentru a permite implementarea diferitelor tipuri de dicționare. În fapt, DTD-ul CONCEDE a fost utilizat pentru codificarea a două dicționare bilingve: un dicționar Sloven-Englez și un dicționar Român-Francez.

Adnotarea XML fiind independentă atât de convențiile tipografice cât și de limba dicționarului, este posibilă căutarea multi-criterială a informației în unul, două sau mai multe dicționare explicative ale unor limbi diferite. De pildă, o căutare multi-criterială ar putea fi parafrazată astfel:

*Gaseste si afiseaza toate intrarile ce corespund **substantivelor feminine, de origine neo-greaca** si al caror cuvinte titlu încep cu secventa de litere ZA.* O astfel de căutare va avea ca rezultat tipărirea cel puțin a intrării corespunzătoare cuvântului titlu ZA²:

ZA², zale, s.f. **1.** Fiecare dintre ochiurile unui lant; *p. gener.* (la pl.) lant. ♦ Lantisor de metal întrebuințat uneori ca podoaba. ♦ Cusatura în forma de lantisor, executată de obicei la broderii. **2.** (La pl.) Împletitura executată din inele mici de fier legate unul de altul; *p. ext.* armura făcută din această împletitură, cu care se îmbracău ostentii în

antichitate si în evul mediu, spre a se apara de loviturile dusmanilor. [Var.: (reg.) *zála, zea* s.f.] - Cf. ngr. *záva*.

3.3. Alte dictionare, lexicoane; indexul interlingual

Unul dintre rezultatele proiectului Multext-East îl constituie un lexicon de forme ocurenta (LFO), cu peste 450.000 de intrari, care contine triplete de tipul < cuvânt, lema, cod_morfo-sintactic >. Acest lexicon va fi completat cu formele flexionare (generate automat) a lemelor din XML-LEX nereprezentate în LFO. Codificarea folosita este compatibila cu recomandarile Eagles (<http://www.ilc.pi.cnr.it/EAGLES/home.html>) pentru adnotarea morfo-sintactica si este documentata pe larg în (Tufis *et al.*, 1997).

O alta resursa lexicala esentiala a fost Dictionarul de Sinonime al Limbii Române – DSLR (Seche, Seche, 1997), care a fost transpus în forma electronica la Facultatea de Informatica a Universitatii "A.I.Cuza" din Iasi. Forma electronica a DSLR a fost convertita în format XML astfel încât aceeasi interfata ce a fost dezvoltata pentru XML-LEX functioneaza si cu XML-DSLR.

Din corpusurile paralele mentionate mai sus si folosind programul ce implementeaza metodologia noastra de extragere a echivalentilor de traducere (Tufis, Barbu, 2001a, 2001b, 2002) s-a construit un dictionar bilingv Român – Englez (de asemenea transpus în format XML). Acest lexicon bilingv a fost validat manual si îmbogătit cu noi intrari din diverse surse publice.

În sfârșit, o resursă extrem de valoroasă a fost și Indexul Interlingual al EuroWordNet, exportat în format XML cu editorul VisDic produs la Universitatea Masaryk din Brno (Pavelek, Pala, 2002).

3.4. Alegerea nucleului lexical

Vom da câteva definitii ale unor notiuni pe care le vom folosi în cele ce urmeaza.

Când ne plasam într-un context monolingv, vorbim despre *sensuri*, *înțeleșuri* si *sinșeturi*. Un cuvânt are unul sau mai multe *sensuri*. Un sens refera un *înteles*. În EuroWordNet sensurile unui cuvânt sunt numerotate în functie de frecventa lor, iar sensul unei leme este denotat adaugând numarul sensului la forma ortografica a acesteia. O multime de sensuri astfel specificate (ex. *action2*, *activity1*, *activiteness1*) care refera acelasi înteles este numit *sinșet* si constituie el însusi denotatia întelesului sensurilor din sinșet. Cu alte cuvinte, un sinșet reprezinta *lexicalizarea unui înteles* în contextul monolingv curent.

Daca abstractizam notiunea de *înteles*, definita ca mai sus, astfel încât sa nu mai facem referirea la un anumit context monolingv, vom vorbi despre *concepte* care sunt referite de *întelesurile* lexicalizate în diferitele limbi. Asadar, putem vorbi despre concepte care au sau nu realizare lingvistica într-o limba sau alta. Un concept este un construct cognitiv, independent de limba, care în EuroWordNet este totdeauna lexicalizat cel puțin într-una dintre limbi. Un concept este mai departe rafinat în termeni de distinctii semantice

elementare (trasaturi semantice), deci putem vorbi despre gruparea conceptelor în functie de trasaturile lor semantice.

În EuroWordNet si deci si în BALKANET, ILI este definit ca o colectie nestructurata de intrari de forma: <ILI-index><descriere ontologica><glosa> {domeniu}. Indexul interlingual initial a fost construit plecând de la versiunea 1.5 a Wordnet-ului si deci glosele pentru fiecare concept au fost importate direct din sinsetul englezesc care se refera la întelesul conceptualizat în ILI.

Pentru a facilita o cât mai buna intercorelare a wordneturilor monolingve din cadrul proiectului si pentru a înlesni extensia lor ulterioara, consortiul proiectului a decis ca procesul implementarilor paralele sa fie centrat pe concepte (independente de limba) selectate de comun acord, la momente succesive de timp.

O prima selectie a constituit-o multimea asa-numitelor „concepte de baza” definite în EuroWordNet ca fiind acele concepte din ILI lexicalizate în limba engleza (în WORDNET) prin sinseturi plasate pe un nivel ierarhic cât mai sus si, în plus, care au un numar mare de hiponimi directi (tot în WORDNET). Ratiunea acestei decizii a constat în faptul ca, aceste concepte fiind foarte generale si totodata productive în definirea unor concepte mai particulare, este foarte probabil ca ele sa fie lexicalizate în majoritatea limbilor de interes. Acest lucru a fost probat atât în EuroWordNet cât si în BALKANET. Multimea *conceptelor de baza* (o motivatie mai detaliata a selectiei lor este prezentata in (Vossen, 1998) în raport cu obiectivele EuroWordNet) contine 1.310 concepte, fiecareia dintre ele fiindu-i atasata o glosa explicativa si o *descriere ontologica* (vezi Rodriguez et al., 1998).

Dupa implementarea, în toate cele 5 limbi ale proiectului, a nucleelor de ontologii lexicele corespunzând conceptelor de baza, s-a facut o noua selectie, de data aceasta, continând 4.000 de noi concepte interlinguale.

Selectia a avut în vedere, pe de o parte maximizarea compatibilitatii cu EuroWordNet, iar pe de alta parte relevanta stocului lexical pentru fiecare limba din perspectiva monolingva. Primul criteriu a fost operationalizat alegându-se acele concepte lexicalizate în cele mai multe limbi din EuroWordNet. Limita inferioara a numarului de limbi a fost fixata la 5, astfel încât dupa implementarea acestor concepte în BALKANET ele sa fie lexicalizate în cel putin 10 limbi.

Criteriul relevantei monolingve a condus la propunerea mai multor multimi candidate de concepte. Pentru fiecare limba a proiectului au fost efectuate analize cantitative în context strict monolingv. Metodele de analiza au diferit de la partener la partener, în raport cu datele si instrumentele disponibile pentru limbile în cauza. Dupa analiza acestor multimi, au fost incluse în multimea finala acele concepte ce au aparut în cel putin doua propuneri. Multimea finala a conceptelor a fost ordonata dupa numarul de limbi din EuroWordNet ce le lexicalizeaza si dupa numarul de limbi din BALKANET care le-au propus. Primele 4000 de noi concepte în aceasta lista au fost de comun acord alese ca tinta comuna pentru cea de a doua etapa a proiectului.

În continuare prezentăm metodologia folosită pentru limba română privind selecția fondului lexical în cadrul BALKANET. Analiza cantitativă s-a efectuat asupra unui corpus foarte mare, format din mai multe romane și dintr-o colecție de texte jurnalistice culese de pe web. Corpusul (conținând mai mult de 100 de milioane de cuvinte) a fost supus unor prelucrări statistice, fiind etichetat și lematizat automat, iar cuvintele care prezentau interes (substantive comune, verbe, adjective și adverbe) au fost sortate în funcție de frecvența lor în texte. Am extras în acest fel o listă de mai mult de 30.000 de leme. În funcție de frecvența acestora în textele din corpus, această listă a fost împărțită în trei părți, corespunzând celor mai frecvente 10.000 de leme (I), următoarele cele mai frecvente 10.000 (II) și restul (III). Frecvența dintr-un corpus este considerată de mulți lexicografi un criteriu subiectiv. Printre cele mai puternice argumente se numără volumul și reprezentativitatea textelor incluse în corpusul folosit la analiza cantitativă. Luând în calcul faptul că din ce în ce mai multe texte sunt disponibile pe web, mărimea corpusului nu mai reprezintă o problemă semnificativă, însă reprezentativitatea rămâne în continuare un punct slab. Definirea exactă a naturii textelor care trebuie incluse într-o analiză cantitativă face obiectul unei îndelungi polemici și nu vom insista asupra ei. Având în vedere că datele noastre constau aproape în întregime din texte jurnalistice, problema reprezentativității poate fi cu îndreptărire ridicată. Dictionarul de Frecvențe al Cuvintelor Românești FDRW (Juilland *et al.*, 1965), publicat cu mult timp în urmă, bazat pe un *corpus balansat* de 500.000 de cuvinte (teatru, nuvele și scurte povestiri, eseuri memorii și corespondențe, texte jurnalistice, literatură tehnică) conține cele mai frecvente 5.000 de leme. Chiar dacă este foarte controversat, FDLW este încă folosit de mulți lingviști români ca o referință. Comparatia pe care am făcut-o a arătat că mai toate cele 5.000 inventariate de FDRW se găsesc și în lista obținută de noi, chiar dacă nu cu aceleași scoruri de frecvență. Pe lângă frecvența în corpus am apelat și la alte două criterii mai puțin controversate și care au putut fi operationalizate în raport cu resursele lingvistice disponibile și instrumentele noastre de analiză a corpusurilor. Primul este numărul de sensuri pe care un cuvânt (împreună cu sintagmele și expresiile în care participă) îl are într-un dictionar. Al doilea este numărul de definiții de dictionar în care apare un anumit cuvânt. Al treilea criteriu, ne-inclus încă în analiză, ar putea fi numărul de derivate lexicale ale unui cuvânt. Pentru o pertinentă analiză din acest punct de vedere, o excelentă lucrare este (Dinu, 1996).

În această fază a proiectului BALKANET, ne-am concentrat atenția asupra substantivelor din limba română, iar datele experimentale raportate mai jos se referă doar la acestea. Având însă în vedere că procedurile tehnice nu depind de categoria gramaticală, metodologia și procedura vor fi aceleași și pentru verbe, adjective și adverbe. Luând în calcul numai primele două clase de frecvență descrise mai sus (primele 20.000 cele mai frecvente din corpusul jurnalistice) am extras din XML-LEX mai mult de 8.000 de intrări de substantive și substantive compuse (care însumează aproximativ 35.000 de sensuri) astfel încât productivitatea definitională PD (numărul de definiții în care participă un substantiv) să fie cel puțin 3. Lista a fost sortată în funcție de productivitatea definitională și numărul de sensuri ale fiecărui cuvânt titlu.

Substantiv	Productivitate definițională	Număr de sensuri	FRECV _{range}
acțiune	2279	13	I
persoană	1979	9	I
parte	1882	94	I
formă	1286	21	I
obiect	1204	16	I
fapt	1044	11	I
...
rasism	3	1	II

Figura 3: Ordonarea candidatilor

Pentru toate aceste substantive am extras traduceri englezesti din dictionarul de echivalenti de traducere. Procedurile pentru extragerea automata a echivalentilor de traducere din corpusuri paralele ca si procedura de discriminare a sensurilor sunt descrise pe larg în (Tufis, Barbu, 2001a,b), (Erjavec et al. 2001), (Tufis, 2002), (Ide *et al.*, 2002). Fiecare substantiv din limba româna a fost pus în corespondenta cu lista tuturor conceptelor din ILI corespunzatoare traducerilor sale în engleza. Conceptele astfel identificate, au fost sortate dupa rangul corelat al substantivelor românești de la care s-a pornit.

Interesant de remarcat ca dintre cele 4000 de concepte selectate în final prin armonizarea propunerilor tuturor partenerilor, circa 2600 s-au regasit si în primele 4000 de concepte ale ierarhiei noastre. Toate cele 4000 de concepte selectate de consortiu se regasesc printre primele 6000 de concepte ale ierarhiei noastre.

Toate substantivele reprezentând potientiale lexicalizari ale celor 4000 de concepte din cea de a doua selectie au fost automat puse în corespondenta cu toate definitiile lor din XML-LEX. De asemenea, ele au fost corelate cu lexicalizarile din limba engleza ale celor 4.000 de concepte. Prin intermediul dictionarului de echivalenti de traducere englez-român, fiecare concept a fost asociat cu lexicalizarea din limba engleza (extrasa din WORDNET) si cu potientialele lexicalizari în limba româna.

Dictionarul de Sinonime al Limbii Române (DSLRL), digitizat si codificat în XML, a fost folosit pentru a extrage seriile sinonimice pentru cuvintele românești selectate. În XML-DSLRL unii membri ai seriilor sinonimice sunt arhaisme sau regionalisme. Discutiile preliminare au condus catre ideea de a elimina toate cuvintele care fac parte din aceste clase (ne-am bazat pe cerinta de a construi un nucleu lexical de uz general în limba româna contemporana). Totusi, pentru eventualitatea în care aceste cuvinte filtrate (împreuna cu informatiile despre uz) vor fi necesare mai târziu, s-a asigurat recuperabilitatea lor. Seriile sinonimice românești au fost considerate ca posibile sinseturi si adaugate la asociatiile descrise mai sus.

4. Instrumente software dezvoltate pentru proiectul BALKANET

Materialul lingvistic de baza descris în secțiunea anterioară, a fost asamblat prin intermediul unor programe unitare, astfel încât toată această informație este disponibilă într-o interfață „prietenoasă”, prin care lexicograful alege echivalentele corecte de sens dintre cele potențiale. Această interfață este generată și „personalizată” automat în funcție de mulțimea conceptelor interlinguale furnizată ca parametru de intrare unui generator de interfețe. Printr-un astfel de model arhitectural, a fost posibil ca sarcina construirii wordnet-ului pentru limba română să fie distribuită între membrii celor două colective românești participante la proiect și judicios controlată. Pentru fiecare dintre aceștia s-a generat o interfață personalizată pentru o submulțime distinctă de concepte dintre cele agreate de consorțiul proiectului. Utilizatorul acestei interfețe, pe care generic îl numim în continuare *lexicograf*, va lucra în mod independent de ceilalți, construind, ca urmare a interacțiunii, fragmente ale wordnetului pentru limba română. La un moment dat, lexicograful alege un concept din mulțimea ce i-a fost repartizată caruia dorește să-i ataseze un sinset românesc. El are la dispoziție simultan, sinsetul ce lexicalizează în limba engleză conceptul respectiv și, pentru fiecare cuvânt englezesc din acest sinset, toate potențialele lui traduceri în limba română, aceste traduceri având atasate toate definițiile conținute în XML-LEX. În plus, fiecare cuvânt românesc are atasate toate seriile sinonimice din XML-DSLRL în care el este prezent. Ceea ce trebuie să decida lexicograful este (vezi figura 4):

- a. care este cuvântul românesc a cărui definiție este cea mai apropiată de definiția conceptului lexicalizat în limba engleză;
- b. care este cea mai bună serie sinonimică a acestui cuvânt;
- c. care dintre definițiile atasate cuvintelor dintr-o serie sinonimică este cea mai adecvată pentru a fi aplicabilă tuturor cuvintelor din seria respectivă.

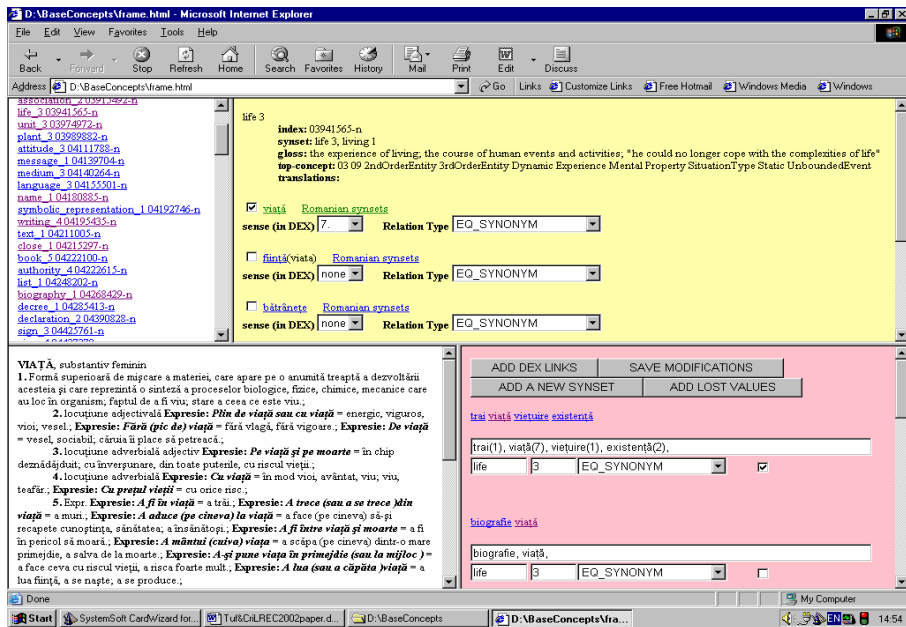


Figura 4: Editorul pentru construirea sinseturilor

În majoritatea cazurilor, definițiile extrase din XML-LEX corespunzând sinonimelor dintr-un sinset nu sunt identice, lexicograful alegând pe cea mai apropiată de definiția conceptului corespunzător (vezi figura 5).

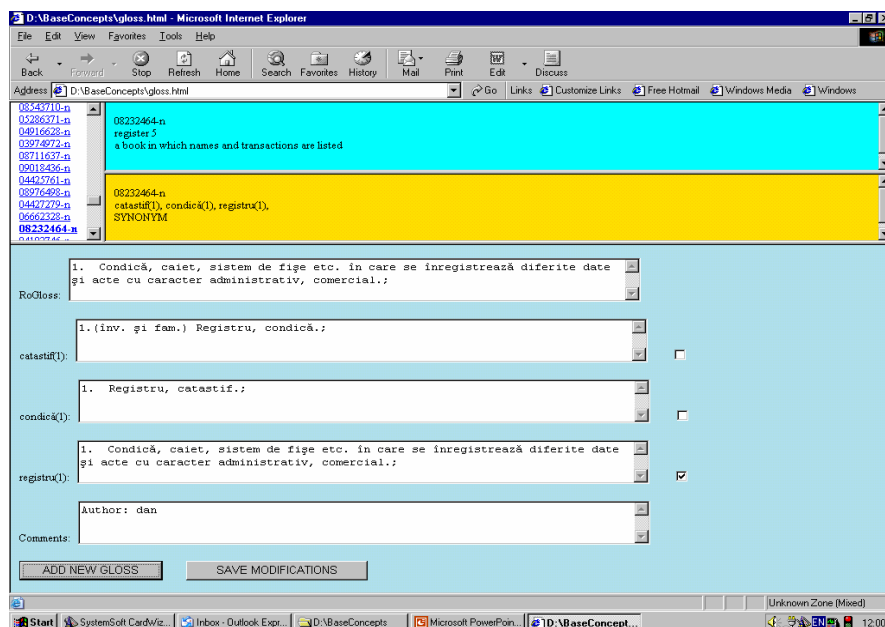


Figura 5: Editorul pentru asignarea gloselor

Merita mentionat ca în faza asocierii gloselor a devenit evidenta incorectitudinea alcatuirii unor sinseturi, ele fiind modificate. În alte cazuri Dictionarul Explicativ al Limbii Române include în aceeași definiție două sensuri care sunt demarcate în ILI ca două concepte diferite. În astfel de situații strategia generală a fost să se despartă definiția românească și să se ataseze ca glosa partea relevantă.

Fragmente create de fiecare lexicograf sunt agregate în mod incremental în structuri din ce în ce mai complexe și mai acoperitoare din punct de vedere lexical. Acest proces de agregare se realizează în mod centralizat, astfel încât corectitudinea structurilor rezultate să poată fi controlată și, în cazul conflictelor, să se poată identifica și corecta sursele de conflict (de exemplu: același sens pus în corespondență cu concepte diferite, sensuri diferite ale aceluiași cuvânt puse în corespondență cu același concept, literali fără identificatori de sens etc.). Corectarea unor conflicte între două porțiuni ale structurii agregate poate să genereze conflicte între alte părți ale sale. Pentru evitarea acestui pericol au fost proiectate mecanisme de control centralizat al unificării subseturilor de wordnet ce gestionează efectul global al oricăror modificări locale.

4.1. Importul relatiilor taxonomice; vizualizare sincronizata a mai multor wordneturi

Constructia sinseturilor si punerea lor în corespondenta cu conceptele interlinguale reprezinta doar una din cele doua dimensiuni fundamentale ale procesului de construire a unei retele semantice lexicale pusa în corespondenta cu indexul interlingual, respectiv cea de implementare a nodurilor si echivalarea acestora cu conceptele interlinguale. Cea de a doua dimensiune a procesului constructiei retelei o constituie definirea relatiilor (intralinguale) între nodurile create si echivalate în prima faza. Deosebit de importante sunt relatiile taxonomice care stabilesc o ierarhie de generic-specific între sinseturile unui wordnet.

Stabilirea relatiilor taxonomice între sinseturile wordnetului pentru limba româna s-a facut automat (urmata de validarea umana) în baza principiului „echivalentei ierarhice interlinguale” (Tufis, Cristea, 2002). În esenta, acest principiu afirma ca:

1. daca sinsetul S_{1LA} din limba LA si sinsetul S_{1LB} din limba LB sunt echivalate cu acelasi concept C_1 din ILI si
2. daca sinsetul S_{2LA} din limba LA si sinsetul S_{2LB} din limba LB sunt echivalate cu acelasi concept C_2 din ILI si
3. daca în limba A sinseturile S_{1LA} si S_{2LA} sunt într-o relatie ierarhica H^+ (H^+ denota compunerea de un numar de ori cel puțin egal cu 1 a relatiei H, în cazul nostru: has-as-hypernym), atunci:

în limba B sinseturile S_{1LB} si S_{2LB} sunt într-o relatie ierarhica similara H^+ (desi lanturile de relatii H pot fi de lungimi diferite în cele doua limbi).

Principiul explicita necesitatea ca interpretarea relatiilor folosite în ontologia multilingva sa fie similara, asadar defineste **coeziunea interpretativa** a relatiilor ontologice în toate limbile participante la proiect. Acest principiu este reprezentat schematic în figura 6:

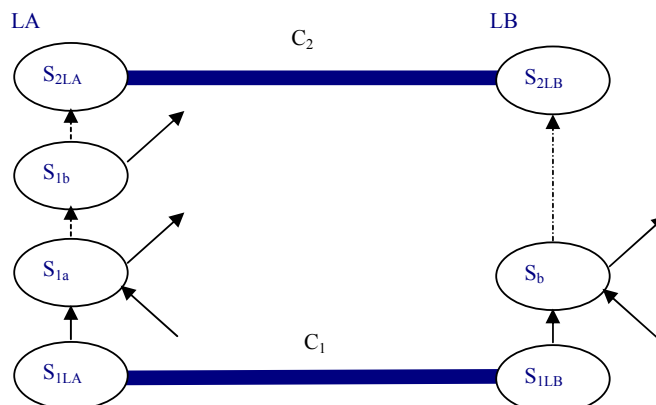


Figura 6: $(S_{1LA} \text{ EQ-SYN } S_{1LB}) \& (S_{2LA} \text{ EQ-SYN } S_{2LB}) \& (S_{1LA} H^+ S_{1LB}) \Rightarrow (S_{2LA} H^+ S_{2LB})$

În secțiunea următoare vom arata pe un caz concret cum poate fi exploatat acest principiu pentru a importa (și eventual valida/corecta manual) relațiile dintr-un wordnet în care structurile ierarhice au fost stabilite, într-un wordnet pentru care au fost stabilite doar relațiile de echivalență translatională cu indexul interlingual (ILI).

Ultima etapă a construirii unui grup de sinseturi este transformarea rezultatelor interacțiunii lexicografului cu interfața descrisă anterior într-un format independent de limbă (codificare XML) și specific editorului multilingual de ontologii lexicale numit VisDic (Pavelek și Pala, 2002). Odată generat acest format, el poate fi încărcat în VisDic, iar wordnetul pentru limba română poate fi vizualizat în mod sincron cu toate celelalte wordneturi încărcate. În figura de mai jos este ilustrată afișarea în mod sincron a sinsetului românesc (*ființă_1*, *forma de viață_1*, *vietuitoare_1*, *vietate_1*) și a celui englezesc (*being_1*, *life form_1*, *living thing_1*, *organism_1*) și a arborilor lor de hiponimi. Cele două sinseturi sunt aliniate via ILI, ambele fiind echivalate independent cu conceptul interlingual cu identificatorul 00002728-n.

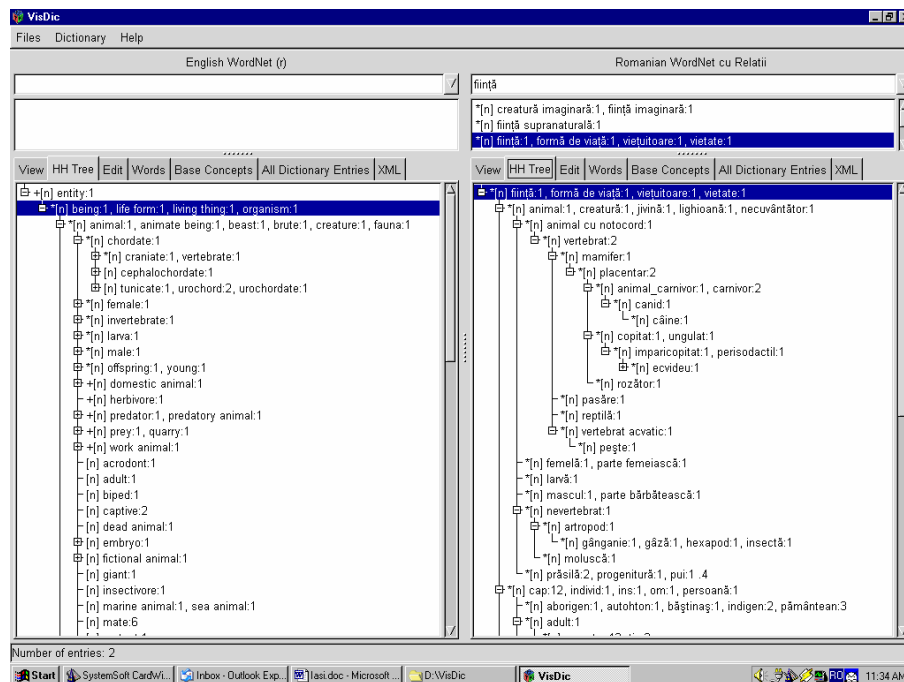


Figura 7: Vizualizarea sincronizată a două ontologii lexicale cu ajutorul VisDic

Editorul de ontologii multilingve, VisDic, a fost dezvoltat în cadrul proiectului BALKANET pentru a substitui funcționalitatea asigurată în cadrul EuroWordNet de

editorul Polaris, dezvoltat de firma Lernout & Hauspie. Implementat initial pentru ca rezultatele proiectului BALKANET sa poata fi utilizate în regim liber de restrictii comerciale (Polaris poate fi utilizat doar contra cost), VisDic este constant îmbunatatit cu facilitati noi a caror necesitate apare pe masura evolutiei proiectului BALKANET, fiind deja unul dintre cele mai puternice instrumente existente pentru gestiunea ontologiilor multilinguale.

5. Principiul conservarii trans-linguale a ierarhiei lexicale. Studiu de caz: Condimente, mirodenii, sosuri si alte ingrediente

Vom considera fragmentele din RO-WordNet si WordNet 1.5 aratate în figura 8. Sagetile reprezinta relatiile de hiponimie (de la hiponime spre hipernime) în cele doua wordneturi. Liniile groase reprezinta relatiile de echivalenta de traducere (EQ-SYN) dintre sinseturile celor doua limbi, aceasta însemnând ca sinseturile respective sunt puse în corespondenta cu acelasi concept din ILI. Linia groasa întrerupta reprezinta o relatie EQ-SYN identificata ca nerespectând principiul conservarii trans-linguale a ierarhiilor lexicale din cele doua wordneturi. Inconsistenta este semnalata deoarece în româna relatiile ierarhice (de hiponimie) dintre *mirodenie*(RO) si *condiment*(RO) ca si dintre *ketchup*(RO) si *sos*(RO) nu sunt verificate de echivalentii lor în limba engleza: *spice*(EN) este frate cu *condiment*(EN) si respectiv *ketchup*(EN) este frate cu *sauce*(EN). Daca structura variantei 1.5 a WordNet este considerata cea corecta, acest exemplu arata ca principiul pastrarii ierarhiei nu este irefutabil. Pe de alta parte, daca ar fi rezonabil sa consideram ca WN 1.5 este amendabil (de exemplu facând *mustard*(EN) si *ketchup*(EN) hiponimii directi ai lui *sauce*(en)) ca în figura 9, atunci principiul pastrarii ierarhiei ar putea fi o puternica proba a consistentei⁷⁶.

În urma restructurarilor ierarhice si de echivalare translationala, necesare pentru respectarea principiului conservarii trans-linguale a ierarhiei lexicale (aratate în figura 9), interesant este faptul ca a disparut relatia de echivalenta între cuvântul românesc *condiment* si cuvântul englezesc *condiment*.

⁷⁶ Consultată recent asupra acestei probleme, Christiane Felbaum a confirmat esistența unei erori în ierarhia WN1.5, probată, de altfel, și de glosa lui ketchup (thick spicy **sauce** made from tomatoes).

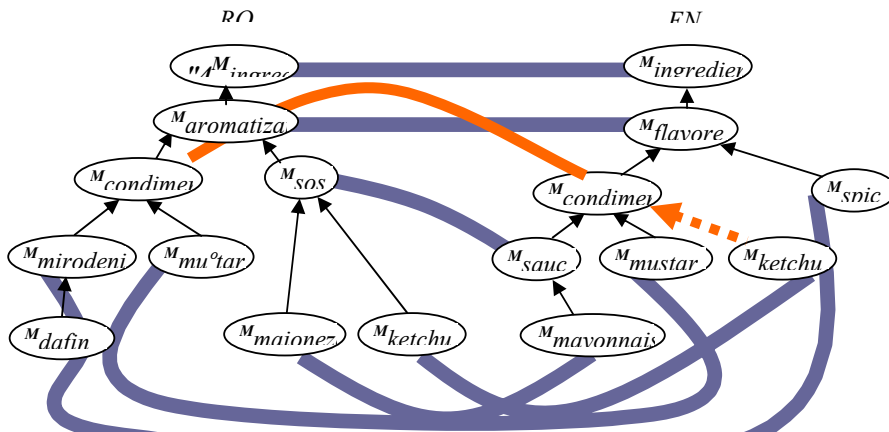


Figura 8: Nerespectarea principiului conservării trans-linguale a ierarhiei lexicale

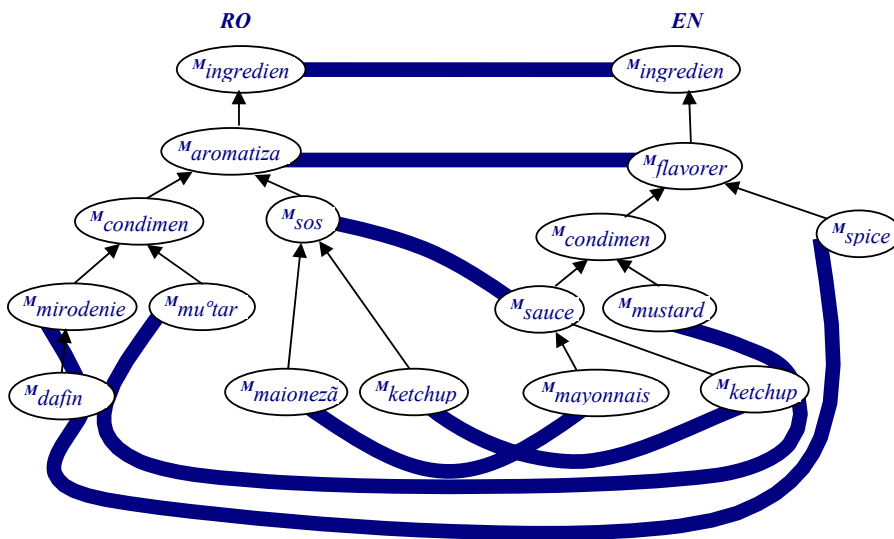


Figura 9: Nerespectarea principiului conservării trans-linguale a ierarhiei lexicale

Pentru ca aceasta echivalenta sa fie posibila, în condițiile principiului conservării trans-linguale a ierarhiei lexicale, ar trebui ori ca în limba engleza *spice* **sa fie** un hiponim al lui *condiment* iar *sauce* **sa nu fie** un hiponim al lui *condiment* ci frate, ori în limba româna *sos* **sa fie** un hiponim al lui *condiment* iar *mirodenie* **sa nu fie** un hiponim al lui *condiment* ci frate. Ambele variante au fost respinse de expertii consultati, lexicografi si vorbitori nativi ai limbii engleze si respectiv române. Singura concluzie posibila este ca în româna si engleza cuvântul *condiment* nu reprezinta exact acelasi lucru.

5. Concluzii

Realizarea ontologiei lexicale pentru limba româna, în contextul multilingual definit de proiecte de tipul EuroWordNet, Balkanet si GlobalWordnet (www.globalwordnet.org), este esentiala pentru procesul de informatizare a limbii române. Experienta internationala arata ca un astfel de proiect nu este niciodata închis, reclamând actualizare si întreținere continua, aparând mereu noi idei de îmbunătățire a performanțelor si noi cerinte de exploatare. Specialistii de la Princeton au anuntat deja versiunea 1.7.1 a Wordnet, mult îmbunătățita. În variantele ce vor urma, pe lângă extensia în continuare a fondului lexical, toate cuvintele nefunctionale aparând în definitii vor contine referinte spre sinsetul corespunzator contextului de utilizare. Cu alte cuvinte, Wordnet va deveni simultan si un dictionar si un corpus adnotat la nivelul sensului. O alta dezvoltare semnificativa o va reprezenta traducerea definitiilor din Wordnet în formule logice, adecvate prelucrarilor inferentiale. Acest proiect, coordonat de Dan Moldovan si Sanda Harabagiu se afla în derulare la Universitatea Texas din Dallas (Moldovan, 2001), (Harabagiu *et.al.*, 1999).

Astfel de extensii vor trebui considerate în viitor si în wordnetul pentru limba româna aflat deocamdata în faza incipienta. Obiectivul final prevazut pentru cei trei ani de derulare ai proiectului BALKANET (septembrie 2004) este realizarea unui nucleu de câte 8.000 de sinseturi în fiecare din limbile proiectului.

În acest moment, la mai puțin de un an de la începerea proiectului, wordnetul românesc se afla cu mult înaintea graficului prevazut, având deja create peste 6.000 de sinseturi. Se poate estima ca, în conditii normale, în cei peste doi ani care au mai ramas wordnetul românesc va ajunge la peste 20.000 de sinseturi, acoperind peste 40.000 de literali. Atingerea unui volum lexical similar cu al altor wordneturi necesita însa continuarea proiectului si dupa anul 2004, atragerea unor noi colective de specialisti în aceasta întreținere si desigur gasirea surselor de finantare, în principal interne, care sa permita dezvoltarea si întreținerea wordnetului românesc. Operationalizarea acestui obiectiv poate fi facilitata de contextul organizatoric creat de curând prin înființarea la Academia Româna a Comisiei de Informatizare pentru Limba Româna (CILR) precum si a Consorțiului de Informatizare pentru Limba Româna (ConsILR: <http://www.consilr.info.uaic.ro/>), for executiv al CILR.

A fost construită o platformă software de dezvoltare incrementală a rețelei semantice ce permite implementarea independentă de regiuni ale rețelei și integrarea ulterioară a acestora. Viabilitatea acestui concept arhitectural și a demersului de dezvoltare distribuită a wordnetului au fost validate prin implicarea în procesul de construire a 10 specialiști, cărora li s-au adăugat încă 12 studenți masteranzi de la Facultatea de Litere a Universității București și Facultatea de Informatică a Universității "A.I. Cuza" (cele două facultăți ce au programe de Master în domeniul prelucrării limbajului natural și al lingvisticii computaționale). Rezultatele produse în mod independent au fost agregate fără nici o dificultate. Mediul lingvare de dezvoltare conține un modul special de verificare a corectitudinii deciziilor lingvistice la crearea sinseturilor românești sau la punerea lor în corespondență cu conceptele indexului interlingual. După cum era de așteptat, procesul de integrare a rezultatelor parțiale furnizate de fiecare membru al celor două echipe de realizare a evidențiat o serie de inconsistente cu explicații diverse:

- neatenție în asignarea sensurilor, generată de oboseala expertului decident uman;
- granularitate semantică diferită între sensurile explicitate în XML-LEX și sensurile conceptelor din ILI;
- absența lexicalizării în limba română a unor concepte existente în ILI și introducerea unor forme perifrastice cu definiții ad-hoc;
- erori sau incompletitudini existente în sursele lingvistice primare folosite în implementare.

Inconsistențele depistate, atât de natură structurală, dar mai ales cele de natură semantică au fost înregistrate, analizate și unele dintre ele corectate. Altele, necesită o analiză mai profundă și rezolvarea lor a fost amânată pentru o etapă ulterioară a proiectului. Aceasta cu atât mai mult cu cât, prin analiză similară pe care am efectuat-o asupra wordneturilor pentru celelalte limbi din proiect, am constatat că există multe similități ale acestor genuri de inconsistente. Sunt puse astfel în evidență o serie de concepte din ILI pentru care diferența semantică dintre ele este prea mică pentru a fi sesizată ușor chiar și de către un vorbitor nativ al limbii respective. Distincții atât de rafinate au, din perspectiva prelucrării automate și mai ales al traducerii automate, o utilitate limitată iar în context multilingv pot fi chiar surse de eroare. Pericolul micșorării distanței semantice (am putea numi acest fenomen pulverizarea conceptuală) între conceptele din ILI este amplificat de adăugarea unor concepte ce au lexicalizări într-o singură limbă sau într-un număr mic de limbi. O soluție pentru evitarea idiosincraziilor lexicale într-un context multilingv și a disparităților de traducere este gruparea conceptelor foarte apropiate semantic în ceea ce s-ar putea numi *concepte agregate*. Lexicalizările înțelesurilor din două sau mai multe limbi, puse în corespondență cu aceleași concepte din ILI sau cu concepte membre ale unui agregat, vor putea fi folosite ca echivalenți de traducere în pofida unor diferențieri semantice specifice unei limbi sau alteia (*ciorba, sarmale, pepper pot, porcupine ball* etc.; vezi și exemplele din secțiunea precedentă). Analiza inconsistentelor interumane în echivalarea înțelesurilor dintr-o limbă cu conceptele interlinguale din ILI, precum și

identificarea conceptelor distincte puse în corespondență cu echivalenți de traducere (extrasi automat din corpusuri paralele sau gasiti într-un dictionar bilingv clasic) pot furniza informatii calitative mult mai interesante (cel puțin din perspectiva psiholingvisticii) și mai demne de încredere decât o analiza statistica. Aceasta este o promitatoare directie de cercetare ce se dezvoltă în paralel cu activitatea principala de constructie a wordnetului pentru limba română.

Referinte bibliografice

- Bloksma, L., Diez-Orzas and Vossen, P. (1996) The User Requirements and Functional Specification of the EuroWordNet-project *EWN-deliverable D.001*, LE-4003
- Danzin, A. (1992) „Towards a European Language Infrastructure” raport al Comisiei Europene
- Dinu, M. (1996). Personalitatea limbii române, Editura ALL, 368 p.
- DEX (1996). Coteanu, I., Seche, L., Seche, M. (coord.). Dicționarul Explicativ al Limbii Române, Ediția a II-a, *Univers Enciclopedic*, București
- Erjavec, T., Ide, N., Tufiş, D.(1997) Encoding and Parallel Alignment of Linguistic Corpora in Six Central and Eastern European Languages” in Michael Levison (ed) *Proceedings of the Joint ACH/ALL Conference* Queen's University, Kingston, Ontario (also on <http://www.qcis.queensu.ca/achallc97>)
- Erjavec, T., Evans, R., Ide, N., Kilgarriff, A. (2000). The CONCEDE Model for Lexical Databases. *Proceedings of the Second Language Resources and Evaluation Conference (LREC)*, Athens, Greece, 355-362.
- Erjavec, T., Ide, N., Tufiş, D.(2001) *Automatic Sense Tagging Using Parallel Corpora*, in Proceedings of the 6th Natural Language Processing Pacific Rim Symposium, Tokyo, Japan, 27-29 November, pp. 212-219, 2001
- Fellbaum, Ch. (ed.) *WordNet: An Electronic Lexical Database*, MIT Press, 423 p.
- Harabagiu, S., Miller, G., Moldovan, D. (1999). „WordNet 2 - A Morphologically and Semantically Enhanced Resource”, in *Proceedings of SIGLEX-99*, Univ. of Maryland, pp 1-8.
- Ide, N. (1998) *Corpus Encoding Standard: SGML Guidelines for Encoding Linguistic Corpora* First International Language Resources and Evaluation Conference, Granada, Spain. See also <http://www.cs.vassar.edu/CES/>.
- Ide, N., Erjavec, T., Tufiş, D. (2002): „Sense Discrimination with Parallel Corpora” in Proceedings of the SIGLEX Workshop on Word Sense Disambiguation: Recent Successes and Future Directions. ACL2002, Philadelphia, pp. 54-60.
- Juilland, A., Edwards, P.M.G, Juilland, I. (1965). The Frequency Dictionary of Rumanian Words. *Mouton & CO.*, London-The Hague-Paris

- Miller, G.A., Beckwith, R., Fellbaum, C., Gross, D., Miller, K.J. (1990) "Introduction to WordNet: An On-Line Lexical Database" 1990 In International Journal of Lexicography, Vol. 3, No. 4 (winter), pp. 235-244
- Moldovan, D. (2001). "Question Answering Systems in Knowledge Management", IEEE Intelligent Systems, vol 16, nr. 6, pp 90 – 92.
- Pavelek, T., Pala, K. (2002) *VisDic: A new Tool for WordNet Editing* in Proceedings of the 1st International Wordnet Conference, Mysore
- Rodriguez, H., Climent, S., Vossen, P., Bloksma, L., Peters, W., Alonge, A., Bertagna, F., Roventini, A.(1998) The Top-Down Strategy for Building EuroWordNet: Vocabulary Coverage, Base Concepts and Top Ontology. In Piek Vossen (ed.) *EuroWordNet: A Multilingual database with lexical semantic networks*, Computers and Humanities, Vol. 32, Nos. 2-3
- Seche, L., Seche, M.(1997) *Dicționarul de sinonime al limbii române*. Univers Enciclopedic, București
- Stamou, S., Oflazer, K., Pala, K., Christoudoulakis, D., Cristea, D., Tufiş, D., Koeva S., Totkov, G., Dutoit, D., Grigoriadou, M. (1997) BALKANET A Multilingual Semantic Network for the Balkan Languages, in *Proceedings of the International Wordnet Conference*, Mysore, India
- Tufiş, D., Bruda, Şt. (1997). Structure Markup in CES and Preliminary Statistics on Romanian Translation of Plato's "Republica", *Proceedings of International Seminar on Encoding*, Ljubljana, February, 1997, also in *TELRI News*, nr. 5
- Tufiş, D. (1999). Tiered Tagging and Combined Classifiers In F. Jelinek, E. Nöth (eds) *Text, Speech and Dialogue*, Lecture Notes in Artificial Intelligence 1692, Springer
- Tufiş, D., Barbu, A.M., Pătraşcu, V., Rotariu, G., Popescu, C. (1997). Corpora and Corpus-Based Morpho-Lexical Processing, in Tufiş D., P. Andersen (eds.) *Recent Advances in Romanian Language Technology*, Editura Academiei
- Tufiş, D., Rotariu, G., Barbu, A.M. (1999) TEI-Encoding of a Core Explanatory Dictionary of Romanian. In Kiefer, F. and Pajzs J. (eds.) *Papers in Computational Lexicography*, Hungarian Academy of Sciences pp. 219-228
- Tufiş, D., Popescu, C., Roşu, R (2000).: Automatic classification of documents by random sampling in *Proceeding of the Romanian Academy*, Series A, vol 1, no. 2, pp. 18-28
- Tufiş, D. (2001), "Promovarea Limbii Române în SI-SC", in "Societatea Informationala - Societatea cunoaşterii: concepte, solutii si strategii pentru România", Florin Gh. Filip (coord.), Ed. Expert, Bucuresti, 2001, ISBN 973-8177-42-1, pp. 131-142
- Tufiş, D., Barbu, A.M.(2001a) *Computational Bilingual Lexicography: Automatic Extraction of Translation Dictionaries*, in International Journal on Science and Technology of Information, Romanian Academy, ISSN 1453-8245, Vol.4, No.3-4, 2001, pp.325-352

- Tufiş, D., Barbu, A.M.(2001b) *Extracting multilingual lexicons from parallel corpora*, in Proceedings of the ACH-ALLC conference, New York, 12-17 June, 4p.
- Tufiş, D., Barbu, A.M(2002). "Lexical token alignment: experiments, results and applications" In Proceedings of LREC2002, Las Palmas, Spain, pp.458-465
- Tufiş, D. (2000). Blurring the distinction between machine readable dictionaries and lexical databases. *Research Report, RACAI-RR56, p. 56*
- Tufiş, D. (2002). "A cheap and fast way to build useful translation lexicons" in *Proceedings of the 19th International Conference on Computational Linguistics, COLING2002*, Taipei, 25-30 August, 2002, 7 p
- Tufiş, D., Cristea, D. (2002). „Methodological issues in building the Romanian Wordnet and consistency checks in Balkanet”, In Proceedings of LREC2002, Las Palmas, Spain, May, pp. 35-41
- Vintilă-Radulescu, I. (2002). "Resurse lingvistice pentru limba română elaborate la Institutul de Lingvistică «Iorgu Iordan»", în acest volum.
- Vossen, P. (ed.) (1998). "A Multilingual Database with Lexical Semantic Networks", Kluwer Academic Publishers, Dordrecht

Algoritmi de segmentare a textului în unitati de tip clauzal

D. GÂLEA,
Institutul de Informatica Teoretica, Academia Româna, Filiala Iasi,
dgalea@iit.tuiasi.ro

N. CURTEANU,
Institutul de Informatica Teoretica, Academia Româna, Filiala Iasi,
curteanu@iit.tuiasi.ro

C. LINTES,
Institutul de Informatica Teoretica, Academia Româna, Filiala Iasi

1. Introducere

Scopul lucrării este dublu: (a) Sa prezinte si sa compare doi algoritmi de segmentare a frazei (românești) în unitati de tip clauzal. (b) Sa întregiasca si sa sustina doua componente de baza ale strategiei lingvistice SCD (Segmentare-Coeziune-Dependentă) [1], [2] de analiza a limbajului natural (LN): procesul de segmentare a textului de LN, si teoria FX-bar [3]. Segmentarea textului poate continua sau interfera cu stabilirea arborilor de dependenta între unitatile clauzale si subclauzale (unitati sintagmatice) ale textului. Unitatile de tip-clauzal corespund, în general, relatiilor retorice dintre unitatile minimale ale discursului, astfel încât algoritmi de segmentare pot fi (si chiar sunt) utilizati în aplicatii ce tin de teoria si procesarea discursului. Primul algoritm este o aplicare la limba româna a segmentarii frazei în unitati de tip-clauzal, algoritm dezvoltat de Daniel Marcu în [4], [5] (si prescurtat în cele se urmeaza “algoritm-Marcu”, sau “algoritm-M”). Al doilea algoritm reprezinta o rafinare a segmentarii în clauze si grupuri sintagmatice din cadrul strategiei lingvistice SCD [1], [2], [3], [6] (prescurtat în cele se urmeaza prin “algoritm SCD”). Acesti algoritmi sunt implementati într-un mediu specializat de procesare (dezvoltat sub C++), si este realizata o comparatie computationala a executiei segmentarii de tip-clauzal pe un set consistent de fraze românești [7].

Segmentarea textului LN a devenit în ultimii ani un subiect intensiv cercetat si cu multiple aplicatii. O atentie speciala a primit segmentarea textului de LN în unitati de discurs, în particular, segmentarea frazei în unitati minimale de discurs, de multe ori si pe buna dreptate asociate cu unitati de tip clauzal în numeroase teorii sintactice, semantice, si de discurs. Unitatile textuale de tip-clauzal obtinute (sau proiectate) prin mijloace orientate

sintactic, pe analiza de suprafață [Eng: *shallow*], s-au dovedit a fi esențiale în numeroase tipuri de procesare a LN: parsare, traducere automată, generare de LN, interpretare de discurs, extragere de date lingvistice, regăsirea informației, rezumare automată, rezoluția anaforei, etc. Un caz special de segmentare a textului este ceea ce numim ‘*chunking*’, un proces dedicat obținerii unor tipuri de “*segmente*” [Eng: *chunks*] dominate de anumite categorii (verb, substantiv, adjectiv-adverb, clauză). În continuare, vom folosi doar termenul de *segmentare* a textului de LN, considerând *chunking*-ul drept un caz particular al procesului de segmentare a LN.

În analiza și implementarea celor doi algoritmi de segmentare, *algoritmul-M* și *algoritmul-SCD*, cel puțin două aspecte le considerăm a fi importante: **(a)** Se demonstrează că *algoritmul-M* de segmentare este scufundat în *algoritmul-SCD*, ceea ce înseamnă că primul dintre cei doi algoritmi poate fi obținut ca un caz particular al claselor de marcheri, ierarhiei acestor clase, și a segmentării (dependentelor) obținute de cel de-al doilea algoritm. **(b)** *Algoritmul-SCD* de segmentare poate fi conceput ca un bun punct de start în proiectarea unui cadru general pentru *algoritmii de segmentare* a textului de LN. Un asemenea cadru ar fi compus din: (b1) mai multe *sisteme de transformare* aplicate în cascada, fiecare sistem component fiind format din *seturi și subseturi specifice de etichete*, (b2) o *ierarhie* stabilită între câteva dintre cele mai importante *clase* ale acestor *etichete*, și (b3) o *gramatică formală* (sau un automat finit) pentru recunoașterea (sub)secvențelor și arborilor de etichete (în concordanță cu ierarhia claselor de etichete). În abordarea prezentă, aceste aspecte sunt exemplificate de către o implementare C++ a celor doi algoritmi într-un mediu specializat, o bază de date a marcherilor (de discurs) românești, și o ierarhie specifică a claselor de marcheri lingvistici. Cei doi algoritmi de segmentare considerați sunt executați și comparați pentru un set consistent de fraze românești. Posibile dezvoltări și aplicații sunt menționate [21], [22].

Importanța segmentării de tip clauzal a frazei în procesul de parsare a textului a fost scoasă în evidență încă de la începutul anilor '80, iar studiile teoretice datează mult mai devreme. În România, primele lucrări științifice și contracte de lingvistică computațională au continuat, printre alte realizări meritorii, și primele încercări de realizare a segmentării automate a frazei în clauze finite (și infinite) [8], [9], [10], [11]. În pofida unor modeste mijloace formale (gramatici formale) și de programare (rețele ATN) disponibile la acel timp, ideile principale pe care se bazează abordările menționate, nu numai că au reprezentat premiere pentru acele timpuri, dar multe din ideile de atunci își păstrează încă o surprinzătoare actualitate, aceste fenomene de *come-back* ciclic fiind frecvente (și perfect explicabile de evoluția tehnologică) în momentul de față. Trebuie menționate aici folosirea intensiva a *marcherilor de discurs* (*cue phrases, connectives*), întâlnită și în [4], [5], [12], [13], a *predicativității* (aparitia categoriilor ‘*deverbale*’) [13], [14], a utilizării automatelor finite în analiza LN, etc.

De fapt, o versiune a gramaticii formale preluată din [8] este folosită în *Pasul 6* al *algoritmului SCD-2002* de segmentare, în concatenarea marcherilor de nivel M3 și M2 (vezi Secțiunea 4), în timp ce rudimente ale unor reguli similare din aceeași gramatică se

regasesc în algoritmul-*M* de segmentare, la compunerea acțiunilor care lucrează cu apariția multiplă a marcherilor [4], [5] (vezi Secțiunea 3).

2. Segmentarea de tip clauzal cu algoritmul *M-1997*

Prescurtat în continuare ca “*algoritmul de segmentare M-1997*”, sau simplu “*algoritmul M-1997*”, algoritmul de *segmentare-Marcu* a frazei în unități de tip-clauzal [4], [5] funcționează ca un automat finit, sau ca o rețea de tranziție, bazat pe un set de *stari* și *acțiuni*. În [4] se face o analiză de corpus a potențialilor *marcheri de discurs*, numiți și “*sintagme indicatoare*” [Eng: *cue phrases*] și “*conective*”, cu scopul de a evalua contribuția potențială a diferitelor marcheri la determinarea (delimitarea) unităților textuale elementare pe care sunt definite *relatiile retorice*, în cadrul unității textuale standard care este fraza. În încercarea de a stabili principalele tipuri de funcții ale marcherilor, și anume de tip clauzal, frazal, de discurs, sau pragmatic, algoritmul de segmentare *M-1997* consideră mai întâi următoarele *trei clase* de marcheri:

- (Mar1) În *prima clasă* sunt cuprinși marcherii (sintagmele indicatoare) care joacă un rol în cadrul discursului pentru majoritatea fragmentelor de text ale corpusului analizat. Elementele din (Mar1) vor fi numite în cele ce urmează “*marcheri de discurs*”, iar specifici acestei prime clase sunt marcheri ca “*deși*” [Eng: *although*], “*pe lângă*” [Eng: *besides*], “*dacă*” [Eng: *if*], “*atunci*” [Eng: *then*], etc.
- (Mar2) Marcherii din *a doua clasă*, numiți “*marcherii de frază/clauză*”, joacă în discurs, pentru majoritatea fragmentelor de text în care apar, rolul de *adiacenți* la alți marcheri de discurs sau clauzali. Un membru specific al clasei (Mar2) este considerat a fi “*și*” [Eng: *and*], deoarece are rol clauzal de fiecare dată când apare înaintea altui marcher de discurs sau clauzal, cu toate că poate avea atât rol de discurs cât și clauzal atunci când apare izolat.
- (Mar3) *A treia clasă* conține marcheri care s-au dovedit că joacă un rol de *delimitare a clauzelor* în majoritatea fragmentelor de text investigate în [4]; ei vor fi referiți, simplu, ca “*marcheri clauzali*” (sau “*de clauză*”). (Mar3) include, de asemenea, acei marcheri pentru care analiza de corpus nu a putut distinge între funcția lor de discurs și cea clauzală. “*După*” [Eng: *after*] este un astfel de element reprezentativ al (Mar3).

Marcu [4] a selectat mai mult de 450 de *marcheri* (pentru engleză) în cadrul analizei sale de corpus pentru marcherii de discurs și de frază/clauză. Marcherii sunt stocați și procesați într-o *bază de date* ale cărei înregistrări conțin următoarele câmpuri:

- a. Câmpul denumit *Exemplu* conține un *fragment de text* din care a fost extras marcherul.
- b. Câmpul *Marker* codifică *marcherul* însuși, împreună cu marcherii de punctuație contextuală și, atunci când este necesar, ceilalți *marcheri adiacenți*.
- c. Câmpul *Usage* furnizează unul sau mai multe dintre *rolurile funcționale* ale marcherului:

- (c1) *Frazal/clauzal* (S), atunci când markerul nu îndeplinește *nici o funcție* în structurarea discursului;
 - (c2) *De discurs* (D), când markerul evidențiază o *relație de discurs* între două unități textuale;
 - (c3) *Pragmatic* (P), dacă există o *relație* între o construcție lingvistică (sau non-lingvistică) care conține markerul, și convingerile, planurile, intențiile și/sau scopurile de comunicare ale vorbitorului.
- d. Câmpul *Break_action* (acțiune de oprire) conține un nume de *acțiune* din mulțimea acțiunilor ce vor fi executate în cadrul procesului de segmentare. Acest proces este controlat de către un set de semnalizatori (*flaguri*). Execuția unei acțiuni din mulțimea {NOTHING, NORMAL, COMMA, NORMAL_THEN_COMMA, END, MATCH_PAREN, COMMA_PAREN, MATCH_DASH, SET_AND, SET_OR, DUAL} are unul dintre următoarele efecte:
- (d1) creează o *margină* pentru unitatea textuală elementară în *string*-ul de intrare;
 - (d2) setează un semnalizator (*flag*).
- e. Câmpul *Position* specifică poziția markerului de discurs în cadrul unității textuale careia îi aparține. Valorile acestui câmp sunt *B*, *M* și *E*, după cum markerul este situat *la început* (*B*), *în mijlocul* (*M*) sau, respectiv, *la sfârșitul* (*E*) unității textuale.

3. Algoritmului de segmentare *M-1997*

Algoritm *M-1997* primește în intrare o frază *S* și masivul *markers[n]* al markerilor potențiali de discurs și clauzali din fraza *S*. Masivul *markers[n]* conține markerii recunoscuți în *S*. Fiecare element al acestui masiv este caracterizat de către următoarea *structură de trasaturi*:

- *Acțiunea* asociată acelui marker;
- *Poziția* markerului în cadrul unității textuale elementare (*B*, *M* sau *E*);
- *Semnalizatorul* *has_discourse_function* care inițial este setat la valoarea “no”.

Câteva dintre variabilele importante cu care lucrează algoritmul *M-1997* sunt: “*status*”, “*parenthetical*” și “*clauses*”.

Algoritm M-1997 pentru identificarea unităților de tip-clauzal din cadrul unei fraze are *două* părți principale:

(1) Când variabila “*status*” este NIL, algoritmul *M-1997* execută acțiuni care pot introduce margini ale unității textuale sau pot modifica variabila, influențând procesarea markerilor ulterioari. Pentru partea (1) a algoritmului *M-1997*, atunci când variabila “*status*” ia valoarea NIL, sunt considerate următoarele situații:

-
- (1a) Dacă tipul de marker este DUAL, determinarea marginilor unitatii textuale depinde de markerul adiacent care precede markerul curent analizat. În aceasta situatie, algoritmul *M-1997* seteaza variabila “*status*” la aceeași valoare ca și în cazul unui marker de tip COMMA.
- (1b) Dacă markerul analizat curent nu este adiacent cu markerul imediat precedent, atunci este identificata o margine a unitatii textuale.
- (1c) Cel mai frecvent tip de marker (și de actiune) este NORMAL, marker care identifica o noua unitate de tip clauzal a carei margine-dreapta este data de markerul curent analizat.
- (1d) Când markerul de tip COMMA este precedat de un marker de discurs, *sau_*
- (1e) Tipul markerului este NORMAL_THEN_COMMA, atunci algoritmul *M-1997* identifica o noua unitate de tip-clauzal ca și în cazul markerului de tip NORMAL.
- În oricare dintre cazurile (1c), (1d), (1e), variabila “*status*” este actualizata astfel încât o margine a unitatii textuale să fie identificata la prima aparitie a unei virgule (COMMA).
- (1f) Pentru markerul de tip NOTHING, singura actiune consta în atribuirea markerului o utilizare specifica discursului.
- (1g) Markerii care introduc posibile aparitii de unitati textuale parantetice (texte între paranteze) au doar efectul de a actualiza variabila “*status*”, ca și în cazul aparitiei markerilor “*și*” și “*sau*”.
- (2) Atunci când variabila “*status*” nu este NIL, algoritmul *M-1997* executa actiuni specifice pentru a realiza:
- (2a) Tratarea informatiei din paranteze. O data identificata o paranteza deschisa, o linie-de-despartire [Eng: *dash*] (între doua asemenea liniute se introduce de obicei o apozitie sau un text explicativ), sau un marker de discurs a carui actiune asociata este COMMA_PAREN, algoritmul *M-1997* cauta prima paranteza închisa, linie-de-despartire, sau virgula, ignorând toti ceilalti markeri întâlniti pe parcurs. Acest tratament atrage dupa sine faptul ca informatiei parantetizate *nu* îi este atribuita nici o stare pentru unitatile textuale elementare. Totusi, algoritmul *M-1997* evita stabilirea de margini parantetizate în cazurile în care prima virgula care urmeaza dupa un marker COMMA_PAREN este imediat urmata de un marker “*și*” ori “*sau*”. De mentionat este, de asemenea, ca tratamentul aplicat informatiei dintre paranteze în algoritmul *M-1997* poate conduce la rezultate eronate, ca în exemplul “*I-am dat lui Ion o racheta de tenis, care i-a placut și o minge de plastic, care nu i-a placut*”. Acest tip de erori poate fi evitat printr-o tratare mult mai adecvata în cadrul algoritmului de segmentare *SCD*.
- (2b) Dacă variabila “*status*” contine actiunea COMMA, aparitia primei virgule care nu este adiacenta unui marker “*și*” ori “*sau*” determina identificarea unei noi

unitati elementare de discurs. Algoritmul *M-1997* nu este, capabil, în general, sa distinga suficient de precis între rolurile de discurs si frazale/clauzale ale markerilor “*si*” si “*sau*”. Anumite situatii sunt totusi recunoscute ca introducând functii de discurs, ca de exemplu aparitia unui marker de discurs imediat dupa un “*si*” ori “*or*”, caz în care valoarea semnalizatorului *has_discourse_function* este stabilita la “*yes*”.

Forma originala a algoritmului *M-1997* [4], [5] este extinsa si îmbunatatita în implementarea noastra pentru limba româna (subsectiunea 5.3) cu o analiza mai detaliata la nivelul ei superior, pentru aparitii multiple si corelate ale markerilor de discurs/clauza.

4. Algoritmul de segmentare *SCD-2002*

Aceasta sectiune prezinta partea de segmentare si dependenta, în principal la nivel de clauza, desprinsa din strategia lingvistica *SCD* (*Segmentare-Coeziune-Dependenta*) [1], [2], [3], [6]. Forma actuala a algoritmului, referita în restul articolului prin prescurtarea *SCD-1994*, este foarte apropiata de versiunea publicata în [1], [2]. Noutatea principala a algoritmului *SCD-2002* fata de *SCD-1994* consta într-o rafinare a claselor de markeri, o noua ierarhie a acestora, si în noul algoritm de stabilire a segmentarii si dependentei (structurarii) clauzelor si grupurilor sintagmatiche. Vom pune în evidenta relatia dintre algoritmul *M-1997* si algoritmi *SCD-1994* si *SCD-2002*, aratând ca primul este scufundat în ceilalti doi.

Rezultatele obtinute prin executia algoritmilor de segmentare *M-1997* si *SCD-2002* pe aceleasi fraze conduc la aceeasi concluzie: *SCD-2002* are o granularitate (mult) mai fina a claselor de markeri în comparatie cu cea a claselor algoritmului *M-1997*, iar rafinarea actiunilor implicate în *SCD-2002* conduce la delimitarea de unitati textuale de tip-clauzal mai precise (de fapt mai corecte si mai adecvate) decât cele obtinute de catre algoritmul *M-1997*, pretul computational ce trebuie platit pentru acest fapt ramânând sa fie analizat.

Este de mentionat ca segmentarea clauzala practicata de *SCD-2002* este doar un aspect particular al segmentarii textului, deoarece se obtin si alte “*bucati*” mai mici de text dominate de nuclele semantice de tip N (Substantiv), V (Verb), A (Adjectiv-Adverb). Segmentarea rezultata din clasele de markeri *SCD-2002* se afla într-o strânsa relatie cu noua teorie *X-bar functionala* (FX-bar) [3], o alta componenta importanta a strategiei lingvistice generale *SCD*.

Din schema generala FX-bar propusa în [3] se detaseaza urmatoarele nivele de proiectie la nivel lexical si gramatical:

Tabelul 4.1.

Nivele de proiectie ale schemei FX-bar (vezi [3])

Markeri	Nivelul de Proiectie	Structura gramaticală	Exemple
---------	----------------------	-----------------------	---------

trăsătura PRED sau EXIST (OBJECT)	nivel de lexicon; prin convenție, (BAR = -1)	forma de dicționar a cuvântului; X = N, V, A, Pron, ...	<i>a ploua</i> <i>conducere</i> (trăsătura PRED) <i>clădire</i> (trăsătura EXIST înțelesul obiectual) <i>clădire</i> (PRED, pentru înțelesul acțional) <i>creion</i> (EXIST)
M0-marcher reprezintă aplicarea inflexiunii M0(X)=X0	X0 (BAR = 0)	forma lexicală (de text) a cuvântului; X=N, V, A, ...	<i>plouă</i>
M1-marcher se aplică nucleului X0 M1(X0)=X1	X1 = CL0; (BAR=1) poate fi identificat și cu nivelul 0 de proiecție a clauzei, BAR-CL = 0	sintagme XG (X=N, V, A), i.e. grupuri nominale, verbale, adjectivale- adverbiale	<i>orice steag alb</i> <i>ploua</i> <i>aleargă repede</i> <i>nu aleargă deloc</i> <i>foarte bine studiat</i>
M2-marcher se aplică proiecției X1 M2(X1)=X2=CL1 M2 se aplică unei singure clauze	proiecția X2 = CL1 BAR = 2 și BAR-CL = 1	clauza finită sau infinită	<i>Maria i-a dat un măr</i> <i>ficei sale.</i> <i>O femeie dăruind un măr</i> <i>unui bărbat conține o clauză</i> <i>infinită.</i>
M3(CL1, CL1)=CL2 marcheri de discurs; M3 se aplică la două sau mai multe clauze	nivelul de proiecție X3 = CL2; BAR = 3 și BAR-CL = 2	relații de discurs între clauze finite	<i>Dacă plouă atunci plec mai</i> <i>devreme și îmi i-au și</i> <i>umbrela.</i>

4.1. Clasele de marcheri pentru algoritmul SCD-2002

Pentru algoritmul de segmentare SCD-2002 propunem o anumită rafinare a claselor de marcheri și a ierarhiilor acestor clase din [1], [2], schimbări ce constau în următorul set de marcheri, în concordanță cu Tabelul 4.1. de mai sus:

M3 = { marcheri (de discurs) inter-clauzali }.

Clasa de marcheri M3 este formată din *funcții* sau *relații* (atunci când marcherii sunt corelați), având ca argumente două sau mai multe clauze finite (unele dintre ele pot fi infinite). Acești marcheri sunt ceea ce [4], [5], precum și alte abordări numesc “*marcheri de discurs*”, și se aplică proiecțiilor sintactice de nucleu X2 = CL1 (și de nivel X3), de tip clauzal în teoria FX-bar (vezi Tabelul 4.1.).

M3 poate fi partitionată în următoarele subclase (în ordinea *descrescătoare* a *priorității* de definire a relațiilor de dependență – vezi Fig. 4.1.1.):

M33 = { marcheri (de discurs) inter-clauzali care introduc o dependență (neambigua) de *supra-ordonare strictă* }. *Supra-ordonarea strictă* înseamnă *ridicarea efectivă a (cel puțin) unui nivel de dependență clauzala*, și este reprezentată de marcheri precum “*atunci*”, “*altfel*”, etc.

M32 = { marcheri (de discurs) inter-clauzali care introduc dependenta de *supra-ordonare*, incluzând *semnele de punctuatie* precum doua puncte, punct-si-virgula, paranteza închisa, linie-de-despartire, etc. }. *Supra-ordonarea* presupune *ridicarea* unuia sau mai multor nivele de dependenta clauzala, sau ramânerea pe acelasi nivel de dependenta în cadrul unei dependente de tip-*coordonare*. Exemple tipice de marcheri din clasa M32 sunt “*dar*”, “*asadar*”, “*chiar*”, “*la_fel_(de)*”, “*în_comparatie_(cu)*”, etc.

M31 = { marcheri (de discurs) inter-clauzali care introduc unul sau mai multe nivele de dependenta de *sub-ordonare*, incluzând *semne de punctuatie* ca paranteza deschisa, linia-de-despartire, etc. } Aceasta este o clasa larga de marcheri de discurs formata din numeroase tipuri de relatii între clauze: logice, sintactice, semantice, pragmatice, etc.

Asa cum a fost mentionat mai sus, fiecare dintre clasele M33, M32 si M31 poate, la rândul ei, sa fie partitionata în subclase care contin marcheri de tip relational (exprimati prin corelatie), ce stabilesc relatii între clauze, sau ca functii de clauze (cu cel puțin doua argumente).

M2 = { *marcheri care introduc o clauza (finita sau infinita), sau un grup sintagmatic al carui nucleu semantic este una din categoriile sintactice N, V, A* }. Compusul sintactic (sau *grupul sintagmatic* în termenii [3]) XG, X = N, V, A, poate fi asimilat unei clauze degenerate, infinite (vezi Tabelul 4.1) în cazul X = N, A.

M2 este divizata în urmatoarele subclase (în ordinea descrescatoare a prioritatii de introducere a relatiilor de dependenta):

M25 = { marcheri care introduc *clauza relativă* }.

Explicatia consta în faptul ca o clauza relativa reprezinta cea mai complexa unitate sintagmatica ce joaca rol de modificador, si care se aplica nucleului NG al clauzei relative:

M24 = { aparitia unui *grup verbal finit* (FVG) sau, echivalent, aparitia valorii FINITE pentru trasatura TENS atribuita unui verb, introducând deci o *clauza finita* }.

Întregul grup verbal poate mosteni valoarea trasaturii FINITE daca nucleul sau V sau alta componenta importanta din VG poarta aceasta valoare a trasaturii TENS (de exemplu, auxiliarul din VG).

M23 = { aparitia unei sintagme *predicationala* XG (sau X1), X=V, N, A, al carei nucleu semantic este o *categorie predicationala*, purtând valoarea PRED = ACT (posibil înca la nivel de lexicon), si introducând astfel o *clauza infinita* }.

Clasele de marcheri M24 si M23 introduc structuri de nivel-X2, si anume clauze finite sau infinite, formate dintr-o sintagma X1 (sau grup XG, X = N, V, A) care reprezinta *nucleul semantic, finit* (TENS = FINITE) sau *predicational* (PRED = ACT), al structurii de nivel-X2, urmata de sateliti (argumente si/sau adjuncti) corespunzatori de tip NG (inclusiv NG-uri prefixate de o prepozitie, deci clasica sintagma PP). Unele dintre argumente, cum este cazul clasic al *subiectului gramatical*, pot preceda nucleul semantic de tip X1 al clauzei careia îi apartin [3]. Sa mai precizam ca exista o *ordine sistemica (canonica)* [18], [19], a satelitelor, sau “*actantilor*” (argumente si adjuncti) dintr-o clauza (finita sau

infinita): ACT(or), PAT(ient), ADDR(essee), ORIG(ine), LOC(ation), etc. Ordinea canonica este specifica fiecarui LN, si se poate obtine în urma unei cercetari statistice si lingvistice foarte atenta [18].

Putem gasi recent un *principii de predicativitate* similare cu cel folosit în *strategia lingvistica* SCD, si aplicat la sintagmele nominale din limba italiana [14], sau la adjectivele “*deverbale*” [14], [16]. În timp ce predicativitatea verbelor este frecventa si naturala, trasatura de *nepredicativitate* [17, p. 22] (de fapt, nepredicationalitate) a verbelor de tip *existential* este si ea la fel de frecventa (formele lui “*a fi*”), valoarea lor FINITE, dublata sau nu de valoarea trasaturii PRED = ACT, anunând totusi aparitia unei clauze finite.

M22 = { marcheri care introduc relatii de tip-JOIN, adica *conjunctii* de tipul “*si*”, “*sau*”, “*la fel ca (si)*”, “*împreuna (cu)*” }.

M21 = { COMMA (sau VIRGULA) }.

Clasele M22 si M21 cuprind marcheri cu un grad important de ambiguitate deoarece pot introduce orice structura de tip X1 (grupuri XG, X = N, V, A) sau X2 (clauze finite sau infinite).

M1 = { *marcheri care delimiteaza (introduc) structuri XG* }.

Conform strategiei SCD si teoriei FX-bar [3], clasa de marcheri M1 consta în *marcheri de nivel-X1*, X = N, V, A, adica marcheri care se aplica constructiilor sintactice de nivel-X1 (denotat si XG, si numit X-grup). Aceste sintagme constau, de fapt, dintr-un *nucleu semantic* înconjurat de *modificatori* (adjective sau adverbe) si/sau *specificatori* (sau *cuantificatori*, unii generalizati, printre cuantificatori incluzându-se determinatorii, negatia, etc.).

Asa cum exista o ordine sistemica a satelitilor unui nucleu semantic într-o clauza (sintagma de nivel-X2), în mod similar exista o “*ordine structurala*”, data de “*distanta*” modificatorilor, cuantificatorilor, prepozitiilor, etc. fata de nucleul X0, pentru constituentii unei sintagme de nivel-X1. Astfel, în limba româna (franceza, engleza), cel mai “apropiat” fata de nucleul X0 trebuie sa fie *modificatorul* (adjectivul sau adverbul), urmeaza apoi *cuantificatorul* (care ocupa locul modificatorului daca acesta lipseste), apoi prepozitia (adpozitia, în general), etc. De exemplu, nu este sintactic corecta sintagma “*frumos orice copil*”, sau “*orice frumos pe copil*”. Nucleul X0 înconjurat de modificatori si/sau specificatori (cuantificatori) poate fi marcat functional prin *pre-pozitii* (în cazul grupului nominal NG din româna, engleza, franceza), dar si prin *post-pozitii* (în cazul NG sau VG din engleza sau germana). Marcarea clitic-functionala (prin particule *pre-* sau *post-pozitionale*) poate exprima *cazul* (pentru NG), sau *timpul, semantica* (pentru VG), etc. Principalele elemente componente ale unei structuri XG corespund si subclaselor de marcheri ai clasei M1.

M1 poate fi divizata în subclase de marcheri, subclase utile în delimitarea substructurilor XG (X1), X = N, V, A, în conformitate cu un criteriu cum este *distanta* dintre nucleul semantic X0 si elementele functionale care îl “*înconjoara*”; un asemenea nucleu este, în ultima instanta, un substantiv comun obiectual (numit si *autosemantic* în

[19]), un nume propriu, sau un substantiv personalizat (dar fara nume propriu, denominalizat).

M14 = { aparitia unui substantiv comun *obiectual* (nepredicational, autosemantic), a unui nume propriu, sau a unui substantiv personalizat denominalizat }

M13 = { aparitia unui *modifier* (adjectiv, adverb, adjectiv pronominal) }

M12 = { aparitia unui *cuantificator* (generalizat) }

M11 = { pre-pozitii sau *post-pozitii* exprimând *cazul* (pentru N), timpul sau *semantismul* (pentru V), etc. }

Ultima clasa de marcheri, notata **M0** (sau M00 pentru uniformitate), si ai carei marcheri se aplica *formei de dictionar* a cuvântului, este reprezentata de *rolul functional* al *flexionarii*.

Recapitulând, clasele de marcheri considerate de strategia lingvistica SCD, în particular de *algoritmul de segmentare SCD-2002*, pot fi reprezentate grafic de urmatoarea ierarhie:

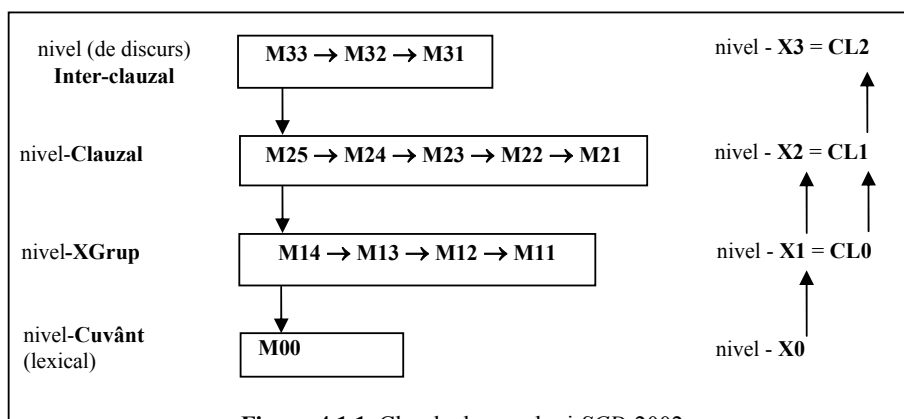


Figura 4.1.1. Clasele de marcheri *SCD-2002* și ierarhia lor

Orientarea arcelor din Fig. 4.1.1., stabilite între clasele și subclasele de marcheri, provine dintr-o *ordine de prioritatea descrescătoare* între marcherii considerați, și este reprezentată mai jos prin relația “ \bowtie ” dintre clasele și subclasele de marcheri. Aceasta ierarhie este o *ipoteza de baza* impusă în *strategia lingvistica SCD* și, prin consecință, și în *algoritmul de segmentare SCD*.

$$(4.1.2) \quad \forall (j = 1 \div 4) \quad M^{(k+1)(j+1)} \bowtie M^{(k+1)j} \quad (k = 0 \div 2);$$

$$(4.1.3) \quad \forall (k = 0 \div 2) \quad M^{(k+1)i} \bowtie M_{kj} \quad (i = 1 \div 5), (j = 0 \div 5).$$

Aceste inegalitati ne spun ca marcherii din subclasa $M(k+1)(j+1)$ sunt de *prioritate mai mare* în comparatie cu marcherii din subclasa $M(k+1)j$, ($k = 0 \div 2$), ($j = 1 \div 4$), în cadrul aceleiasi clase $M(k+1)$ de marcheri aflata pe *acelasi nivel* de proiectie lingvistica, iar marcherii din aceeasi clasa $M(k+1)$ au o *prioritate mai mare* fata de marcherii din clasa M_k de pe nivelul *inferior* de proiectie lingvistica.

Aceasta ierarhie a marcherilor si claselor de marcheri este considerata de noi ca fiind valida pentru *limba româna*. Probabil ca anumite modificari vor fi necesare când se trece de la un LN la altul. Daca ne situam în domeniul mai restrâns al limbajelor indo-europene (cum sunt franceza, engleza, germana, italiana, spaniola, posibil rusa), se poate aprecia ca structurile si clasele de marcheri propuse în Tabelul 4.1. si Fig. 4.1.1. ramân aceleasi sau foarte asemanatoare, cu anumite modificari parametrizate în functie de limbaj.

4.2. Algoritmul SCD-2002 de segmentare si stabilire a dependentelor

Urmând algoritmi de segmentare si dependenta (numiti si *meta-algoritmi SCD*) propusi în [1] si [2] (denotati în continuare *SCD-1994*), rafinati cu clasele de marcheri considerate în subsectia precedenta, se obtine forma prezenta a *algoritmului de segmentare SCD* (denotata *SCD-2002*). Dezvoltam aici forma *secvential-liniara* a acestui algoritm, însa în [1] sunt expuse si o forma *secvential-recursiva*, ca si o versiune *paralela* a algoritmului. O forma "*inversata*" (pentru care în intrare avem un arbore de derivare sau o formula logica, iar în iesire – ca si în intrarea în algoritmul standard – avem o fraza) poate fi folosita pentru a ghida procesul de generare a unei fraze de LN [2], schimbând operatia de recunoastere a marcherilor cu cea de generare a lor, si analiza (parsarea) compusilor sintactici cu generarea lor.

În *descrierea algoritmului* de segmentare *SCD-2002* sunt folosite câteva *operatii* al căror înțeles este bine să fie precizat de la început.

- (4.2.a) *Recunoasterea marcherilor* înseamna inserarea în text a unor etichete adecvate, ce corespund marcherilor care realizeaza delimitarea unitatilor textuale sintactice, semantice, si de discurs.
- (4.2.b) *Verificarea marcherului* înseamna preluarea, din baza de date a marcherilor, a celor mai importante valori din structura de trasaturi a acelu marcher.
- (4.2.c) *Segmentarea* implica a analiza liniara (parsare) a secventei de etichete de marcheri, si *recunoasterea* unei subsecvente (eventual discontinua) care face parte din secventa originala de etichete de marcheri.
- (4.2.d) *Recunoasterea structurii sintactice* înseamna *segmentarea* si *recunoasterea* structurilor sintactice elementare cum sunt NG, VG, AG, clauza infinita, si clauza finita.
- (4.2.e) *Compunerea structurilor (de dependenta)* consta în stabilirea dependentelor (sub-ordonare, co-ordonare, supra-ordonare) succesive

dintre structurile sintactice recunoscute, pe baza rolului functional specific al markerilor care delimiteaza aceste structuri, si utilizând ierarhia corespunzatoare dintre clasele carora le apartin acesti narcheri (vezi Fig. 4.1.1. si relatiile (4.1.2.-4.1.3.).

Algoritmul de segmentare SCD-2002

Step01. Recunoasterea pe text a markerilor din clasa M3;

Step01. Recunoasterea pe text a markerilor din clasa M2;

Step03. Verificarea contextuala si recunoasterea aparitiei corelate a markerilor de tip M3 si M2⁽¹⁾;

Step04. Segmentarea frazei în clauze finite;

Step05. Segmentarea (chunking), daca este necesar, a clauzelor finite în clauze infinite;

[Stop: Daca scopul procesarii este de a obtine o structura liniara a clauzelor finite si/sau infinite din fraza].

Step06. Verificarea markerilor M3 si stabilirea relatiilor de dependenta inter-clauzala⁽²⁾;

[Stop: Daca scopul procesarii este doar de a obtine arborele de dependenta a clauzelor finite (si infinite) din fraza].

Step07. Recunoasterea pe text a markerilor din clasa M1;

Step08. Verificarea contextuala si recunoasterea (eventualei aparitii corelate) a markerilor M1⁽³⁾;

Step09. Recunoasterea structurilor XG ($X = N, V, A$)⁽⁴⁾;

Step10. Verificarea markerilor M24 si M23, si stabilirea relatiilor de dependenta dintre structurile infinite, intra-clauzale de tip XG⁽⁵⁾;

[Stop].

Indicii superiori (**n**) care apar în algoritmul de mai sus corespund urmatoarelor *remarci*:

(1) Marcherii corelati pot fi reprezentati ca upluri ordonate (liste) de marcheri.

(2) Relatiile de dependenta clauzala pot fi stabilite (ca în [8, Anexa 9, p. 108], de exemplu) prin utilizarea unei *gramatici formale* (ambigue) definita pe secvente de marcheri din (sub)clasele M3, M25, M22, si M21.

(3) Marcherii complecsi pot fi sintagme sau expresii de tipul gradelor de comparatie a adjectivelor, diferiti cuantificatori generalizati, etc.

(4) În executia acestui pas se realizeaza parsarea sintagmelor XG dintr-o clauza finita si infinita.

- (5) Dependentele dintre structurile de tip XG sunt stabilite în principal prin utilizarea trasaturilor și valorilor de trasaturi TENS = FINITE sau INFINITE, și PRED = ACT sau EXIST, pe care le posedă nucleele semantice ale sintagmelor XG, X = N, V, A (a se vedea [3]). Aceste valori pot fi *mostenite* din reprezentarea de lexicon a cuvintelor care poartă aceste trasaturi și care formează XG, sau pot fi *dobândite* de către nucleul semantic al XG în procesul de recunoaștere (parsare) a structurii.

5. Compararea algoritmilor de segmentare

5.1. Algoritmii de segmentare SCD-1994 și SCD-2002

Algoritmii SCD-1994 expuși în [1], [2] se bazează pe *patru* (sub)clase principale de marcheri, denotate acolo prin (clasele de) “1-marcheri” până la “4-marcheri”. Aceste subclase de marcheri din SCD-1994 corespund următoarelor (sub)clase de marcheri din prezentul algoritm SCD-2002:

- (5.1.1) 1-marcheri = $M3 \cup M25 \cup M22$;
 2-marcheri = $M24$;
 3-marcheri = $M23$;
 4-marcheri = $M21 \cup M1$

Prezentăm în continuare algoritmul de segmentare SCD-1994 (în forma secvențial-recursivă), așa cum a fost expusă în [1, p.68-69], având ca scop parsarea LN. Algoritmul SCD-1994 (în forma secvențial-liniară) și destinat sarcinii de *generare* a LN este prezentat în [2, p.172-173].

Algoritm de segmentare SCD-1994 în forma secvențial-recursivă (SR)

- Step01.** Recunoașterea marcherilor de clauza.
Step02. Recunoașterea sintagmelor VG (grupuri verbale) finite și infinite.
Step03. Verificarea contextuală a marcherilor.
Step04. Segmentarea clauzala.
Step05. Segmentarea sub-clauzala.
Step06. Recunoașterea 1-marcherului;
 Recunoașterea 1-structurii:
Wait-until 1-structura este completă.
Step07. Recunoașterea 2-marcherului;
 Recunoașterea 2-structurii:

Wait-until structura de nivel-X2 este completa*.

Step08. Recunoasterea 3-marcherului;

Recunoasterea 3-structurii.

Step09. Recunoasterea 4-marcherului;

Procesarea 4-structurii.

Step10. 3-structura completa?

Nu: *Go-to* Step08.

Da: Compune 3-structuri; *Go-to* Step11.

Step11. 2-structura completa ?

Nu: *Go-to* Step07.

Da: Compune 2-structuri; *Go-to* Step12.

Step12. 1-structura completa ?

Nu: *Go-to* Step06.

Da: Compune 1-structuri; *Go-to* Stop.

Stop.

* *Structuri AX-bar* (în original, în [1]), înțelegând structuri sintactice derivate din *schemele X-bar augmentate*, definite în [20] și extinse în [3]. Scopul acestui pas al algoritmului este de a completa clauza finită introdusă printr-un grup verbal finit.

Principala problemă cu algoritmul de segmentare și dependența *SCD-1994* (forma SR) este că sunt necesare “multiple nivele de recursie pentru a completa și compune structurile” [1, p.69].

5.2. Algoritmii de segmentare *M-1997* și *SCD-2002*

În această subsecțiune vom arăta că algoritmul de segmentare *M-1997* este *scufundat* în algoritmul *SCD-2002* (de fapt, și în *SCD-1994*).

M-1997 este un algoritm de “*suprafata*” destinat segmentării discursului în unități textuale de tip-clauzal. În timp ce, pentru acest scop, *M-1997* folosește numai *marcheri de discurs* (“*cue phrases*” sau *conective*), algoritmul *SCD-2002* utilizează un set de clase de marcheri mai larg și în același timp mai rafinat, set care include clasele de marcheri din *M-1997* ca un caz particular. Mai precis, relațiile dintre clasele de marcheri Mar1, Mar2, și Mar3 (vezi Secțiunea 2) utilizate pentru *M-1997*, și clasele de marcheri Mkj ale algoritmului *SCD-2002* sunt următoarele:

$$(5.2.1) \quad \text{Mar1} \cup \text{Mar2} \cup \text{Mar3} \subseteq \text{M3} \cup \text{M25} \cup \text{M22} \cup \text{M21}$$

sau, posibil, mai precis:

$$(5.2.2) \quad \text{Mar1} \cup \text{Mar3} \subseteq \text{M3} \cup \text{M25} \text{ și } \text{Mar2} \subseteq \text{M22} \cup \text{M21}$$

Diferența dintre algoritmi *M-1997* și *SCD-2002* nu constă doar în faptul că al doilea algoritm are un număr mai mare de clase, care sunt mai fine (mai precise), ci, mai important este faptul că aceste clase formează un *sistem ierarhic* (expus în Fig. 4.1.1.) ce este utilizat în procesele de *segmentare* și de *stabilire a dependentelor*. *SCD-2002* furnizează noi clase de marcheri, cum sunt M23 și M24 (aparitia categoriilor predicative și/sau având un timp finit), precum și clasa M1, cu subclasele sale (aparitia unor componente ale sintagmei XG, X = N, V, A). Acesta este un *prim argument* din care rezultă că *M-1997* este *scufundat* în *SCD-2002*. “*Scufundarea*” este un termen care reflectă, de fapt, un proces de rafinare și de creștere a preciziei în calculul marginilor (limitelor) unităților textuale și a dependentelor dintre ele, pentru *SCD-2002* în comparație cu *M-1997*.

Al doilea argument important care susține validitatea relației afirmate între cei doi algoritmi este următorul: fiecare *actiune* din *M-1997* are un corespondent într-o *operatie* (sau o multime de *operatii*) din algoritmul *SCD-2002* (subsecțiunea 4.2).

Pentru segmentare, algoritmul *M-1997* asociază fiecărui marcher, în baza de date a marcherilor, o anumită *actiune* ce este statistic determinată de către analiza de corpus efectuată în [4]. Corespondența dintre operațiile algoritmului *SCD*, și o *actiune* din algoritmul *M*, se face în felul următor:

- (5.2.a) *Actiunea* (și marcherii) NORMAL din algoritmul *M* are același efect cu operațiile de procesare a marcherilor de discurs din clasa M3 a algoritmului *SCD*. Când este întâlnit un asemenea marcher, aceasta înseamnă că o clauză (în *SCD-2002*) sau o unitate de tip-clauzal (în *M-1997*) este pe cale de a se încheia și o alta clauză, respectiv unitate de tip-clauzal, este probabil că va începe.
- (5.2.b) *Actiunile* COMMA, SET_AND, și SET_OR din algoritmul *M* sunt folosite pentru a dezambigua rolul unor marcheri din M3 pentru care nu se poate aplica întotdeauna regula generală (*actiunea* NORMAL). Acești marcheri sunt următorii pentru limba română: “,“ [Eng: *comma*], “și”, și “sau”. Rolul acestor marcheri este ambiguu deoarece comportamentul lor nu este uniform în cadrul delimitării unităților textuale. *SCD-2002* rezolvă aceste cazuri cu ajutorul utilizării unei gramatici formale de marcheri care descrie principalele reguli de delimitare și dependență a clauzelor (în limba română). Această gramatică (vezi indicele superior (2) din *SCD-2002* și remarca corespunzătoare) are ca scop să recunoască secvențele cele mai frecvente de marcheri din clasele M3 și M2. Numai *câteva* dintre aceste reguli sunt încorporate în mod explicit în algoritmul *M-1997* original.
- (5.2.c) O unitate de tip-clauzal din *M-1997* nu este în mod necesar o clauză finită în sensul gramatical al notiunii, așa cum este folosit în algoritmul *SCD*. O asemenea unitate de tip-clauzal, în sens *M-1997*, poate fi o întreagă frază, formată din mai multe clauze finite. *M-1997* folosește, de fapt, pentru segmentarea liniară a frazei în unități de tip-clauzal numai *trei reguli* din

cele folosite de *SCD-2002*, iar aceste reguli sunt sintetizate de catre *actiunile* *COMMA*, *SET_AND*, *SET_OR*.

- (5.2.d) *Actiunile* *MATCH_PAREN*, *MATCH_DASH*, *COMMA_PAREN* sunt utilizate de catre *M-1997* pentru a delimita acele întinderi de text care pot fi omise atunci când fraza este segmentata în unitati de tip-clauzal. Aceste parti “explicative” din text, considerate a nu fi importante, sunt, în text, puse între *paranteze*, (perechi de) *liniute-de-despartire*, sau (perechi de) *virgule*. Algoritmii *M-1997* nu trateaza aceste întinderi “parantetizate” de text ca fiind unitati de tip-clauzal propriu-zise, ci le considera ca doar ca fiind scufundate în unitatea de tip-clauzal de care apartin. Pentru *SCD-2002*, aceste *actiuni* *M-1997* nu au un corespondent specific deoarece *paranteza* (închisa si deschisa), *virgula*, si *liniuta-de-despartire* sunt tratate ca *marcheri* de discurs (*M3*), si fac parte din *gramatica de marcheri compusi* (concatenati) care este asociata cu algoritmul *SCD-2002* de segmentare si dependenta a clauzelor dintr-o fraza.
- (5.2.e) Din acelasi motiv ca cel mentionat mai sus, în (5.2.d), *actiunile* *DUAL*, *NORMAL_THEN_COMMA* din *M-1997* nu au, nici ele, un corespondent în *SCD-2002*; aceste doua actiuni sunt de asemenea înglobate în *gramatica formala de secvente de marcheri de discurs*, care se dovedeste a fi, în mod clar, mai generala, usor de extins (sau de restrâns), este dependenta de LN specific analizat, si modeleaza comportamentul *marcherilor simpli* si *compusi* (concatenati) de tip *M3* si *M2*.

Relatiile (5.2.1-2) si observatiile (5.2.a-e) demonstreaza ca algoritmul de segmentare *M-1997* este (chiar strict) scufundat în algoritmiile-*SCD* (atât *SCD-2002* cât si *SCD-1994*). Acest fapt, stabilit teoretic aici, este confirmat de catre rezultatele empirice ale implementarilor, prezentate în subsectiunea care urmeaza.

5.3. Executia segmentarii pentru algoritmiile *M-1997* si *SCD-2002*

Actuala etapa de implementare a algoritmilor de segmentare este prezentata în exemplele care urmeaza. *Step06* din *SCD-2002*, si *Step12* din *SCD-1994* stabilesc relatiile de *dependenta inter-clauzala*, folosind o gramatica formala pentru *marcherii* de discurs, simpli si compusi (concatenati), din clasele *M3* si *M2*. Aceasta faza a algoritmului nu este încă implementata, în prezent. Sa mentionam ca stabilirea *dependentelor intra-clauzale* este (partial) implementata prin utilizarea, pentru moment, (numai) a subclaselor *M2* si *M1* de *marcheri*. Marginile inter-clauzale din text sunt reprezentate prin *paranteze patrate*, în timp ce pentru marginile si dependentele intra-clauzale sunt folosite *parantezele rotunde* (obisnuite). Indicii inferiori ai *parantezelor patrate* arata numarul curent al unitatilor textuale de tip-clauzal din algoritmul *M-1997*, respectiv numarul curent al clauzei obtinute din algoritmul *SCD-2002*.

Exemplul 5.3.1.

Ex.5.3.1.Tag. (Etichetarea morfologica realizata cu mediul *TexTag* – vezi Fig. 5.4.1. si Fig. 5.4.2.)

<NSRY,23,0>Câmpul</NSRY,23,0> <V3,24,0>era verde</V3,24,0>
 <CR,25,0>si</CR,25,0> <NSRY,26,0>vita</NSRY,26,0>
 <S,27,0>de</S,27,0> <NSRN,28,0>vie</NSRN,28,0>
 <PXA,29,0>se</PXA,29,0> <V3,30,0>acoperise</V3,30,0>
 <S,31,0>cu</S,31,0> <NPN,32,0>lastari</NPN,32,0>
 <APN,33,0>verzi</APN,33,0><COMMA,34,0>,</COMMA,34,0>
 <NPRY,35,0>copacii</NPRY,35,0> <S,36,0>de pe</S,36,0>
 <NSRY,37,0>marginea</NSRY,37,0>
 <NSOY,38,0>soselei</NSOY,38,0> <V3,39,0>înfrunziseră</V3,39,0>
 <CR,40,0>si</CR,40,0> <NSRY,41,0>briza</NSRY,41,0>
 <V3,42,0>suflă</V3,42,0> <S,43,0>dinspre</S,43,0>
 <NSRN,44,0>mare</NSRN,44,0><POINT,45,0>.</POINT,45,0>

Ex.5.3.1.Mar. (Rezultatul segmentarii (fara dependente), obtinut prin aplicarea algoritmului *M-1997* în cadrul mediului *ClauSEGM* – vezi Fig. 5.4.3.)

[Câmpul era verde si vita de vie se acoperise cu lastari verzi, copacii de pe marginea soselei înfrunziseră si briza suflă dinspre mare.]₁

Ex.5.3.1.SCD. (Rezultatul segmentarii (fara dependente), obtinut prin aplicarea algoritmului *SCD-2002* în cadrul mediului *ClauSEGM* – vezi Fig. 5.4.4.)

[(Câmpul) era verde]₁ si[(vita) (de (vie)) se acoperise (cu (lastari) (verzi))]₂ , [(copacii) (de pe (marginea (soselei))) înfrunziseră]₃ si[(briza) suflă (dinspre (mare))].]₄

Exemplul 5.3.2.

Ex.5.3.2.Tag.

<S,1,0>În</S,1,0> <NSN,2,0>întuneric</NSN,2,0> <V2,3,0> ai fi zis</V2,3,0>
 <C,4,0>ca</C,4,0> <V3,5,0>fulgera </V3,5,0> <R,6,0>ca</R,6,0>
 <NSRY,7,0>vara</NSRY,7,0> <COMMA,8,0>,</COMMA,8,0> <C,9,0>dar</C,9,0>
 <NPRY, 10,0>noptile</NPRY,10,0> <V3,11,0>erau reci</V3,11,0>
 <CR,12,0>si</CR,12,0> <QZ,13,0>nu</QZ,13,0> <PPSD, 14,0>ti</PPSD,14,0>
 <PXA,15,0>se</PXA,15,0> <V3,16,0> parea</V3,16,0> <R,17,0>deloc</R,17,0>
 <C,18,0>ca</C, 18,0> <PXA,19,0>se</PXA,19,0> <V3,20,0>apropie</V3,20,0><NSRY,21,0>furtuna</NSRY,21,0><POINT,22,0>.</POINT,22,0>

Ex.5.3.2.Mar.

[În întuneric ai fi zis]₁ [ca fulgera ca vara,]₂ [dar noptile erau reci si nu ti se parea deloc]₃ [ca se apropie furtuna.]₄

Ex.5.3.2.SCD.

[(În (întuneric)) ai fi zis]₁ [ca fulgera (ca (vara))]₂ , [dar (noptile) erau reci]₃ si [nu (ti) se parea (deloc)]₄ [ca se apropie (furtuna)].]₅

Example 5.3.3.**Ex.5.3.3.Tag.**

<NSRY,46,0>Poarta</NSRY,46,0> <V3,47,0>era deschisa </V3,47,0>
 <COMMA,48,0>,</COMMA,48,0> <TSR,49,0>un</TSR,49,0> <NSN,50,0>soldat
 </NSN,50,0> <V3,51,0>sedea</V3,51,0> <S,52,0>la</S,52,0> <NSN, 53,0>soare<
 /NSN,53,0> <S,54,0>pe</S,54,0> <TSR,55,0> o</TSR,55,0> <NSRN,56,0>
 banca</NSRN,56,0><COMMA, 57,0>, </COMMA,57,0> <TSR,58,0>o</TSR,58,0>
 <NSRN, 59,0>ambulanta</NSRN,59,0> <V3,60,0>astepta</V3,60,0> <S,61,0>la </S,61,0>
 <NSRY, 62,0>usa</NSRY,62,0> <S, 63,0>de</S,63,0> <NSN,64,0> serviciu</NSN,64,0>
 <CR, 65,0>si</CR,65,0> <VG,66,0>intrând</VG,66,0> <V1,67,0> am simtit</V1,67,0>
 <NSRY,68,0>mirosul</NSRY,68,0> <NSOY,69,0> pardoselii</NSOY,69,0>
 <S,70,0>de</S,70,0> <NSRN,71,0>marmura</NSRN,71,0> <S,72,0>si</S,72,0>
 <S,73,0>de</S,73,0> <NSN,74,0>spital</NSN,74,0><POINT, 75,0>.</POINT,75,0>

Ex.5.3.3.Mar. (întindere de text între paranteze acolade {...})

[Poarta era deschisa, {un soldat sedea la soare pe o banca,} o ambulanta astepta la
 usa de serviciu si intrând am simtit mirosul pardoselii de marmura si de spital.]₁

Ex.5.3.3.SCD.

[(Poarta) era deschisa]₁ , [(un (soldat)) sedea (la (soare) (pe (o (banca))))]₂ , [(o
 (ambulanta)) astepta (la (usa) (de (serviciu)))]₃ si [intrând am simtit (mirosul (pardoselii))
 (de (marmura) (si (de (spital))))]₄

Example 5.3.4.**Ex.5.3.4.Tag.**

<NPRY,1,0>Trupele</NPRY,1,0> <V3,2,0>treceau</V3,2,0> <S,3,0>pe
 lângă</S,3,0><NSRN,4,0>casa</NSRN,4,0><COMMA,5,0>,</COMMA,5,0> <S,
 6,0>pe</S,6,0> <NSRN,7,0>sosea</NSRN,7,0><COMMA,8,0>,</COMMA,8,0>
 <CR,9,0>si</CR,9,0> <NSRY,10,0>praful</NSRY,10,0> <RELO,11,0>pe care</
 RELO,11,0><Z,12,0>-</Z,12,0><PPSA,13,0>l</PPSA,13,0> <V3,14,0>ridicau</
 V3,14,0> <PXA,15,0>se</PXA,15,0> <V3,16,0>asternea</V3,16,0> <S,17,0>pe
 </S,17,0> <NPRY,18,0>frunzele</NPRY,18,0> <NPOY,19,0>copacilor</NPOY,
 19,0><POINT,20,0>.</POINT,20,0>

Ex.5.3.4.Mar.

[Trupele treceau pe lângă casa, pe sosea, si praful]₁ [pe care-l ridicau se asternea
 pe frunzele copacilor.]₂

Ex.5.3.4.SCD. (clauza relativa – atributiva)

[(Trupele) treceau (pe lângă (casa)) , (pe (sosea))]₁ , [si (praful) [pe care-(l)
 ridicau se asternea (pe (frunzele (copacilor)))]₂]₃

Example 5.3.5.**Ex.5.3.5.Tag.**

<QZ,76,0>Nu</QZ,76,0> <PPSA,77,0>m</PPSA,77,0><Z,78,0>-</Z,78,0>
 <V3,79,0>a vazut</V3,79,0> <CR,80,0>si</CR,80,0> <QZ,81,0>n</QZ,81,0><Z, 82,0>-
 </Z,82,0><V1,83,0>am stiut</V1,83,0> <C,84,0>daca</C,84,0> <V3,85,0> e</V3,85,0>
 <NSRY,86,0>cazul</NSRY,86,0> <C,87,0>sa</C,87,0> <PPSA, 88,0>ma</PPSA,88,0>
 <V1,89,0>duc</V1,89,0> <S,90,0>la</S,90,0> <PPS,91,0> el</PPS,91,0>
 <C,92,0>sa</C,92,0><Z,93,0>-</Z,93,0><PPSA,94,0>i</PPSA, 94,0>
 <V1,95,0>raportez</V1,95,0> <C,96,0>ca</C,96,0> <V1,97,0>am sosit</ V1,97,0>
 <C,98,0>sau daca</C,98,0> <QZ,99,0>nu</QZ,99,0> <V3,100,0>e mai bine</V3,100,0>
 <C,101,0>sa</C,101,0> <PPSA,102,0>ma</PPSA,102,0> <V1, 103,0>duc</V1,103,0>
 <C,104,0>sa</C,104,0> <PPSA,105,0>ma</PPSA,105,0> <V1,106,0>aranjez</V1,106,0>
 <R,107,0>putin</R,107,0><POINT,108,0>.</POINT, 108,0>

Ex.5.3.5.Marc.

[Nu m-a vazut si n-am stiut]₁ [daca e cazul]₂ [sa ma duc la el]₃ [sa-i raportez]₄
 [ca am sosit sau]₅ [daca nu e mai bine]₆ [sa ma duc]₇ [sa ma aranjez putin.]₈

Ex.5.3.5.SCD.

[Nu (m)-a vazut]₁ si[n-am stiut]₂ [daca e (cazul)]₃ [sa (ma) duc (la (el))]₄ [sa-
 (i) raportez]₅ [ca am sosit]₆ [sau daca nu e mai bine]₇ [sa (ma) duc]₈ [sa (ma) aranjez
 (putin).]₉

Example 5.3.6.**Ex.5.3.6.Tag.**

<NSRY,109,0>Fereastra</NSRY,109,0> <V3,110,0>era deschisa</V3,
 110,0><COMMA,111,0>,</COMMA,111,0> <NSRY,112,0>patul</NSRY,112,0>
 <PSS,113,0>meu</PSS,113,0> <V3,114,0>era acoperit</V3,114,0> <S,115,0>
 cu</S,115,0> <NSRY,116,0>patura</NSRY,116,0><COMMA,117,0>,</COMMA,
 117,0> <NSRY,118,0>masca</NSRY,118,0> <S,119,0>de</S,119,0> <NPN,
 120,0>gaze</NPN,120,0> <S,121,0>cu</S,121,0> <NSRY,122,0>cutia</NSRY, 122,0>
 <PSS,123,0>ei</PSS,123,0> <ASN,124,0>lunguiata</ASN,124,0> <S,125,
 0>de</S,125,0> <NSRN,126,0>tinichea</NSRN,126,0> <CR,127,0>si</CR,127, 0>
 <NSRY,128,0>casca</NSRY,128,0> <S,129,0>de</S,129,0> <NSN,130,0>
 otel</NSN,130,0> <V3,131,0>erau agatate</V3,131,0> <S,132,0>pe</S,132,0>
 <DMSR,133,0>acelasi</DMSR,133,0> <NSN,134,0>cuier</NSN,134,0><POINT,
 135,0>.</POINT,135,0>

Ex.5.3.6.Mar. (întindere de text între paranteze acolate {...})

[Fereastră era deschisă, {patul meu era acoperit cu patura}, masca de gaze cu cutia
 ei lunguiată de tinichea și casca de otel erau agățate pe același cuier.]₁

Ex.5.3.6.SCD.

[(Fereastra) era deschisa]₁ ,[(patul) (meu) era acoperit (cu (patura))]₂ ,[(masca) (de (gaze) (cu (cutia) (ei (lunguiata)) (de (tinichea)))) si (casca) (de (otel)) erau agatate (pe (acelasi (cuier))).]₃

Example 5.3.7.**Ex.5.3.7.Tag.**

```
<V1,1,0>As    vrea</V1,1,0>    <C,2,0>sa</C,2,0><Z,3,0>-</Z,3,0><PPSD,
4,0>ti</PPSD,4,0> <V1,5,0>spun</V1,5,0> <C,6,0>ca</C,6,0> <CR,7,0>si</CR, 7,0>
<R,8,0>mai</R,8,0>    <R,9,0>târziu</R,9,0><COMMA,10,0>,</COMMA,10,0>
<CR,11,0>si</CR,11,0>    <S,12,0>într</S,12,0><Z,13,0>-</Z,13,0><ASN,14,0>alta
</ASN,14,0>    <NSRN,15,0>parte</NSRN,15,0><COMMA,16,0>,</COMMA,16,0>
<V1,17,0>am    vazut</V1,17,0>    <C,18,0>ca</C,18,0>    <NPRY,19,0>lucrurile</
NPRY,19,0> <PXA,20,0>se</PXA,20,0> <V3,21,0>întâmpla</V3,21,0> <R,22,0> tot
asa</R,22,0><COMMA,23,0>,</COMMA,23,0> <C,24,0>dar</C,24,0> <V3,25,0>ar fi
nevoie</V3,25,0>    <S,26,0>de</S,26,0>    <PI,27,0>oarecari</PI,27, 0>
<NPN,28,0>precizari</NPN,28,0> <CR,29,0>si</CR,29,0> <V1,30,0>simt</ V1,30,0>
<C,31,0>ca</C,31,0> <QZ,32,0>nu</QZ,32,0> <PPSD,33,0>mi</PPSD, 33,0><Z,34,0>-
</Z,34,0><V3,35,0>ar    ajunge</V3,35,0>    <NSRY,36,0>respiratia</
NSRY,36,0><COMMA,37,0>,</COMMA,37,0> <C,38,0>ca</C,38,0> <V1,39,0>as
ocoli</V1,39,0>    <R,40,0>prea</R,40,0>    <R,41,0>mult</R,41,0><POINT,42,0>.</
POINT,42,0>
```

Ex.5.3.7.Mar. (întindere de text între paranteze acolade {...})

[As vrea]₁ [sa-ti spun]₂ [ca si mai târziu, {si într-alta parte,} am vazut]₃ [ca lucrurile se întâmpla tot asa,]₄ [dar ar fi nevoie de oarecari precizari si simt]₅ [ca nu mi-ar ajunge respiratia,]₆ [ca as ocoli prea mult.]₇

Ex.5.3.7.SCD.

[As vrea]₁ [sa-(ti) spun]₂ [ca si (mai (târziu)) , si (într-(alta (parte))) , am vazut]₃ [ca (lucrurile) se întâmpla (tot asa)]₄ ,[dar ar fi nevoie (de (oarecari (precizari)))]₅ si[simt]₆ [ca nu (mi)-ar ajunge (respiratia)]₇ ,[ca as ocoli (prea (mult))].]₈

5.4. Programele *TexTag* și *ClauSEGM*

În cele ce urmeaza sunt prezentate câteva imagini de executie în cadrul programelor *TexTag* si *ClauSEGM*, scrise în Visual C++ 5.0, si utilizate pentru a eticheta si segmenta texte de LN (limba româna). Figurile 5.4.1. si 5.4.2. se refera la *TexTag*, Figura 5.4.3. contine executia algoritmului de *segmentare M-1997* în cadrul *ClauSEGM*, iar Figura 5.4.4. contine o executie a algoritmului de *segmentare SCD-2002* sub mediul *ClauSEGM*. Stabilirea relatiilor de *dependentă inter-* si *intra-clauzale*, pentru aceleasi doua tipuri de algoritmi, urmeaza sa fie implementata în cadrul aceluiasi mediu *ClauSEGM*.

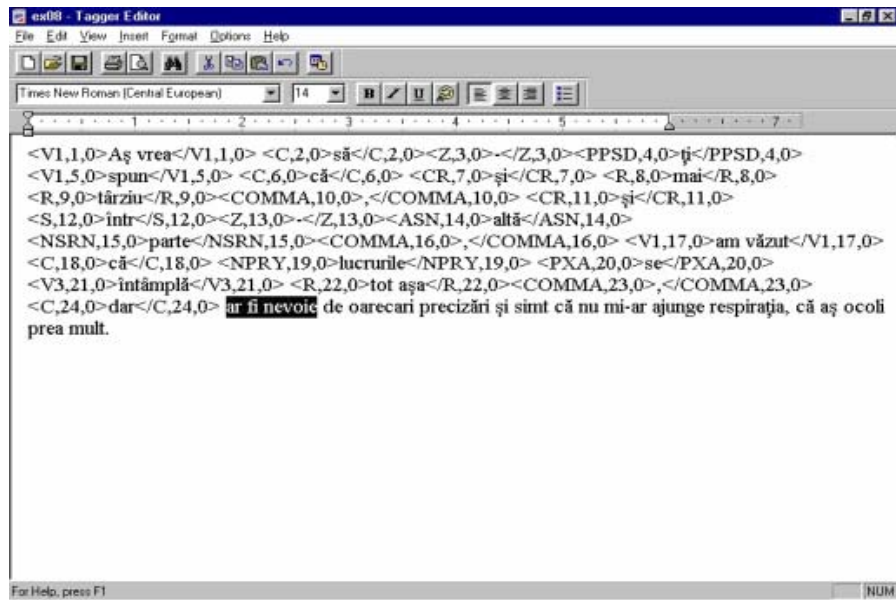


Figura 5.4.1. Rezultatul etichetării morfologice sub *TexTag*



Figura 5.4.2. Lista de etichete selectata cu un meniu din *TextTag*

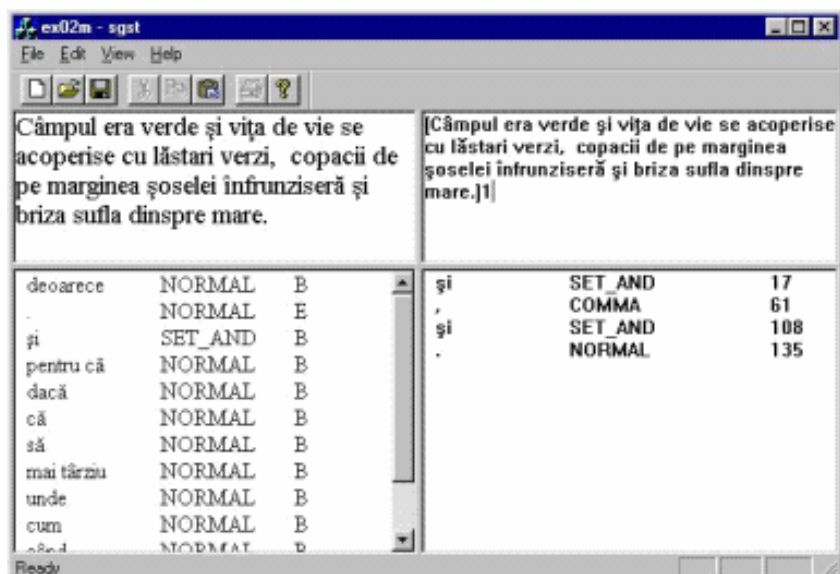


Figura 5.4.3. O execuție a algoritmului de segmentare *M-1997* cu programul *ClauSEGM*

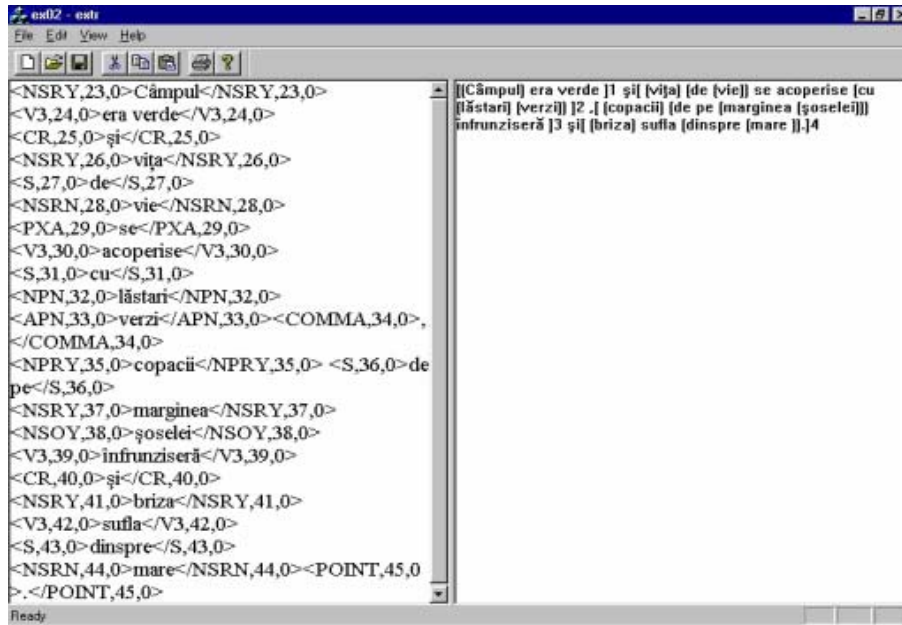


Fig. 5.4.4. O execuție a algoritmului de segmentare SCD-2002 sub mediul *ClauSEGM*

6. Concluzii

Rezultatele obținute în această lucrare nu se referă strict la compararea și implementarea celor doi algoritmi de segmentare. Avem, de fapt, două tipuri de algoritmi de segmentare (și dependentă), și fiecare din cele două tipuri reprezintă linii specifice de cercetare, cu importante consecințe asupra domeniilor de procesare a LN cărora se adresează: algoritmul *M-1997* este destinat (teoriei și) aplicațiilor de procesare a discursului, generare automată a LN, și rezumării automate, în timp ce algoritmul *SCD-2002* se încadrează mai curând în teoriile sintactice ale LN, cum sunt teoria *FX-bar* [3], parsarea bazată pe teorii (principii) sintactico-semantice ale LN, dar și punerea în evidență a structurilor (segmentelor) și relațiilor de discurs [6].

Demonstrarea relației (de scufundare) dintre cele două tipuri de algoritmi de segmentare, schițarea (în secțiunea 1) a unui *cadru formal general* pentru *algoritmii de segmentare* a LN, în particular a segmentării de tip *chunking*, propunerea (în cadrul algoritmilor-*SCD*) unei metode generale de segmentare în unități textuale a LN și de stabilire a dependențelor între ele, toate acestea constituie posibile noi perspective pentru

abordările teoretice și aplicative curente în procesarea automată a LN, inclusiv, și mai ales, pentru limba română.

Revenind la aspectele concrete expuse în acest articol, extinderea algoritmilor către analiza complexă a structurilor semantico-discursive antrenate de clasele de marcheri, și perfecționarea actualelor implementări rămân principalele direcții de continuare a prezentei abordări.

Referințe bibliografice

- [1] Neculai Curteanu (1994). *From Morphology to Discourse Through Marker Structures in the SCD Parsing Strategy*, Language and Cybernetics, Akademia Libroservo, Prague, p. 61-73.
- [2] N. Curteanu, G. Holban (1996). *Strategia lingvistică SCD aplicată la analiza și generarea limbii române*, Limbaj și Tehnologie (D. Tufis, Ed.), Editura Academiei Române, p. 169-176.
- [3] Neculai Curteanu (2000). "Towards a Functional X-bar Theory", Technical Report, Institute of Theoretical Informatics, Romanian Academy, Iasi Branch, 32p.
- [4] Daniel Marcu (1997). *The Rhetorical Parsing, Summarization, and Generation of Natural Language Texts*, Ph.D. Thesis, Univ. of Toronto, Canada, 331 p.
- [5] Daniel Marcu (2000). *The Theory and Practice of Discourse Parsing and Summarization*, The MIT Press, Cambridge.
- [6] O. Popârda, N. Curteanu (2002). *L'évolution du discours juridique français analysé par la stratégie linguistique SCD*, LINCOS Studies in Theoretical Linguistics, Lincom Europa (va apărea).
- [7] N. Curteanu, C. Lintes (2002). *Segmentation Algorithms for Clause-Type Textual Units*, Research Report, Institute of Theoretical Informatics, Romanian Academy.
- [8] N. Curteanu, D. Cristea, P. Mihaescu (1982). *Cercetări în domeniul comunicării om-calculator prin intermediul limbajului natural*. Contract de cercetare nr. 4774/1982, Universitatea Iasi - ICI București.
- [9] Neculai Curteanu (1983). *Algoritmi de analiză sintactică a frazei și propoziției românești*. Lucrările Conferinței INFO-IASI'83, p. 553-548.
- [10] N. D. Cristea, N. Curteanu, P. Mihaescu (1983). *Implementarea analizorului morfologic și definitivarea proiectului de analiză sintactică*. Contract de cercetare nr. 1906/1983, Universitatea Iasi - ICI București.
- [11] N. Curteanu (1984). *Aspecte ale analizei logice a limbajului natural*. Contract de cercetare nr. 4709/1984, Universitatea Iasi - ICI București.
- [12] Rebecca Passonneau, Diane Litman (1997). *Intention-based segmentation: human*

-
- reliability and correlation with linguistic cues*, in Proc. 31th Annual Meeting of ACL, Ohio, p. 148-155.
- [13] Lance Ramshaw, Michel P. Marcus (1999). *Text Chunking Using Transformation-based Learning*, in (S. Armstrong *et al.*, Eds.) "Natural Language Processing Using Very Large Corpora", Kluwer Acad. Publ., p. 157-176.
- [14] Victor Raskin, S. Nirenburg (1999). "Lexical Rules for Deverbal Adjectives", in E. Viegas (Ed.) *Breadth and Depth of Semantic Lexicons*, Kluwer Acad. Publ., p. 99-119.
- [15] M. Johnson, Federica Busa (1999). "Qualia Structure and Compositional Interpretation of Compounds", in E. Viegas (Ed.) *Breadth and Depth of Semantic Lexicons*, Kluwer Acad. Publ., p. 167-186.
- [16] Denis Bouchard (2001). *La source sémantique des facteurs hétérogènes qui régissent la distribution des adjectifs*, Conferinta Internationala "Representations du Sens Linguistique", Bucuresti.
- [17] Dumitru Irimia (1997). *Morfo-sintaxa verbului românesc*. Editura Universitatii "Al. I. Cuza", Iasi.
- [18] Eva Hajicova, H. Skoumalova, P. Sgall (1995). *An Automatic Procedure for Topic-Focus Identification*. Computational Linguistics, 21(1): 81-94.
- [19] P. Sgall, E. Hajicova, J. Panevova (1986). *The Meaning of the Sentence in Its Semantic and Pragmatic Aspects*. Kluwer Academic Publishers, Dordrecht.
- [20] Neculai Curteanu (1988). *Augmented X-bar Schemes*. COLING'88 Proceedings, Budapest, p. 130-132.
- [21] Dan Tufis (2000). *Using a Large Set of EAGLES-compliant Morpho-Syntactic Descriptors as a Tagset for Probabilistic Tagging*, in Proceedings of the LREC'2000 International Conference, Athens.
- [22] Dan Tufis, A.M. Barbu (2001). *Computational bilingual lexicography: automatic extraction of translation dictionaries*, In Romanian Journal on Information Science and Technology, vol. 4, no. 3.

O metoda automata pentru inserarea diacriticelor în texte în limba româna

Rada F. MIHALCEA
University of Texas at Dallas, Richardson, Texas, U.S.A.
rada@utdallas.edu

Vivi A. NASTASE
University of Ottawa, Ottawa, Canada
vnastase@site.uottawa.ca

1. Introducere

Problema restaurării diacriticelor constă în inserarea diacriticelor într-un text în care lipsesc. Creșterea continuă a numărului de texte disponibile prin Internet face ca metodele automate de inserare a diacriticelor să devină o componentă esențială în multe aplicații importante, cum ar fi extragerea de informații, traducerea automată, colecționarea de texte, construirea dictionarelor electronice și multe altele. Corectarea erorilor ortografice poate să aibă un impact major asupra calității rezultatelor obținute în aceste aplicații. De exemplu, în absența unei metode de restaurare a diacriticelor, unele cuvinte devin ambigue, cum este cazul cuvintelor din limba română *peste*, *peste* sau *paturi*, *paturi*. O căutare bazată pe astfel de cuvinte poate returna multe texte irelevante (de exemplu, o căutare pentru *peste* ar returna și documente conținând *peste*). De asemenea, traducerea unor astfel de cuvinte într-o limbă străină poate fi eronată (de exemplu, traducerea corectă a cuvântului *paturi* în limba engleză este *blankets*, dar în absența diacritice este tradus greșit ca și *beds*).

Metodele dezvoltate până în prezent pentru rezolvarea acestei probleme se bazează în general pe dictionare și pe diverse procesoare lexicale și/sau sintactice. Multe dintre limbile care se confruntă cu problema restaurării diacriticelor nu beneficiază însă de astfel de resurse, și ca urmare aplicabilitatea acestor metode este limitată la limbi bine studiate care dispun de suficiente resurse. Lucrarea de față prezintă o metodă automată de reînserare a diacriticelor în text care necesită doar o colecție de texte de dimensiuni modeste. Spre deosebire de alte metode dezvoltate anterior, metoda introdusă în această lucrare nu necesită nici un fel de dictionare sau procesoare morfologice și/sau sintactice, și prin urmare poate fi folosită pentru prelucrarea de texte în orice limbă care dispune de un număr minim de texte cu diacritice. Datorită lipsei de restricții, metoda propusă este foarte generală și ușor aplicabilă pentru orice limbă. Pentru a demonstra această afirmație, după

ce vom prezenta experimentele pentru texte în limba româna, vom arata câteva rezultate obținute pentru limbile ceha, poloneza și maghiara.

2. Experimente anterioare

Restaurarea diacriticelor nu este în sine o problema dificila. Experimente efectuate până în prezent au demonstrat ca folosirea de dictionare electronice poate duce la o acuratete de peste 90% în restaurarea accentelor pentru limbile franceza și spaniola [9],[11],[5]. Metoda descrisa de Michael Simard în [9] este o îmbunatatire adusa unei metode propusa anterior de El-Bêze [4]. Aceasta metoda se bazeaza pe Hidden Markov Models și învata folosind cuvintele învecinate. Precizia raportata este de 99%. Tufis și Chitu [10] propun o metoda similara pentru inserarea diacriticelor în texte în limba româna. Yarowsky prezinta în [11] un set de metode folosite pentru restaurarea accentelor în limbile franceza și spaniola. Majoritatea algoritmilor pe care îi prezinta se bazeaza pe dictionare și cuvinte învecinate pentru a decide asupra ortografiei potrivite pentru fiecare cuvânt ambiguu. Yarowsky compara N-gram taggers, clasificatoare Bayesiene și liste de decizii cu metoda de baza care consta în folosirea unui dictionar. Pentru cele doua limbi considerate în experimentele raportate, listele de decizii duc la performantele cele mai ridicate. Toate aceste tehnici se bazeaza însa pe context, dictionare, și în unele cazuri pe informatii aditionale de natura morfologica și sintactica. Nagy et al. prezinta în [7] o abordare diferita a problemei, în care siruri de litere sunt extrase din fiecare cuvânt și folosite pentru a obtine statistici. Folosind metoda propusa, s-a observat o precizie foarte buna obtinuta pe texte în limba franceza. Experimentele prezentate în [7] sunt asemanatoare cu cele raportate în [1], unde masuri de similaritate între trigrame sunt folosite pentru a automatiza corectarea greselilor de ortografie.

Majoritatea studiilor efectuate până în acest moment pe aceasta tema, s-au ocupat de limbi bine cunoscute și raspândite, cum ar fi franceza și spaniola. Foarte putine studii s-au concentrat pe limbi mai puțin mediatizate cum ar fi ceha, slovena, turca sau alte limbi care folosesc diacritice. Tabelul 1⁷⁷ prezinta diacriticele folosite în limbile europene cu alfabet latin. Dupa cum rezulta din aceasta lista, numeroase limbi se confrunta cu problema restaurarii diacriticelor. Din setul de 36 de limbi cuprinse în tabel, engleza pare sa fie singura limba pentru care diacriticele nu constituie o problema. Cuvintele din engleza care contin diacritice au fost împrumutate din alte limbi, și varianta acestora fara diacritice nu are un corespondent care sa duca la ambiguitate. Diacriticele par însa sa aiba un rol

⁷⁷ Tabelul cuprinde numai litere mici. Fiecarei litere mici îi corespunde o litera mare. Informatia din acest tabel a fost agregata din liste de diacritici în limbi europene, disponibile la adresa www.tiro.com/di_intro.html

important în diferențierea cuvintelor. Engleza, care după cum spuneam nu are diacritice pe se, are în schimb o ambiguitatea semantică mai ridicată⁷⁸

Tabel 1

Diacritice din limbile europene cu alfabet latin

Limba	Diacritice	Limba	Diacritice
Albaneză	ç ë	Malteză	c g h ż
Bască	ñ ü	Norvegiană	á ć ř
Bretonă	â ç ñ ũ ü	Olandeză	á f â ä é ç ę ë i î ð í ó ñ ô ö ú û ũ ü
Catalană	í ç č é í đ' ? ñ ó ú ü	Poloneză	ą ć ę ł ń ó ś ź ż
Cehă	á č đ' é e í ñ ó ř š ť ú ů ý ž	Portugheză	â ã ç ę ó ô õ ö
Daneză	Í ć ř	Română	â ä î ș ț
Engleză	None	Sami (Laponă)	á đ' č đ' n ŋ s t ž
Estoniană	ä ç õ ö š ü ž	Serbo-croată	ć č đ š ž
Faroeză	á ć đ' í ó ř ú ý	Slovacă	á ä ç đ' é í ľ' n ó ô ř š ť ú ý ž
Finlandeză	ä í ö š ž	Slovena	č š ž
Franceză	í â á ç ç é ę ë î đ' ô s ũ ü	Spaniolă	á é í ó ú ü ñ
Galițiană	á é í ó ú	Suedeză	ä å ö ø
Germană	ä ö ü ß	Turcă	ç ğ ö ş ü
Islandeză	á ć đ' é í ó ö ú ý ț	Sorbiană (1)	ć č e ł n ř ś š ž ž
Italiană	í é č í è đ' ó ñ ú ũ	Sorbiană (2)	ć č e ł n ó ř š ž
Maghiară	á é í ó ö ő ú ü ű	Welsh	â ę î ô ũ w y

Aplicabilitatea metodelor menționate anterior este limitată în următoarele cazuri:

1. Dictionarele electronice nu sunt disponibile, sau doar dictionare de dimensiuni relativ mici sunt făcute publice. Mai mult decât atât, în cazul în care dictionarul însuși nu are diacritice, metodele care se bazează pe această resursă pentru restaurarea diacriticelor devin inaplicabile.

⁷⁸ Studii efectuate pe corpusuri bilingve paralele, ar arătat că vocabularul construit dintr-un text în limba engleză este aproximativ jumătate din vocabularul construit pe baza aceluiași text într-o altă limbă. Competiția SENSEVAL [6] raportează de asemenea precizii mult mai mici pentru engleză comparativ cu alte limbi în rezolvarea ambiguității semantice. Lipsa diacriticelor în limba engleză ar putea constitui o explicație a acestui fenomen.

2. Procesoarele folosite pentru analiza morfologica si/sau sintactica, considerate folositoare pentru problema restaurarii diacriticelor, nu exista sau nu sunt public disponibile.
3. Numarul de texte disponibile continând diacritice este relativ mic. Marimea corpusurilor publice sau disponibile prin Internet influenteaza marimea vocabularului care poate fi construit ad-hoc pe baza acestor texte. În plus, majoritatea siturilor care publica texte pe Internet prefera în multe cazuri sa evite diacriticele din motive de simplitate, uniformitate, sau pur si simplu lipsa de mijloace necesare pentru codificarea diacriticelor.

Lucrarea de fata prezinta o metoda de restaurare a diacriticelor bazata pe învatarea la nivel de litera, si nu la nivel de cuvânt. Avantajul principal al acestei metode este faptul ca ofera posibilitatea de generalizare dincolo de cuvinte. Metoda este folositoare mai ales pentru limbile pentru care resursele disponibile sunt limitate, în speta limbi care nu au dictionare electronice mari cu diacritice. Limbi cunoscute si bine studiate, precum franceza si spaniola, pot de asemenea beneficia de aceasta metoda pentru procesarea cuvintelor necunoscute.

Experimentele prezentate în aceasta lucrare adreseaza în principal problema restaurarii diacriticelor în texte în limba româna. Precizia observata pe limba româna este de 99%, masurata la nivel de litera. Experimente similare au fost efectuate pe alte trei limbi, si anume poloneza, maghiara si ceha, de asemenea cu rezultate foarte bune. Avantajul principal al metodei este faptul ca nu necesita nici o etapa de preprocesare, ci numai un corpus relativ mic format din texte cu diacritice. Datorita simplitatii algoritmului, viteza de procesare este foarte mare, de aproximativ 20 pagini de text pe secunda, masurata pe un calculator cu un procesor Pentium III cu frecventa de 500MHz si 250MB memorie.

Practic, metoda propusa încearca sa învete reguli aplicabile la nivel de litera. În loc de a învata reguli care se aplica la nivel de cuvânt, cum ar fi „*anuncio se scrie anunció atunci când are functia de verb*”, dorim sa învatam reguli aplicabile la nivel de litera, cum ar fi „*s urmat de i si spatiu si precedat de spatiu se scrie s*”. Astfel de reguli, învatae la nivel de litera, sunt mai generale si au aplicabilitate mai mare, în special în cazurile în care dictionarele disponibile sunt de dimensiune redusa, când se întâlnesc multe cuvinte necunoscute în textul dat, sau când procesoare pentru analiza morfologica sau sintactica nu sunt la îndemâna.

Este evident ca în analiza limbajului literele constituie nivelul cu granularitatea cea mai scazuta, si de aceea au si cel mai mare potential de generalizare. În loc de aproximativ 150.000 de unitati candidate potentiale pentru algoritm (marimea aproximativa a vocabularului de uz general a unei limbi), vom avea mai mult sau mai putin 26 caractere pe baza carora se vor constitui datele de intrare pentru algoritmul de dezambiguare⁷⁹.

⁷⁹ Numărul de litere depinde de limba care se analizează. S-a arătat de exemplu că aproximativ 85% dintre cuvintele în limba franceză nu au o formă ortografică cu diacritice, și deci numai 20.000 de cuvinte sunt potențial ambigue. Pe de altă parte, numai 7 litere sunt ambigue în limba franceză.

3. Experimente

Scopul experimentelor descrise în această lucrare este de a arata ca învățarea la nivel de literă este posibilă și poate rezolva, cu precizie mare, problema restaurării diacriticelor. Pe lângă faptul că metoda propusă constituie o problemă de cercetare, scopul învățării la un nivel de granularitate atât de scăzut este de a oferi o metodă viabilă pentru limbile pentru care resursele lexicale și semantice disponibile sunt limitate, și pentru care restaurarea diacriticelor prin învățare la nivel de cuvânt este greu de realizat.

3.1. Date

Prezentăm în primul rând experimentele efectuate pe texte în limba română. Limba română nu este o limbă foarte răspândită, și în consecință nu are foarte multe resurse publice disponibile pentru pre-procesare, iar dicționarele electronice sunt de dimensiuni relativ mici. În al doilea rând, am avut de rezolvat o problemă specifică de restaurare a diacriticelor într-un dicționar electronic român-englez care conține aproximativ 75.000 de cuvinte, dar are dezavantajul că diacriticele lipsesc. Am considerat că este avantajos să studiem problema restaurării diacriticelor și să folosim acest dicționar, în loc să ne bazăm pe alte dicționare cu diacritice de dimensiuni reduse. În plus, pentru procesoarele pe care am dori să le dezvoltăm pentru limba română avem nevoie de numeroase texte electronice în limba română. De obicei aceste texte nu au diacritice, și deci reinserarea diacriticelor este din nou necesară. Avem de asemenea posibilitatea de a compara eficacitatea acestei metode cu rezultate obținute în experimente efectuate pe aceeași limbă constând în metode în care învățarea se face la nivel de cuvânt [10].

Pentru a aplica metoda descrisă în lucrarea de față, avem deci nevoie de o colecție de texte românești cu diacritice. În acest scop, am colectat articole din „România Literară”⁸⁰, un ziar românesc publicat săptămânal, cu articole legate în general de literatură. Ziarul are o versiune care conține diacritice începând din anul 2000. Colecția disponibilă on-line la data colectării datelor (august 2001) cuprindea 2780 articole. În pasul următor, textul a fost extras din fișierele HTML. Atenție deosebită a fost acordată doar caracterelor românești. Alte caractere cu diacritice întâlnite ocazional, cum ar fi *c, é, etc.* au fost transformate în forma lor echivalentă, fără diacritice, având în vedere că suntem interesați doar de caracterele românești, și nu de caractere franceze sau din alte limbi. După toate aceste faze premergătoare, am obținut un corpus conținând aproximativ 3 milioane de cuvinte.

Literele mari au fost transformate în litere mici. Cazul literelor *â* și *î* este special în limba română: deși pronunția lor este identică, folosirea lor este guvernată de reguli bazate pe poziția lor în cuvânt. La începutul cuvântului se folosește întotdeauna *î*, iar *â* se folosește în interiorul cuvântului. Este bine cunoscut faptul că folosirea acestor litere a fost controversată de-a lungul timpului. O lege din anii '60 a schimbat ortografierea de la *â* la *î*, singura excepție fiind cuvintele derivate din rădăcina *român*. La începutul anilor '90 ortografia veche a fost reintrodusă, și astfel s-a ajuns la cazuri de texte inconsistente, în care se întâlnesc scrieri

⁸⁰ Accesibil prin <http://www.romlit.ro>

diferite ale aceluiași cuvânt. De exemplu, *cîntec* și *cântec* sunt forme ale aceluiași cuvânt care pot fi întâlnite în același text. Ziarul „*România Literară*” păstrează încă ortografia cu *î*, cu mici excepții (de exemplu, articole scrise de scriitori invitați care preferă să scrie folosind *â* în loc de *î*).

3.2. Algoritmi de învățare

Pentru a rezolva problema restaurării diacriticelor, am ales să folosim un algoritm bazat pe învățarea de instanțe (IBL). Există două motive importante care au stat la baza luării acestei decizii. În primul rând, este faptul demonstrat că excepțiile au un rol important în procesarea limbajelor naturale. Algoritmii de tip IBL sunt recunoscuți pentru faptul că iau în considerare fiecare exemplu de antrenament în luarea unei decizii de clasificare [2], și deci folosirea acestui tip de algoritmi prezintă un avantaj deosebit în probleme de limbaj natural. În al doilea rând, acest gen de algoritmi sunt foarte eficienți relativ la timpul de antrenament și testare.

Învățarea pe baza de instanțe se desfășoară în felul următor: în pasul de antrenament, toate exemplele de intrare sunt memorate. În faza de testare, fiecare exemplu din set este comparat cu exemplele memorate, și va primi clasificarea dată de exemplul memorat de care este cel mai apropiat, distanța fiind dată de măsura specifică aleasă în implementarea folosită. Pentru efectuarea experimentelor propuse, am folosit implementarea TiMBL [3] a acestor algoritmi. În plus, am efectuat experimente asemănătoare și cu un clasificator pe baza de arbori de decizie, și anume C4.5 [8]. Arborii de decizie sunt construiți din setul de exemple de antrenament. La fiecare pas este ales un atribut care discriminează cel mai bine exemple din clase diferite (prin valorile sale). Grupele obținute prin diviziunea după acest atribut vor fi din nou împărțite în grupe mai mici și mai pure, prin alegerea unui nou atribut care discriminează cel mai bine exemplele din grupă. Acest proces continuă până când grupele obținute au un grad de puritate acceptabil, sau mărimea arborelui depășește un prag ales inițial. Rezultatele obținute folosind C4.5 sunt asemănătoare cu cele obținute folosind TiMBL, însă C4.5 are capacitatea de a genera reguli expresive, folosite pentru implementări practice.

Având în vedere că lucrăm la nivelul literelor, atributul care trebuie învățat este constituit de litera ambiguă. Acesta poate fi oricare din literele ambigue enumerate în Tabelul 1. Pentru limba română avem 4 perechi de litere ambigue: *s - ș*, *t - ț*, *a - ă*, *i - î*. Literele mari au fost convertite în prealabil în litere mici. Datorită faptului că datele folosite aplică ortografia cu *î*, nu avem ambiguitatea *a - ă*, ci doar ambiguitatea *i - î*. Aceasta nu implică însă o pierdere de generalitate. Conversia între cele două forme de ortografie este simplă și se poate realiza folosind doar poziția literei în cuvânt, și prin urmare scrierile diferite nu afectează rezultatul algoritmului.

3.3. Atribute

Atributele folosite în orice algoritm de învățare au un impact foarte mare asupra eficacității algoritmului. După cum am menționat și în introducere, nu avem posibilitatea de a folosi procesoare care determină partea de vorbire a cuvintelor, și nici un alt fel de analize morfologice sau sintactice. În plus, nu dorim să ne bazăm pe cuvintele

încecate, deoarece avem un numar limitat de date, si în consecinta exista sansa de a întâlni un numar mare de cuvinte necunoscute. Prin urmare, ne-am decis asupra folosirii unor atribute foarte simple, pentru extragerea carora nu este nevoie de nici un fel procesare speciala. Vom folosi litere învecinate, cu o notatie speciala atribuita spatiilor, virgulelor si punctelor (aceste caractere pot afecta procesul de învățare, fiind considerate caractere speciale de catre C4.5 si/sau TiMBL).

Daca X este litera a carui ambiguitate trebuie rezolvata, atributele folosite sunt N litere la stânga si la dreapta literei ambigue:

$$L_{-N}, L_{-(N-1)}, \dots, L_{-1}, X, L_1, L_2, \dots, L_{(N-1)}, L_N$$

Acest set de atribute se comporta surprinzator de bine, relativ la acuratete, dupa cum vom arata în cele ce urmeaza.

Dupa cum am mentionat anterior, am ales sa nu ne bazam pe nici un tag obtinut cu procesoare lexicale sau morfologice, ci doar pe informatia care se poate extrage din text neprelucrat. De asemenea, suntem interesati sa gasim posibilitati de generalizare, astfel încât un corpus limitat sa poata fi folosit pentru a genera reguli de reinserare a diacriticelor. În loc de a învăța reguli bazându-ne pe cuvinte, dupa cum s-a procedat până acum, dorim sa învățam reguli bazate pe litere, pentru ca acestea constituie cele mai mici unitati în limbaj, si ofera posibilitatea învățarii chiar si dintr-o colectie mica de texte.

Pentru fiecare pereche ambigua de litere, parcurgem textul si generam toate exemplele posibile întâlnite în text. Atributele într-un exemplu sunt formate folosind N litere la stânga si la dreapta literei ambigue, si atributul tinta este însasi litera ambigua. Forma generala a exemplelor generate este:

$$L_{-N}, L_{-(N-1)}, \dots, L_{-1}, L_1, L_2, \dots, L_{(N-1)}, L_N, X$$

unde ca si în exemplul anterior, X este litera ambigua. Prezentam mai jos exemple de vectori de atribute care constituie date de intrare pentru algoritmul de învățare pentru rezolvarea ambiguitatii perechii $s - s$. CO, DO si SP sunt codurile care înlocuiesc virgula, punctul si spatiul.

$$\begin{aligned} & l, i, n, SP, (, u, b, SP, i, n, s. \\ & e, CO, SP, r, o, -, g, a, r, d, \$. \\ & g, a, r, d, i, t, u, l, CO, SP, s. \\ & e, SP, o, r, a, DO, SP, t, o, t, \$. \end{aligned}$$

Învatarea se reduce la detectarea corelatiilor între valorile atributelor care caracterizeaza exemplele de antrenament si valorile atributelor tinta, si utilizarea acestora pentru stabilirea valorii atributului tinta din exemplele de testare.

Numarul de exemple extrase din corpus depinde de perechea de litere. Din întregul set de 3 milioane de cuvinte, am obtinut 2.161.556 exemple pentru perechea ambigua $a - a$, 2.055.147 pentru perechea $i - i$, 1.257.458 exemple pentru $t - t$, si în final 866.964 exemple pentru perechea $s - s$. În fiecare din aceste cazuri, spatiul exemplelor este împartit în doua

clase, date de cele 2 variante ale literei ambigue. Metoda de învățare automată va folosi atributele date pentru a găsi reguli de clasificare a exemplurilor în cele 2 clase.

3.4. Rezultate

Precizia cea mai ridicată s-a obținut pentru o fereastră de 10 litere în vecinătatea literei ambigue ($N = 5$). Data fiind această observație, am considerat ca este important să studiem mai în detaliu acest caz, și să determinăm ratele de învățare pentru cele 4 perechi de litere ambigue. Cu toate acestea, prezentăm rezultate pentru ferestre de diverse dimensiuni, pentru comparație.

Tabelul 2 arată rezultatele obținute pentru $N=5$. Preciziile raportate în acest tabel sunt obținute folosind algoritmul bazat pe învățarea de instanțe. Am condus experimente cu seturi de antrenament de diverse dimensiuni, variind de la 2.000.000 exemple până la 10 exemple, pentru a determina rata de învățare și dimensiunea minimă a corpusului necesară pentru a obține o precizie satisfăcătoare. În toate aceste experimente s-au folosit seturi de testare conținând 50.000 exemple. Pentru a obține rezultate cât mai acurate am folosit validare încrucișată folosind 10 seturi diferite de test. Tabelul indică de asemenea baza de comparație, definită aici ca fiind precizia obținută când se folosește implicit litera cea mai frecventă din fiecare pereche ambiguă.

Rezultatele prezentate în Tabelul 2 sunt reprezentate grafic în Figura 1. Este interesant de observat că cea mai importantă fază a procesului de învățare are loc când se folosesc primele 10.000 exemple. În conformitate cu măsurătorile efectuate, a rezultat că aproximativ 100.000 – 250.000 caractere (aproximativ 25-60 pagini de text) sunt necesare pentru a genera 10.000 exemple cu diacritice, ceea ce constituie un corpus de dimensiune relativ mic. Mai departe, pentru a obține îmbunătățiri de numai 1% este necesar un număr semnificativ de exemple. Tabelul 2 indică de asemenea, în caractere groase, prima precizie care depășește baza de comparație, ca o indicație a dimensiunii minime a setului de antrenament pentru care se observă o formă minimă de învățare. După cum se observă din tabel, și numai 1.000 exemple sunt suficiente pentru învățare.

Tabel 2

Rezultate obținute în rezolvarea ambiguității literelor cu diacritice în limba română

	Pereche ambiguă				
	a - a	a - a(2)	i - î	s - ș	t - ț
Nr. total exemple	2.161.566	1.369.517	2.055.147	866.964	1.157.458
Baza comparație	74.70%	85.90%	88.20%	76.53%	85.81%
Exemple de Antrenament	Precizie obținută pe date de test (50.000 exemple)				
2,000,000	96.14%	-	99.69%	-	-
1,000,000	95.10%	99.14%	99.58%	-	98.75%
750,000	94.83%	98.97%	99.53%	99.07%	98.63%
500,000	94.57%	98.79%	99.46%	98.86%	98.40%

250,000	94.00%	98.37%	99.28%	98.87%	98.26%
100,000	93.03%	97.56%	98.96%	98.54%	97.81%
50,000	92.10%	96.86%	98.57%	98.13%	97.40%
25,000	90.99%	95.75%	98.11%	97.58%	96.92%
10,000	88.99%	93.75%	97.31%	96.53%	96.20%
5,000	87.56%	92.76%	96.65%	95.61%	95.10%
4,000	86.91%	91.86%	96.49%	94.99%	94.53%
3,000	86.39%	90.99%	96.19%	94.18%	94.30%
2,000	85.81%	89.93%	95.49%	93.47%	93.56%
1,000	83.49%	88.36%	93.78%	92.31%	91.85%
500	80.61%	85.66%	93.07%	90.75%	89.74%
250	77.89%	83.17%	92.75%	87.41%	87.23%
100	74.80%	84.04%	91.41%	82.13%	84.46%
50	72.79%	82.73%	88.05%	86.53%	77.54%
25	72.45%	81.34%	88.15%	78.26%	78.52%
10	73.38%	85.90%	88.20%	75.88%	85.81%

Folosind întregul set de exemple extrase din corpus, rezolvarea ambiguitatii perechii $i - \hat{i}$ este aproape 100% corecta. Pentru aceasta diacritica, avem acum o instanta gresita din 300 instante, în timp ce baza de comparatie implica o instanta gresita din fiecare 8 instante, deci o îmbunatatire semnificativa.

Cel mai slab rezultat este obtinut în cazul perechii $a - \hat{a}$. Dupa o analiza a rezultatelor, reiese ca principalul motiv care cauzeaza aceasta precizie scazuta este faptul ca multe substantive în limba româna au forma nearticulata terminata în a si forma articulata terminata în \hat{a} . De exemplu, *masa* si *masa* reprezinta forma articulata si respectiv nearticulata a substantivului *masa*. De asemenea, timpuri diferite ale aceluasi verb se disting numai prin terminatia în a sau \hat{a} . Algoritmul de învățare este deci indus în eroare din cauza folosirii acestor litere în contexte identice. O solutie simpla consta în evitarea în procesul de învățare a exemplurilor care contin a sau \hat{a} la sfârșitul unui cuvânt. Rezultatele obtinute sub aceasta ipoteza simplificatoare sunt raportate în Tabelul 2, în coloana $a-\hat{a}(2)$. Dupa cum se arata în tabel, câștigul este de mai mult de 4% în precizie folosind doar aceasta conditie simpla (câștig care se traduce într-o reducere a erorii de 87%).

Am folosit de asemenea si algoritmul de învățare bazat pe arbori de decizie C4.5, cu aceleasi date de antrenament, fara a observa însa nici o îmbunatatire comparativ cu rezultatele raportate în Tabelul 2. Dezavantajul folosirii C4.5 pentru aceasta problema este faptul ca faza de învățare este mult mai lenta decât în cazul folosirii algoritmului TiMBL. Pe de alta parte, C4.5 are capacitatea de a genera reguli expresive. „Daca $L_1=e$ si $L_2=spatiu\ atunci\ s$ ” (99.5%), „Daca $L_1=t$ si $L_2=spatiu\ atunci\ s$ ” (98.7%), „Daca $L_{i-4}=p$ si $L_i=v$ si $L_1=t$ si $L_2=e$ atunci s ” (95.5%), sunt exemple de astfel de reguli. L_i denota o litera învecinata în pozitia i relativ la litera ambigua. Se observa ca aceste reguli nu tin cont de faptul ca literele folosite în clasificare apartin aceluasi cuvânt sau nu. Algoritmul de

învatare se bazeaza pur si simplu pe litere, indiferent de cuvântul caruia îi apartin. În consecinta, pseudo-omonimele (cum ar fi *peste* si *peste*), sunt adresate în mod egal de aceasta metoda, pentru ca algoritmul are capacitatea de a se extinde dincolo de cuvinte.

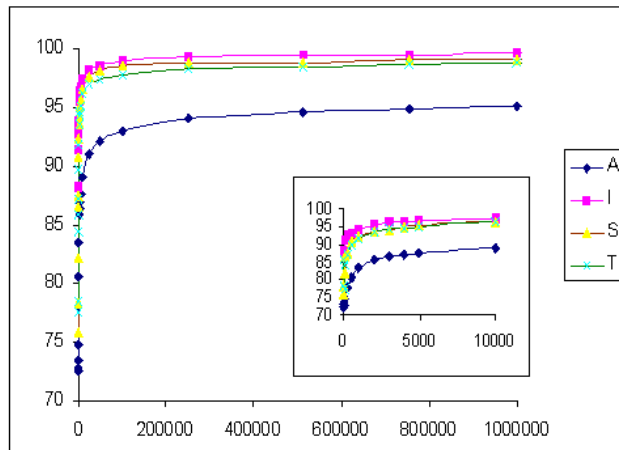


Figura 1. Rate de învățare pentru diacriticele în limba română.
din mijloc este o reprezentare marita a zonei 0-10.000

Graficul

3.5. Ferestre de dimensiune diferita

Am efectuat experimente cu ferestre de diverse dimensiuni, pentru a determina dimensiunea contextului care modeleaza cel mai bine problema noastra. Pentru aceasta, am considerat ferestre de dimensiune doi, sase, zece, patrusprezece si optsprezece litere învecinate (i.e. $N = 1,3,5,7,9$). Rezultate comparative sunt prezentate în Tabelul 3. Aceste numere trebuie comparate cu primul rând din Tabelul 2 (coloana corespunzatoare valorii $N=5$ în tabelul de fata).

Tabel 3

Rezultate comparative obtinute cu ferestre de dimensiuni diferite

Pereche ambiguă	Dimensiune fereastra				
	N=1	N=3	N=5	N=7	N=9
a - ă(2)	88.63%	98.79%	99.14%	99.10%	99.10%
i - î	94.18%	99.13%	99.69%	99.68%	99.43%
s - ș	88.09%	99.06%	99.07%	99.02%	99.00%
t - ț	89.45%	98.57%	98.75%	98.67%	98.25%

Când nu exista suficient context disponibil, o fereastra de dimensiune $N=3$ poate fi folosita fara a pierde mult din precizie. Însa, dupa cum am specificat si înainte, cea mai ridicata acuratete se obtine pentru o fereastra de zece litere înconjurate (N=5).

3.6. Comparatie cu experimente asemanatoare

Rezultatele prezentate în lucrarea de fata se pot compara cu rezultatele raportate de Tufis si Chitu [10], care au folosit tot limba româna în experimentele lor. Tufis si Chitu mentioneaza ca sarcina recuperarii diacriticelor în limba româna este mai dificila decât în alte limbi, deoarece în româna diacriticele sunt mai intens folosite. Dupa cum raporteaza în experimentele lor, numai 60% din cuvintele din limba româna nu au diacritice, comparat cu studii mentionate în [9] care arata ca aproximativ 85% dintre cuvintele limbii franceze se scriu fara accent.

Abordarea prezentata de Tufis si Chitu foloseste dictionare, un analizor morfologic, iar învatarea se face la nivel de cuvinte. Folosind aceste resurse, au obtinut o precizie globala de 97.4%. Nu putem efectua o comparatie directa a rezultatelor noastre, având în vedere ca atât metodele, cât si modul de evaluare, sunt fundamental diferite. Precizia medie de 99% pe care noi o raportam este masurata la nivel de litera, pe când acuratetea raportata in [10] este determinata la nivel de cuvânt⁸¹.

Metodologia noastra depaseste abordarile anterioare, prin faptul ca s-au obtinut precizii si viteze de procesare ridicate fara a folosi nici un fel de resurse aditionale cum ar fi procesoare pentru analiza morfologica sau sintactica sau dictionare. Din aceste motive, algoritmul se poate aplica oricarei limbi, singura cerinta fiind un corpus relativ mic de texte cu diacritice.

4. Alte limbi

Pentru a demonstra generalitatea algoritmului pa care l-am propus, am condus experimente pe texte în alte trei limbi europene care fac uz de diacritice: ceha, poloneza si maghiara. Limbile considerate pentru aceste experimente sunt limbi cu raspândire restrânsa, pentru care resursele publice sunt limitate.

Pentru fiecare dintre aceste limbi am colectat texte cu diacritice disponibile prin Internet. Principalele surse folosite pentru formarea setului de date sunt dupa cum urmeaza: (1) pentru ceha, am folosit arhiva ziarului *Lidovky* si texte literare de *Kafka*, *Hašek* si *Capek*; (2) pentru maghiara, arhiva furnizata de catre *Digitális Irodalmi Akadémia* si un roman de *Petőfi Sándor*; (3) pentru poloneza, arhiva ziarului *Wiedza i zycie*. Pe lângă acestea, am mai folosit texte aditionale colectate de pe diverse situri, astfel încât sa obtinem un corpus de minim un milion de cuvinte pentru fiecare limba. Asemănător cu procesarea aplicata limbii române, datele au fost convertite în fisiere text, iar literele mari au fost transformate în litere

⁸¹ Diferența dintre precizia raportată la nivel de literă și precizia raportată la nivel de cuvânt rezultă practic din diferența de granularitate dintre litere și cuvinte. Presupunând că un cuvânt conține L litere ambigue, o singură literă din acest set L a cărui ambiguitate este rezolvată greșit face ca întreg cuvântul să fie considerat greșit, pe când la nivel de litere avem doar o singură eroare din setul L . Pe de altă parte, chiar dacă mai multe litere din setul L sunt rezolvate greșit, avem tot o singură eroare la nivel de cuvânt, dar mai multe erori la nivel de literă. Nu este deci foarte clar care ar fi modalitatea corectă de a compara aceste două metode care lucrează la nivele de granularitate diferite.

mici. În urma acestei etape de pre-procesare, am obținut un corpus de 1.46 milioane cuvinte pentru ceha, 1.72 milioane cuvinte în maghiara și 2.5 milioane cuvinte în poloneza.

Algoritmii de învățare și atributele folosite în procesul de învățare sunt identice cu cele folosite în experimentele efectuate pe limba română, raportate în detaliu în secțiunea precedentă. Tabelul 4 prezintă rezultatele obținute pentru cele trei limbi. Pentru fiecare set de litere ambigue, sunt prezentate în tabel: (1) numărul de exemple obținute din corpusul limbii respective, (2) baza de comparație, măsurată ca fiind precizia ce se poate obține dacă pentru fiecare set ambiguu se folosește implicit litera cu frecvența de apariție cea mai ridicată, și (3) precizia obținută prin aplicarea metodei propuse în lucrarea de față.

Media obținută pentru toate patru limbile studiate (cele trei limbi a căror rezultate sunt prezentate în Tabelul 4, și limba română) este de 98.17%. Precizia medie măsurată pe fiecare limbă în parte este influențată de mărimea setului de date folosit. Textele colectate pentru ceha și maghiara conțin aproximativ 1.4-1.7 milioane cuvinte, și prin urmare precizia obținută în aceste două limbi este mai joasă decât pentru poloneza și română, pentru care am reușit să colectăm un corpus de 2.5-3 milioane cuvinte. Estimăm deci posibilitatea creșterii preciziei ca urmare a creșterii dimensiunii corpusului de antrenament.

Tabel 4

Rezultate obținute în restaurarea diacriticelor în trei limbi europene

Set litere ambigue	Număr exemple	Baza comparație	Metodă propusă
Cehă			
a á	649,886	75.01%	96.96%
c è	217,570	72.21%	97.08%
d d'	271,070	99.05%	99.86%
e é e	768,051	74.59%	97.02%
i í	504,298	60.43%	96.29%
n ò	439,552	98.97%	99.71%
o ó	566,521	99.08%	99.86%
r ø	319,352	65.55%	97.60%
s š	380,805	84.44%	98.88%
t ť	387,214	99.05%	99.85%
u ú ù	264,408	80.89%	93.51%
y ý	191,317	65.55%	95.06%
z ž	219,082	66.49%	98.70%
Medie			97.83%
Maghiară			
a á	1,198,294	73.51%	96.91%
e é	1,306,944	76.34%	96.40%
i í	647,137	89.14%	99.49%
o ó ö õ	678,012	71.15%	96.10%
u ú ü û	207,753	56.00%	97.31%

		Medie	97.04%
Poloneză			
a ¹	1,387,019	88.83%	97.07%
c æ	657,669	91.50%	99.42%
e ê	1,305,584	89.23%	98.47%
l ³	506,041	59.29%	98.80%
n ñ	878,824	96.75%	99.85%
o ó	1,230,389	88.67%	99.87%
s œ	688,677	88.67%	99.83%
z Ÿ ĭ	896,909	86.26%	99.73%
		Medie	99.02%

Este interesant de observat ca numarul de diacritice într-o limba nu influenteaza precizia medie obtinuta. Precizia care se obtine în cazul limbii maghiare, care are un total de 5 seturi de litere ambigue, este mai scazuta decât precizia care se obtine pentru limba ceha, care are un numar impresionant de diacritice (treisprezece). Si aceasta cu toate ca datele colectate pentru limba maghiara sunt mai numeroase decât datele colectate pentru limba ceha.

5. Concluzii

Am descris în lucrarea de fata o metoda de restaurare a diacriticelor bazata pe tehnici de învățare la nivelul de litera. Avantajul principal al metodei consta în capacitatea ei de generalizare dincolo de cuvinte. Nu este necesara nici un fel de analiza a textului, si nu se folosesc nici un fel de procesoare de limbaj sau dictionare. Singura cerinta este un corpus relativ mic de texte cu diacritice.

Metoda este folositoare în special pentru limbi pentru care nu sunt disponibile dicționare electronice de dimensiune adecvate, și nici procesoare pentru analiză morfologică și/sau sintactică. Mecanismul de învățare folosește date de intrare extrase din texte neprelucrate, și generează rezultate cu o precizie ridicată. Experimente detaliate efectuate pe texte în limba română au arătat că restaurarea diacriticelor în această limbă se poate efectua folosind metoda propusă cu o precizie de peste 99% la nivel de literă. Rezultatele au fost validate prin experimente efectuate pe alte trei limbi europene care fac uz de diacritice: cehă, poloneză și maghiară. Precizie medie măsurată pe cele patru limbi de studiu este de 98.14%, fapt care demonstrează că metoda este independentă de limbă. În plus, un alt avantaj al metodei este faptul că, datorită simplității sale, viteza de procesare este foarte mare, de până la 20 pagini de text pe secundă.

Referinte bibliografice

- [1] Angell, R., Freund G., Willett, P. Automatic spelling correction using a trigram similarity measure. *Information Processing and Management* 19, 4 (1983), 255-261.
- [2] Daelemans, W., van den Bosch, A., Zavrel, J. Forgetting exceptions is harmful in language learning. *Machine Learning* 34, 1-3 (1999), 11-34.
- [3] Daelemans, W., Zavrel, J., van der Sloot, K., van den Bosch, A. TiMBL: Tilburg memory based learner, version 4.0, reference guide. Tech. Rep., University of Antwerp, 2001.
- [4] El-Bcze, M., Merialdo, B., Rozeron, B., Derouault, A., Accentuation automatique des textes par des methodes probabilistes. *Techniques et sciences informatique* 16, 6 (1994), 797-815.
- [5] Galicia-Haro, S., Bolshakov, I., Gelbukh, A. A simple Spanish part of speech tagger for detection and correction of accentuation error. In *Text, Speech and Dialogue – Second International Workshop, TSD'99, September 1999, Proceedings (Plzen, Czech Republic, 1999)*, vol 1692 of *Lecture Notes in Computer Science*, Springer, pp. 219-222.
- [6] Kilgariff, A., Ed. , *Proceedings of SENSEVAL-2*, 2002.
- [7] Nagy, G., N., N., and Sabourin, M. Signes diacritiques: perdus et retrouvés. In *Actes du 1er Colloque International Francophone sur l'Écrit et le Document CIFED '98 (Quebec, Canada, 1998)*, pp. 404-412.
- [8] Quinlan, J. *C4.5: programs for machine learning*. Morgan Kaufman, 1993.
- [9] Simard, M. Automatic insertion of accents in French text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing EMNLP-3 (Granada, Spain, 1998)*.
- [10] Tufis, D., Chitu, A. Automatic diacritics insertion in Romanian texts. In *Proceedings of the International Conference on Computational Lexicography COMPLEX'99 (Pecs, Hungary, June 1999)*.
- [11] Yarowsky, D. Corpus-based techniques for restoring accents in Spanish and French texts. In *Natural Language Processing Using Very Large Corpora*. Kluwer Academic Publisher, 1999, pp 99-120.

Contributii privind structura statistica de cuvinte în limba româna scrisa^{*}

Adriana VLAD si Adrian MITREA
Universitatea "POLITEHNICA" din Bucuresti
Facultatea de Electronica si Telecomunicatii
B-dul. Iuliu Maniu, 1-3, Bucuresti, România
vadriana@vala.elia.pub.ro

1. Introducere

Aceasta lucrare apartine unui studiu mai larg dedicat de autori descrierii limbii române ca sursa de informatie. Punctul de plecare al acestui studiu a fost presupunerea generala conform careia limba naturala este bine aproximata de un lant Markov ergodic multiplu, cu ordin de multiplicitate mai mare decat 30, [1]. Descrierea acestei surse Markov multiple se realizeaza prin aproximatii succesive. Investigatia noastra statistica pâna în prezent a descris structura de litere, digrame, trigrame, tetragrame, precum si probabilitatile conditionate de o litera precedenta, [2]-[8].

Obiectivul principal ale prezentei lucrari este descrierea sursei de informatie fara memorie având ca simboluri cuvintele limbii române scrise. Aceasta presupune determinarea probabilitatii unui cuvânt (oricare ales), în caz ca aceasta probabilitate exista. Determinarea probabilitatii a însemnat implicit si o verificare a ipotezei de stationaritate a limbii române scrise pe baza structurii de cuvinte; verificarea s-a facut utilizând o procedura similara cu cea pe care am dezvoltat-o pentru m -grame, [3]-[8] (m -grama este o succesiune de m litere consecutive în texte naturale). Metoda noastra statistica de a determina probabilitatile cuvintelor a combinat urmatoarele tipuri de inferente statistice: teoria estimarii cu multiple intervale de încredere statistica; test al ipotezei ca probabilitatea apartine unui interval; test de egalitate între doua probabilitati.

Primele doua tipuri de inferente statistice mentionate (teoria estimarii cu multiple intervale de încredere statistica; test al ipotezei ca probabilitatea apartine unui interval) au folosite pentru a decide care este intervalul de încredere statistica "*reprezentativ*" pentru probabilitatea cuvântului investigat în textul natural. Simultan a aparut si o multime "*reprezentativa*" de date *i.i.d.* extrase din textul natural, corespunzatoare cuvântului

^{*} O parte din acest studiu s-a desfasurat în cadrul unui *Grant CNCSIS-MEC (2001-2002)* cu tema: "Descrierea limbii române scrise ca sursa de informatie"

investigat (modelul statistic *i.i.d.* presupune ca datele provin din variabile aleatoare independente statistic si identic distribuite).

Ultimele doua tipuri de inferente statistice mentionate (test al ipotezei ca probabilitatea apartine unui interval; test de egalitate între doua probabilitati) ca si intervalele de încredere statistica “*representative*” si multimile de date “*representative*” obtinute în prealabil au fost folosite pentru comparatii matematice între texte naturale. Aceste comparatii matematice (dincolo de valoarea lor ca atare) au avut scopul principal de a vedea daca putem vorbi de un model matematic al sursei de cuvinte pentru limba ca ansamblu, pe domenii ale limbii, pe autori, etc. Comparatiile s-au facut în doua moduri:

- urmarind probabilitatile unui cuvânt (acelasi) în texte naturale diferite;
- urmarind probabilitatile cuvintelor situate pe acelasi rang în texte naturale diferite (se compara probabilitatile asociate unui aceluasi rang în ierarhiile frecventelor relative).

Rezultatele experimentale au adus probe noi în sprijinul ipotezei de stationaritate a limbii române scrise în cadrul unui aceluasi domeniu punctând catre unele diferente între domenii diferite.

Investigatia noastra (atât privind intervalele de încredere statistica “*representative*”, cât si comparatia matematica dintre texte) a avut în vedere si eroarea statistica de ordinul al doilea. Acest tip de eroare are un rol special în dimensionarea unui nou corpus lingvistic care sa satisfaca acuratetea dorita pentru descrierea modelului matematic (sursa de informatie de cuvinte).

Lucrarea mai contine si un studiu experimental al uneia dintre cele mai cunoscute legi de tipul rang – frecventa, legea lui Zipf. Este analizat si un corolar al acesteia, de interes lingvistic.

Analiza experimentală s-a bazat pe corpusul lingvistic global pe care l-am alcatuit în prealabil pentru studiul structurilor de litere, digrame, trigrame si tetragrame (vezi spre exemplu [6]). Acest corpus este format din 93 de carti în limba româna, scrise cu noua ortografie (introdusa dupa 1993). Cartile reprezinta: literatura scrisa de autori români (11 carti: romane si nuvele), literatura straina tradusa în româna (47 de romane si nuvele), carti stiintifice (drept, medicina, silvicultura, istorie, sociologie, etc.) si altele. Au fost considerate doar cele 31 de litere ale limbii române (A A Â B C D E F G H I Î J K L M N O P Q R S S T T U V W X Y Z) precum si caracterul spatiu (blanc); orice alte simboluri (cifre, elemente de ortografie sau punctuatie) au fost eliminate (suprimate).

Rezultatele experimentale au fost obtinute pe diverse corpusuri organizate pe baza celor 93 de carti:

- Corpusul Mixt Global (#CMG) – obtinut prin concatenarea aleatoare a celor 93 de carti; acesta contine un numar de $L_c = 8806433$ cuvinte dintre care $N_c = 202403$ sunt distincte.

- Cele doua jumatați ale Corpusului Mixt Global: prima jumatațe (#1JCMG) și a doua jumatațe (#2JCMG); acestea contin un numar de $L_c = 4403217$ cuvinte și respectiv $L_c = 4403216$ dintre care $N_c = 148853$ și respectiv $N_c = 137845$ sunt distincte.
- Corpusul Literar Global (#CLG) – obtinut prin concatenarea aleatoare a 58 de carti (romane și nuvele scrise de autori români sau traduse în româna); acesta contine un numar de $L_c = 6255235$ cuvinte dintre care $N_c = 162124$ sunt distincte.
- Cele doua jumatați ale Corpusului Literar Global: prima jumatațe (#1JCLG) și a doua jumatațe (#2JCLG); acestea contin un numar de $L_c = 3127618$ cuvinte și respectiv $L_c = 3127617$ dintre care $N_c = 116247$ și respectiv $N_c = 116860$ sunt distincte.
- Corpusul Stiintific Global (#CSG) – obtinut prin concatenarea aleatoare a 11 de carti; acesta contine un numar de $L_c = 1049969$ cuvinte dintre care $N_c = 59093$ sunt distincte.

Au fost facute determinari atât pe o singura carte cât și pe grupuri de carti scrise de acelasi autor. Dintre acestea mentionam:

- #1. George Calinescu, *Bietul Ioanide*, Editura Minerva, Bucuresti, 1995, ISBN 973-21-0432-5 (vol. 1, ISBN 973-21-0431-7, pag. 1-214), (vol. 2, ISBN 973-21-0433-3, pag. 5-256), (vol. 3, ISBN 973-21-0434-1, pag. 5-238).
- #2. Radu Anton Roman, *Precum fumul*, Editura Cartea Româneasca, Bucuresti, 1996, ISBN 973-23-0274-7, pag. 5-283.
- #3. Radu Anton Roman, *Zile de pescuit*, Editura Metropol, Bucuresti, 1996, ISBN 973-562-073-1, pag. 11-302.
- #4. John le Carré, *Casa Rusia*, Editura Univers, Bucuresti, 1997, ISBN 973-34-0457-8, pag. 9-355.
- #5. John le Carré, *Spionul care venea din frig*, Editura Univers, Bucuresti, 1996, ISBN 973-34-0355-5, pag. 9-252, cu ortografie actualizata.
- #6. John Le Carré, *Micuta tobosareasa*, Editura Univers, Bucuresti, 1998, ISBN 973-34-0430-6, pag. 7-443, cu ortografie actualizata.
- #7. Alexandr Soljenitîn, *Arhipelagul Gulag*, Editura Univers, Bucuresti, (vol. I, 1997, ISBN 973-34-0454-3, pag. 7-432), (vol. II, 1997, ISBN 973-34-0480-2, pag. 5-474), (vol. III, 1998, ISBN 973-34-0497-7, pag. 5-414), cu ortografie actualizata, fara note.

Primul pas în analiza noastra a fost evaluarea frecventelor relative ale cuvintelor din corpusurile mentionate anterior. Tabelul 1 contine primele 55 de cuvinte din ierarhia frecventelor relative din diverse corpusuri.

Un alt rezultat experimental interesant este identificarea unui numar de 162 de cuvinte care se regasesc în toate cele 93 de carti ce alcatuiesc corpusul (fie ca este vorba de literatura, medicina, drept, etc.). Desi sunt doar 162, aceste cuvinte au o pondere importanta în textul global #CMG acoperind circa 45% din totalul celor 8806433 cuvinte. Aceste cuvinte comune împreuna cu rangul lor în ierarhie si frecventele lor relative în întreg textul #CMG sunt continute în Tabelul 2.

Tabel 1

Ierarhia frecventelor relative în câteva corpusuri
0. Rang; 1. Cuvânt; 2. Frecventa relativa (în %)

0	#CMG		#CLG		#1		#4+#5+#6		#6		#CŞG	
	1	2	1	2	1	2	1	2	1	2	1	2
1	de	4,10	de	4,02	de	4,17	de	4,17	de	4,12	de	4,87
2	şi	3,20	şi	3,12	şi	2,65	în	2,55	în	2,58	în	3,47
3	în	2,67	în	2,44	în	2,50	şi	2,39	şi	2,58	şi	3,07
4	să	1,62	să	1,87	cu	1,75	să	1,81	o	1,94	a	2,35
5	a	1,47	la	1,52	o	1,62	o	1,73	să	1,69	la	1,52
6	la	1,46	cu	1,50	a	1,47	la	1,55	cu	1,52	se	1,46
7	se	1,39	pe	1,45	la	1,43	cu	1,48	la	1,46	cu	1,21
8	cu	1,38	se	1,43	se	1,42	nu	1,41	se	1,44	care	1,17
9	o	1,30	o	1,41	pe	1,39	pe	1,39	pe	1,41	o	0,87
10	nu	1,28	nu	1,33	nu	1,37	se	1,35	nu	1,27	din	0,85
11	pe	1,27	a	1,17	să	1,33	un	1,18	un	1,25	pe	0,82
12	care	0,98	că	1,05	un	1,26	că	1,08	a	1,04	este	0,79
13	că	0,97	un	0,99	că	1,04	a	1,05	care	0,95	mai	0,75
14	mai	0,95	mai	0,97	lui	0,88	care	0,95	că	0,95	nu	0,73
15	din	0,91	din	0,94	mai	0,86	din	0,93	din	0,91	sau	0,71
16	un	0,87	care	0,89	din	0,86	ce	0,85	mai	0,84	să	0,70
17	ce	0,66	ce	0,69	care	0,84	mai	0,84	ce	0,79	pentru	0,67
18	ca	0,60	ca	0,58	ioanide	0,74	lui	0,68	pentru	0,71	că	0,54
19	pentru	0,54	lui	0,54	era	0,66	era	0,63	lui	0,71	al	0,53
20	lui	0,49	dar	0,51	ce	0,64	pentru	0,63	era	0,64	un	0,50
21	dar	0,45	era	0,51	e	0,54	ca	0,55	charlie	0,59	prin	0,44
22	fî	0,42	pentru	0,48	ca	0,53	el	0,53	ei	0,57	ca	0,43
23	este	0,42	fî	0,42	fî	0,52	dar	0,50	dar	0,54	fî	0,35
24	era	0,39	când	0,39	pompo- nescu	0,44	fî	0,49	ca	0,53	sunt	0,33
25	sau	0,35	el	0,38	pentru	0,40	îi	0,43	îi	0,53	ale	0,32
26	e	0,34	e	0,37	când	0,35	ei	0,42	ea	0,51	poate	0,29
27	el	0,34	am	0,35	el	0,33	ea	0,38	el	0,50	sa	0,29
28	al	0,33	ei	0,32	prin	0,27	când	0,34	fî	0,45	au	0,28
29	când	0,33	nici	0,30	am	0,26	nici	0,33	kurtz	0,36	ce	0,27
30	ei	0,29	îi	0,29	după	0,26	cum	0,31	când	0,34	art	0,27

31	am	0,28	mă	0,28	il	0,26	e	0,30	nici	0,32	fost	0,24
32	nici	0,28	cum	0,28	al	0,26	aŞa	0,30	cum	0,30	după	0,24
33	prin	0,28	sau	0,25	nici	0,25	dacă	0,29	iŞi	0,28	dacă	0,21
34	sa	0,26	fost	0,25	fără	0,25	il	0,29	il	0,28	c	0,19
35	sunt	0,25	după	0,24	ii	0,25	charlie	0,29	al	0,28	când	0,19
36	cum	0,25	sa	0,24	avea	0,24	fost	0,28	aŞa	0,26	m	0,18
37	fost	0,25	dacă	0,24	ar	0,23	barley	0,28	dacă	0,25	unei	0,17
38	dacă	0,24	al	0,24	dacă	0,23	al	0,27	e	0,25	cele	0,17
39	după	0,24	ea	0,23	gaittany	0,22	spuse	0,26	joseph	0,24	pot	0,16
40	au	0,23	aşa	0,22	însă	0,22	iar	0,26	sau	0,24	are	0,16
41	ii	0,23	il	0,22	foarte	0,21	iŞi	0,26	ai	0,23	penală	0,16
42	mă	0,21	iŞi	0,21	spre	0,20	este	0,26	după	0,23	trebuie	0,16
43	ea	0,21	este	0,21	ei	0,20	după	0,25	te	0,22	această	0,16
44	iar	0,19	au	0,21	aşa	0,20	am	0,25	sa	0,22	lui	0,16
45	poate	0,19	sunt	0,20	sunt	0,19	sau	0,24	fără	0,22	acest	0,16
46	aşa	0,18	iar	0,20	cum	0,19	ai	0,24	fost	0,22	iar	0,15
47	ar	0,18	fără	0,19	dar	0,19	ar	0,24	ar	0,22	lor	0,15
48	fără	0,18	prin	0,19	hagienuş	0,19	sa	0,23	este	0,21	numai	0,15
49	iŞi	0,17	ar	0,19	sau	0,18	te	0,20	iar	0,21	dar	0,15
50	il	0,17	le	0,18	fost	0,18	avea	0,20	le	0,21	mare	0,15
51	le	0,17	asta	0,18	toate	0,18	le	0,20	spuse	0,21	cel	0,14
52	ale	0,17	tot	0,18	este	0,18	leamas	0,20	apoi	0,20	unor	0,14
53	toate	0,17	eu	0,18	sa	0,18	timp	0,20	timp	0,20	fie	0,14
54	va	0,16	acum	0,17	iŞi	0,17	apoi	0,19	lor	0,19	va	0,14
55	decât	0,16	până	0,17	gonzalv	0,17	au	0,18	săi	0,19	între	0,13

Tabel 2.

Lista cuvintelor comune în toate cele 93 de carti

1. Cuvânt; 2. Rangul cuvântului în ierarhia frecvențelor relative în textul mixt global, #CMG; 3. Frecvența relativă a cuvântului în textul mixt global, #CMG, (în %)

	1	2	3	1	2	3	1	2	3	1	2	3
de	1	4,10	iar	44	0,19	unei	93	0,10	sar	185	0,05	
și	2	3,20	poate	45	0,19	atunci	94	0,10	una	187	0,05	
în	3	2,67	aşa	46	0,18	două	95	0,10	început	188	0,05	
să	4	1,62	ar	47	0,18	doar	96	0,10	încât	193	0,05	
a	5	1,47	fără	48	0,18	dintre	100	0,10	alte	196	0,04	
la	6	1,46	iŞi	49	0,17	are	101	0,10	acestea	198	0,04	
se	7	1,39	il	50	0,17	face	102	0,10	facă	199	0,04	
cu	8	1,38	le	51	0,17	sub	104	0,09	altă	200	0,04	
o	9	1,30	ale	52	0,17	nimic	106	0,09	același	204	0,04	
nu	10	1,28	toate	53	0,17	fel	107	0,09	deși	206	0,04	
pe	11	1,27	va	54	0,16	ia	108	0,09	fac	213	0,04	
care	12	0,98	decât	55	0,16	puțin	109	0,09	printre	220	0,04	

că	13	0,97	tot	56	0,16	între	110	0,09	pare	224	0,04
mai	14	0,95	lor	57	0,16	întrun	111	0,09	partea	225	0,04
din	15	0,91	spre	58	0,15	cea	112	0,09	afară	226	0,04
un	16	0,87	până	59	0,15	i	113	0,08	sus	240	0,04
ce	17	0,66	chiar	60	0,15	săl	116	0,08	faptul	246	0,03
ca	18	0,60	mult	61	0,14	aceea	117	0,08	locul	252	0,03
pentru	19	0,54	cel	63	0,14	ci	119	0,08	adevărat	260	0,03
lui	20	0,49	fie	65	0,14	față	126	0,08	tuturor	264	0,03
dar	21	0,45	ne	66	0,14	unul	127	0,08	măcar	266	0,03
fi	22	0,42	ai	67	0,14	astfel	128	0,08	primul	268	0,03
este	23	0,42	acum	68	0,14	parte	129	0,08	aceeași	275	0,03
era	24	0,39	trebuie	69	0,14	înainte	132	0,07	altfel	277	0,03
sau	25	0,35	cele	70	0,13	pot	134	0,07	nouă	299	0,03
e	26	0,34	numai	72	0,13	ele	138	0,07	acela	302	0,03
el	27	0,34	despre	73	0,12	totul	140	0,07	trebui	307	0,03
al	28	0,33	avea	74	0,12	dată	141	0,07	dintro	330	0,03
când	29	0,33	atât	75	0,12	toți	143	0,07	dă	358	0,02
ei	30	0,29	această	76	0,12	loc	144	0,07	afla	364	0,02
nici	32	0,28	putea	78	0,12	fiecare	153	0,06	rămâne	371	0,02
prin	33	0,28	unde	80	0,12	orice	155	0,06	alt	373	0,02
sa	34	0,26	într-o	81	0,11	spune	165	0,06	pus	377	0,02
sunt	35	0,25	acest	82	0,11	asemenea	166	0,06	întâi	387	0,02
cum	36	0,25	noi	84	0,11	sale	167	0,06	rând	397	0,02
fost	37	0,25	săi	86	0,11	acesta	168	0,06	alta	404	0,02
dacă	38	0,24	cât	87	0,11	lucru	169	0,06	legătură	415	0,02
după	39	0,24	mare	88	0,11	către	174	0,05	măsură	429	0,02
au	40	0,23	apoi	89	0,11	multe	175	0,05	rândul	601	0,01
îi	41	0,23	ceva	91	0,11	celor	178	0,05			
ea	43	0,21	însă	92	0,10	totuși	184	0,05			

2. Descrierea structurii statistice de cuvinte. Studiu bazat pe multiple intervale de încredere statistica

Fie un text natural **considerat ca succesiune de cuvinte** pe care îl esantionam cu o perioada suficient de mare astfel încât sa rupem practic dependenta dintre observatiile succesive. Initial în investigatia noastra statistica am considerat aceasta perioada ca fiind de 200 cuvinte. La fiecare moment de esantionare am înregistrat observatia facuta (cuvântul respectiv), conform Fig. 1. Multimea de date obtinute în acest fel contine N cuvinte unde $N=L_c/200$, unde L_c este lungimea textului în cuvinte.

CÂND GAITTANY AMINTI LUI ... GAITTANY TACU CACI STIA CA ... ÎNSA
GAITTANY LIPSIT DE

1.	CÂND	CACI	LIPSIT...
2.	GAITTANY	STIA	DE
...			
200.		TACU	GAITTANY ...

Figura 1. 200 de multimi de date (cuvinte) în model statistic *i.i.d.* obtinute prin esantionare periodica a textului natural

Deplasând originea esantionarii în textul natural apar 200 de astfel de multimi de date experimentale, fiecare în parte de volum N , Fig. 1.

Fiecare multime de N observatii astfel obtinuta satisface modelul statistic *i.i.d.*, model necesar în aplicarea inferentelor statistice utilizate. Independenta este asigurata de marimea perioadei de esantionare; distributia identica este un rezultat al ipotezei de stationaritate a limbii naturale.

Acceptând ipoteza de stationaritate a limbii, toate cele 200 de multimi de date experimentale (compatibile cu modelul *i.i.d.*) extrase din textul natural conform Fig. 1, trebuie sa contina aceeasi informatie despre probabilitatea cuvântului investigat (oricare ar fi acesta).

! Atentie însa, aceste multimi de date nu sunt independente între ele.

Un prim obiectiv al studiului nostru a fost de a vedea daca într-adevar cele 200 multimi de date confirma sau nu aceeasi probabilitate p teoretica (necunoscuta) a cuvântului investigat.

Un raspuns afirmativ ne-ar permite sa obtinem un model matematic pentru sursa de informatie de cuvinte asociata limbii române. Pentru a da un raspuns am extins o procedura statistica pe care am dezvoltat-o în [3]-[8] pentru m -grame. Prin aceasta procedura cele 200 de multimi de date experimentale se compara între ele aplicând repetat un test statistic al ipotezei ca probabilitatea apartine unui interval dat, vezi Anexa 1.

Mentionam ca nu am putut face o comparatie pe baza unui test mai des folosit, anume acela privind egalitatea între doua probabilitati, întrucât multimile de date care se compara nu sunt independente între ele.

Procedura a permis în final determinarea unui interval de încredere statistica optim care a fost denumit în continuare “*reprezentativ*” pentru cuvântul urmarit si textul natural. Simultan a aparut si multimea de date experimentale *i.i.d.* “*reprezentativa*” pentru cuvântul respectiv si textul natural, multime ce va fi folosita în comparatii matematice între texte naturale.

2.1. Intervale de încredere statistica “*reprezentative*” pentru probabilitatile cuvintelor. Metoda de determinare si rezultate experimentale

Scopul acestui subcapitol este de a determina probabilitatea p a unui cuvânt urmărit.

Fie m_i numărul de aparitii ale cuvântului în multimea i de date experimentale *i.i.d.* de volum N , $i = 1 \div 200$. (Aceste multimi sunt extrase din textul natural conform Fig. 1.)

Aplicând teoria estimării, fiecare din cele 200 de multimi de date conduce la o estimatie $\hat{p}_i = m_i/N$ a probabilitatii p necunoscute si la un interval de încredere statistica al probabilitatii $I_i = (p_{1,i}; p_{2,i})$, $i = 1 \div 200$. Considerând N suficient de mare astfel încât condiția de Moivre – Laplace sa fie satisfăcută, $Np(1-p) \gg 1$, limitele intervalului de încredere statistica (inferioara si superioara) se calculeaza conform relatiei (1), [9], [10]:

$$p_{1,i} \cong \hat{p}_i - z_{\alpha/2} \sqrt{\hat{p}_i(1-\hat{p}_i)/N} \quad p_{2,i} \cong \hat{p}_i + z_{\alpha/2} \sqrt{\hat{p}_i(1-\hat{p}_i)/N} \quad (1)$$

unde $z_{\alpha/2}$ este $\alpha/2$ cuantila legii normale de medie 0 si dispersie 1. În determinarile noastre experimentale am lucrat cu un nivel de încredere statistica de 95%; rezulta $z_{\alpha/2} = 1.96$.

Cu alte cuvinte putem spune ca probabilitatea adevarata p se afla în intervalul $I_i = (p_{1,i}; p_{2,i})$, cu o încredere statistica egala cu 0,95.

Într-o prima etapa a analizei noastre, pentru un anumit eveniment urmărit (aparitia unui cuvânt), s-au folosit urmatoarele marimi (a se vedea Fig. 2):

- p^* – frecventa relativa a cuvântului pe întreg textul natural considerat (ceea ce înseamna masurare din date corelate); p^* este raportul între numărul de aparitii ale cuvântului în textul natural si lungimea L_C a textului respectiv (numarul total de cuvinte). Se observa ca p^* este media aritmetica a celor 200 de estimatii. Subliniem ca p^* este o marime importanta pentru orice experimentator.
- $\hat{p}_{\min} = \min_i \hat{p}_i$, $i = 1 \div 200$ – valoarea minima a estimatiilor;
- $\hat{p}_{\max} = \max_i \hat{p}_i$, $i = 1 \div 200$ – valoarea maxima a estimatiilor;
- $\Delta_M = \max_i p_{2,i} - \min_i p_{1,i}$, $i = 1 \div 200$ – reuniunea celor 200 de intervale de încredere statistica;

- $\Delta_M^c = \max_i \hat{p}_i - \min_i \hat{p}_i$, $i = 1 \div 200$ – diferența maximă între două estimatii (intervalul de împrastiere al estimatiilor);
- $\delta_M = \max_i |\hat{p}_i - p^*|$, $i = 1 \div 200$ – diferența maximă între estimatiile \hat{p}_i și frecvența relativă p^* ;
- $\delta_m = \min_i |\hat{p}_i - p^*|$, $i = 1 \div 200$ – diferența minimă între estimatiile \hat{p}_i și frecvența relativă p^* .

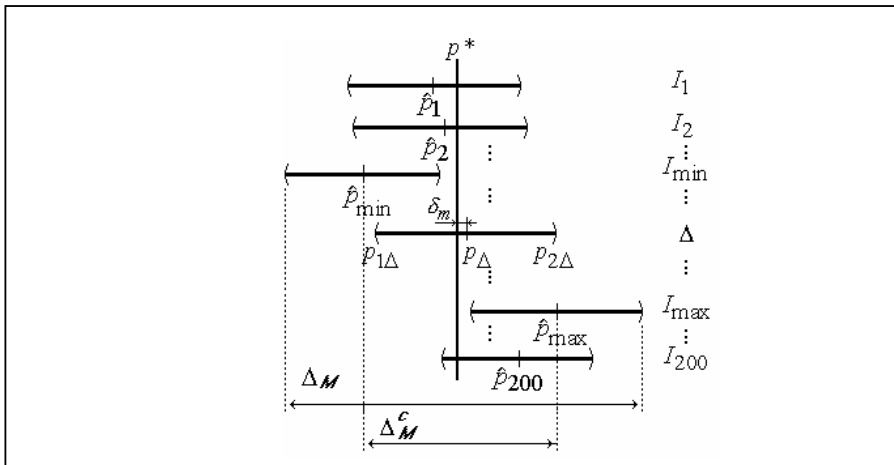


Figura 2. Marimi utilizate în obținerea intervalului de încredere statistică “reprezentativ” pentru probabilitate

Urmatoarele întrebări (probleme) au ghidat analiza noastră teoretică și experimentală:

1. **Cât de largi sunt intervalele Δ_M^c , δ_M și Δ_M ?** Intervalele Δ_M^c și δ_M sunt importante în analiza împrastierii estimatiilor în jurul valorii p^* . Intervalul Δ_M ne da o idee despre cel mai mare interval în care s-ar afla p , probabilitatea adevărată, banuită ca există.
2. **Exista valori \hat{p}_i foarte apropiate de p^* și cât de apropiate?**
3. Pentru a răspunde la această întrebare a fost urmărită experimental mărimea δ_m . δ_m conduce la estimatia \hat{p}_i care este cea mai apropiată de p^* ,

estimatie care va fi în continuare notata cu p_{Δ} . S-a notat cu Δ intervalul de încredere statistica asociat estimatiei p_{Δ} conform relatiei (1).

4. **Cât de multe intervale de încredere statistica I_i contin (îmbraca) p^* ?**
Prin presupunerea de stationaritate, ne asteptam ca un mare numar de intervale de încredere statistica I_i sa se intersecteze, continându-l în acelasi timp pe p^* . Nu ne asteptam la o proportie de 95%, întrucât cele 200 de multimi de date *i.i.d.* nu sunt independente între ele.
5. Putem gasi un interval de încredere statistica pentru probabilitatea adevarata p , interval care sa fie acceptat de toate cele 200 de multimi de date experimentale? Se pot gasi mai multe astfel de intervale? Daca ipoteza de stationaritate este adevarata atunci astfel de intervale trebuie sa existe. Daca intervalul Δ (definit mai sus în întrebarea 2) este unul dintre aceste intervale, atunci el va fi preferat de experimentator si va fi considerat ca "reprezentativ" pentru probabilitatea cuvântului si textul analizat.

Pentru a înțelege metoda dezvoltata de noi care raspunde la aceste întrebări si care conduce la obtinerea intervalului de încredere statistica "reprezentativ" exemplificam pentru cuvântul DE în corpusul mixt global #CMG.

În Tabelul 3 cele 200 de rânduri corespund celor 200 de multimi de date experimentale *i.i.d.* fiind explicitate atât estimatiile \hat{p}_i cât si intervalele de încredere statistica I_i , $i = 1 \div 200$. Succesiv, fiecare interval I_i a fost considerat ca interval de referinta si am aplicat 199 de teste ale ipotezei ca probabilitatea apartine intervalului mentionat, test descris în Anexa 1. Fiecare dintre cele 199 de teste este aplicat pe o singura multime de date experimentală. (Numarul 199 se explica prin faptul ca nu se testeaza si multimea care a produs intervalul de referinta.)

În primul rând al tabelului 3, intervalul $I_1 = (3,89; 4,27) \times 10^{-2}$ este intervalul de referinta fata de care se aplica testele de apartenenta a probabilitatii. Se testeaza daca probabilitatea cuvântului DE apartine sau nu intervalului I_1 pe baza unei singure multimi de date *i.i.d.*; aceasta înseamna ca verificam succesiv fiecare din restul multimilor de date, anume $i = 2 \div 200$. Acceptarea ipotezei ca probabilitatea cuvântului DE apartine intervalului I_1 este marcata cu "DA" pentru respectiva multime de date (Tabelul 3, rândul 1). În caz contrar, pe pozitia respectiva este completat "NU". Numarul total de multimi de date care trec testele este continut în ultima coloana din dreapta. Aceasta procedura se repeta alegând ca referinta pe rând toate cele 200 de intervale de încredere statistica I_i , $i = 1 \div 200$.

Tabel 3

Determinarea intervalului de încredere statistica "reprezentativ" Δ pentru probabilitatea cuvântului DE în corpusul mixt global #CMG. Este îngrosat rândul 3 care corespunde intervalului "reprezentativ" Δ

	\hat{p}_i	$I_i (\times 10^{-2})$	Multime i	Total

i	$(\times 10^{-2})$	$p_{1,i}$	$p_{2,i}$	1	2	3	...	99	...	97	...	200	"DA"
1	4,08	3,89	4,27		"DA"	"DA"		"DA"		"DA"		"DA"	199
2	4,06	3,87	4,25	"DA"		"DA"		"DA"		"DA"		"DA"	199
3	4,10	3,91	4,29	"DA"	"DA"			"DA"		"DA"		"DA"	199
...													
94	4,31	4,12	4,51	"DA"	"DA"	"DA"				"NU"		"DA"	182
...													
99	3,78	3,60	3,96	"DA"	"DA"	"DA"		"NU"				"DA"	121
...													
200	4,18	3,99	4,37	"DA"	"DA"	"DA"		"NU"		"DA"			198

Pentru cuvântul DE s-a obținut $p^* = 0,040986$, iar estimatia cea mai apropiată de p^* a fost $\hat{p}_3 = p_\Delta = 0,040992$, $\delta_m / p^* = 0,0002$. Pentru estimatia p_Δ se obține intervalul de încredere statistică 95% $\Delta = (0,0391; 0,0429)$. Din rândul 3 se observă că intervalul Δ trece toate cele 199 de teste ale ipotezei că probabilitatea cuvântului DE este cuprinsă în interiorul său. Sunt multe intervale I_i care au compatibilitate cu toate multimile de date *i.i.d.* (în ultima coloană numărul 199 a apărut de 101 ori). Dintre aceste 101 intervale am ales $\Delta = (0,0391; 0,0429)$ ca fiind interval de încredere statistică 95% *reprezentativ* pentru probabilitatea cuvântului DE întrucât este ușor de determinat de oricare experimentator. Multimea de date *i.i.d.* specificată de indicele $i = 3$ va fi numită multime de date "*reprezentativă*" pentru cuvântul DE în corpusul #CMG.

Tabelul 4 conține informații despre elementele analizei pentru primele zece cuvinte din ierarhia frecvențelor relative în corpusul #CMG. Exemplificăm pentru cuvântul DE care având frecvența relativă $p^* = 4,10 \times 10^{-2}$ este pe primul loc în ierarhie. Valoarea p^* este cuprinsă în $N(p^*) = 192$ de intervale de încredere statistică din cele 200 considerate (coloana 3); reuniunea celor 200 intervale de încredere statistică raportată la p^* este $\Delta_M / p^* = 22,24 \times 10^{-2}$, (coloana 4); diferența maximă între două estimatii raportată la p^* este $\Delta_M^c / p^* = 13,08 \times 10^{-2}$, (coloana 5); diferența maximă între o estimatie \hat{p}_i și p^* raportată la p^* este $\delta_M / p^* = 7,80 \times 10^{-2}$, (coloana 6); diferența minimă între o estimatie \hat{p}_i și p^* raportată la p^* este $\delta_m / p^* = 0,02 \times 10^{-2}$, (coloana 7); lățimea intervalului de încredere statistică "*reprezentativ*" Δ raportată la p^* este

$\Delta/p^* = 9,23 \times 10^{-2}$, (coloana 8); exista $N(\Delta) = 101$ intervale de încredere statistica la fel de bune ca intervalul Δ "reprezentativ", (coloana 9); Aceste $N(\Delta)$ intervale sunt confirmate de toate cele 199 de teste de apartenență a probabilității la interval, prin care s-a făcut verificarea staționarității.

Numarul relativ mare de intervale de încredere statistica confirmate practic de toate cele 199 de teste de apartenență a probabilității la interval – $N(\Delta)$ din coloana 9 a tabelului 4 – este o susținere a ideii de staționaritate.

Tabel 4

Rezultate numerice privind marimile din Fig. 2 pentru cele mai frecvente 10 cuvinte în #CMG. Valorile din coloanele 2, 4-8 sunt multiplicare cu 100

Cuvânt	p^*	$N(p^*)$	Δ_M/p^*	Δ_M^c/p^*	δ_M/p^*	δ_m/p^*	Δ/p^*	$N(\Delta)$
<i>I</i>	2	3	4	5	6	7	8	9
DE	4,10	192	22,24	13,08	7,80	0,02	9,23	101
ȘI	3,20	198	25,42	14,91	8,77	0,04	10,45	102
ÎN	2,67	194	24,88	13,43	7,09	0,04	11,43	172
SĂ	1,62	189	35,74	21,02	11,40	0,05	14,69	122
A	1,47	191	37,11	21,66	11,40	0,05	15,42	124
LA	1,46	185	38,07	22,52	12,89	0,00	15,45	104
SE	1,39	190	39,97	24,27	13,92	0,02	15,89	79
CU	1,38	195	37,40	21,55	11,30	0,05	15,92	132
O	1,30	191	37,52	21,29	12,38	0,02	16,39	120
NU	1,28	189	40,20	23,62	12,78	0,01	16,54	120

Prin centralizarea acestor tipuri de rezultate pentru toate corpusurile analizate și pentru toate cuvintele pentru care s-a putut face analiza a rezultat Tabelul 5. Concret, studiul experimental a cuprins toate corpusurile prezentate în Introducere. Am putut aplica inferențele statistice doar pentru acele cuvinte pentru care am avut suficiente date; anume $N p^*(1-p^*) > 20$, unde N este volumul multimii de date *i.i.d.* (forma experimentală pentru condiția DeMoivre - Laplace). Cuvintele au fost sortate în ordine descrescătoare a frecvențelor de apariție p^* . Aceasta sortare a permis organizarea studiului pe clase de frecvență. Am ales ca limite ale claselor următoarele valori: 5%, 2%, 1%, 0,5%, 0,2%, 0,1% și 0,05%.

În studiul nostru experimental în aproape toate situațiile (oricare cuvânt urmărit și orice corpus lingvistic investigat) am găsit o estimatie p_Δ practic egală cu p^* . Acest lucru se vede în Tabelul 5, coloana 8 urmărind raportul dintre δ_m și p^* . Pentru toate situațiile analizate am obținut $\delta_m/p^* \leq 2,23\%$. Având în vedere că studiul experimental a condus și la obținerea de intervale Δ "reprezentative" în toate situațiile analizate, rezulta că aceste intervale de încredere statistica 95% pot fi scrise sub forma:

$$\Delta = (p_{1\Delta}; p_{2\Delta}) \cong p^*(1 \mp \varepsilon_r), \quad \varepsilon_r \cong 1.96 \times \sqrt{(1-p^*)/(N p^*)} \quad (2)$$

ε_r este eroarea relativa cu care se determina probabilitatile.

Exemplificam citirea Tabelului 5 pentru corpusul #CMG si clasa a doua de frecventa. Exista 8 cuvinte (coloana 3) care au frecventele relative cuprinse între (0,01; 0,02). Aceste 8 cuvinte acopera 11,17% (coloana 4) din totalul aparitiilor de cuvinte din #CMG, $L_C = 8806433$. Celelalte coloane, 5-9, contin informatii referitoare la marimile din Fig. 2. Astfel coloana 9 contine raportul dintre lungimea intervalului Δ si p^* pentru cuvintele existente în clasa respectiva (limita minima si maxima). Acest raport este practic dublul erorii relative, ε_r , în determinarea probabilitatii cuvântului; se observa o precizie relativ buna a determinarilor din aceasta clasa, $\varepsilon_r \leq 8.5 \times 10^{-2} = 17 \times 10^{-2} / 2$.

În total în #CMG au fost $194 = 3 + 8 + 8 + 24 + 59 + 92$ cuvinte pentru care s-a putut determina intervalul Δ "reprezentativ". Desi cele 194 cuvinte reprezinta o mica pondere din totalul cuvintelor distincte posibile, ele acopera 48,87% din $L_C = 8806433$, totalul aparitiilor de cuvinte în corpusul mixt global, #CMG.

Tabel 5

Rezultate experimentale organizate pe clase de frecvente relative. Valorile din coloanele 4-9 au fost înmultite cu 100

Clasa de frecvente	Corpus	Nr.	Aco- perire	Δ_M / p^*	Δ_M^c / p^*	δ_M / p^*	δ_m / p^*	Δ / p^*
1	2	3	4	5	6	7	8	9
$2 \times 10^{-2} \leq p^* < 5 \times 10^{-2}$	#CMG	3	9,97	22-25	13-15	7-9	0,02-0,04	9-11
	#1JCMG	3	10,04	31-39	18-22	10-13	0,01-0,07	13-16
	#2JCMG	3	9,90	30-41	17-25	9-14	0,01-0,06	13-16
	#CLG	3	9,58	27-33	16-19	8-11	0,03-0,04	11-14
	#1JCLG	3	9,60	39-47	23-27	12-14	0,05-0,10	16-20
	#2JCLG	3	9,55	38-46	22-25	12-13	0,03-0,10	16-20
	#CSG	4	13,76	52-83	27-52	14-27	0,05-0,37	25-35
$10^{-2} \leq p^* < 2 \times 10^{-2}$	#CMG	8	11,17	36-40	21-24	11-14	0,00-0,05	15-17
	#1JCMG	9	12,07	49-59	28-33	14-18	0,00-0,20	21-27
	#2JCMG	8	11,28	46-60	24-37	12-19	0,10-0,16	21-24
	#CLG	9	12,73	39-52	23-30	12-17	0,00-0,09	16-22
	#1JCLG	9	12,69	55-78	30-48	16-25	0,02-0,26	23-30
	#2JCLG	10	13,77	51-75	28-44	14-24	0,00-0,25	23-31
	#CSG	4	5,36	110-130	67-78	34-48	0,24-0,59	44-50
$5 \times 10^{-3} \leq p^* < 10^{-2}$	#CMG	8	6,48	42-61	23-36	12-20	0,00-0,16	19-26
	#1JCMG	8	5,93	60-99	33-62	17-34	0,11-0,29	27-38
	#2JCMG	8	6,51	58-84	31-52	15-29	0,01-0,26	27-36
	#CLG	9	6,62	53-77	30-45	15-27	0,03-0,27	22-31
	#1JCLG	8	6,15	73-112	40-68	23-42	0,01-0,54	31-44
	#2JCLG	9	6,14	72-107	40-66	24-39	0,06-0,42	32-44
#CSG	12	8,46	135-186	71-117	36-64	0,35-1,45	58-78	

$2 \times 10^{-3} \leq p^* < 5 \times 10^{-3}$	#CMG	24	7,36	59-103	31-63	16-33	0,01-0,43	27-41
	#1JCMG	24	7,17	90-158	51-100	27-59	0,01-0,82	39-60
	#2JCMG	23	7,05	91-143	49-87	25-45	0,04-0,84	38-60
	#CLG	25	7,02	74-123	42-75	21-40	0,05-0,73	32-50
	#1JCLG	23	6,87	116-186	65-116	37-66	0,02-1,33	46-71
	#2JCLG	23	6,40	105-186	57-115	29-66	0,02-1,20	48-71
	#CŞG	2	0,87	198-198	111-112	69-72	1,08-2,14	82-83
$10^{-3} \leq p^* < 2 \times 10^{-3}$	#CMG	59	7,78	92-156	46-97	26-61	0,02-1,09	43-61
	#1JCMG	57	7,40	135-224	74-141	41-80	0,02-2,07	61-87
	#2JCMG	61	8,13	135-217	75-134	40-77	0,04-2,23	61-87
	#CLG	58	7,74	114-190	62-120	33-67	0,01-1,41	51-73
	#1JCLG	33	5,21	158-231	81-143	42-85	0,07-1,84	73-91
	#2JCLG	32	5,06	160-220	85-129	46-82	0,04-2,11	72-91
	#CŞG	0	0,00	-	-	-	-	-
$5 \times 10^{-4} \leq p^* < 10^{-3}$	#CMG	92	6,11	134-224	72-142	37-82	0,02-2,16	62-90
	#1JCMG	8	0,72	185-221	101-134	52-76	0,30-1,89	88-90
	#2JCMG	9	0,81	194-242	111-150	57-95	0,25-1,89	88-91
	#CLG	52	4,11	162-224	88-131	45-82	0,05-1,95	73-90
	#1JCLG	0	0,00	-	-	-	-	-
	#2JCLG	0	0,00	-	-	-	-	-
	#CŞG	0	0,00	-	-	-	-	-
Total	#CMG	194	48,87					
	#1JCMG	109	43,33					
	#2JCMG	112	43,68					
	#CLG	156	47,80					
	#1JCLG	76	40,52					
	#2JCLG	77	40,92					
	#CŞG	22	28,45					

Tabelul 5 indica și precizia determinărilor (eroarea relativă ε_r) pentru cuvintele analizate. Aceasta precizie este relativ bună pentru determinările făcute pe corpusul mixt global #CMG (pentru cuvinte din primele patru clase de frecvență, $\varepsilon_r \leq 20,5 \times 10^{-2} = 41 \times 10^{-2} / 2$).

Aplicând procedura descrisă în Cap. 2.1 pentru toate corpusurile lingvistice și pentru toate cuvintele care au satisfăcut condiția de Moivre - Laplace au rezultat probe în sprijinul ipotezei de staționaritate a limbii române scrise.

Intervalele "reprezentative" precum și multimile de date *i.i.d.* "reprezentative" determinate pentru un cuvânt anumit și textul natural considerat au fost în continuare folosite în Cap. 2.2, pentru a analiza dacă putem vorbi despre un model matematic al sursei de cuvinte pentru limba ca ansamblu, pentru diverse domenii ale limbii, pentru diversi autori, etc.

Acuratețea în determinarea probabilității cuvintelor este dată de:

- încrederea statistică (95%);
- erorile relative, ε_r , cu care s-au obținut intervalele Δ "reprezentative" conform Tabelului 5;
- mărimea celor două tipuri de erori statistice care apar în testul de apartenență a probabilității la interval, întrucât acest tip de test se bazează pe validarea intervalului Δ ca "reprezentativ".

În ceea ce privește testul de apartenență a probabilității la interval acesta a fost aplicat pentru un prag statistic $\alpha = 0,05$. Întrucât testul a fost trecut de fiecare dată (fapt pentru care Δ a fost validat ca "reprezentativ" este important de a avea un control asupra mărimii β , probabilitatea de a accepta date false. Dacă am impune valori mici pentru β am avea nevoie de un corpus mai mare. Spre exemplu, conform [6, Tabel 4], dacă se dorește $\beta \leq 0,3$ și $\delta = 0,15$ (Anexa 1) pentru a investiga cuvinte din primele patru clase de frecvență am avea nevoie de un corpus de circa 30 de milioane de cuvinte.

2.2 Comparații matematice între diverse texte naturale pe baza structurii de cuvinte

Investigația noastră (privind staționaritatea) a fost completată cu comparații matematice privind probabilitățile cuvintelor, pe care le-am organizat după următoarele criterii:

- a) Se verifică dacă un același cuvânt are aceeași probabilitate în cele două texte naturale care se compară. Această comparație va fi numită în continuare *comparație între cuvinte ca atare*.
- b) Se verifică dacă probabilitățile cuvintelor situate pe un același rang în ierarhia frecvențelor relative din cele două texte sunt egale. Spre exemplu, pe rangul 20 în corpusul literar global, se află cuvântul DAR, iar în corpusul științific cuvântul UN, vezi Tabelul 1. La comparația între cele două domenii se va urmări dacă probabilitatea celor două cuvinte (DAR și UN) este aceeași. În cele ce urmează numim acest criteriu *comparație pe baza rangului*.

Toate comparațiile matematice, atât pe baza criteriului a) cât și pe baza criteriului b) au fost făcute folosind următoarele teste statistice:

- $T1$ – test al ipotezei ca probabilitatea aparține unui interval, (Anexa 1);
- $T2$ – test de egalitate între două probabilități, (Anexa 2).

Pentru fiecare din cele două texte naturale care se compară și pentru fiecare cuvânt investigat s-au determinat în prealabil intervalele "reprezentative" precum și multimile de date *i.i.d.* "reprezentative".

Când aplicăm testul $T1$, intervalul $(a;b)$ este intervalul "reprezentativ" Δ din primul text natural implicat în comparație, iar mulțimea $[x_1, x_2, \dots, x_N]$ de date experimentale *i.i.d.* este mulțimea de date "reprezentativă" din cel de-al doilea text natural.

Testul a fost aplicat în ambele situații: corpus1 *versus* corpus2 și corpus2 *versus* corpus1, Tabel 6.

Când aplicăm testul $T2$ se considera pentru comparație cele două mulțimi de date *i.i.d.* "reprezentative" extrase din cele două texte naturale pentru cuvintele care se compară.

Toate testele au fost aplicate pentru un prag de semnificație statistică $\alpha = 0,05$. Cu alte cuvinte probabilitatea de a respinge date corecte este mai mică decât 0,05.

Tabel 6

**Comparații între texte naturale pe baza probabilității cuvintelor.
Coloanele 4-9 conțin numărul de cuvinte rejectate de testele statistice**

Texte comparate		Nr.	Comparație între cuvinte ca atare			Comparație pe baza rangului		
Corpus 1	Corpus 2		Test $T1$		Test $T2$	Test $T1$		Test $T2$
1	2		1 versus 2	2 versus 1		1 versus 2	2 versus 1	
<i>l</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>
#1JCLG	#2JCLG	72	0	0	0	0	0	0
#1JCMG	#2JCMG	104	0	0	0	0	0	0
#CLG	#CSG	22	10	18	13	1	16	10

Rezultatele experimentale sunt sintetizate în Tabelul 6. Comparațiile făcute în cadrul domeniului literar, când se compară cele două jumătăți de corpus între ele (#1JCLG și #2JCLG) nu indică diferențe între probabilități indiferent de testul utilizat ($T1$ sau $T2$) sau de criteriul utilizat (comparații pe baza aceluiași cuvânt sau pe baza aceluiași rang).

Același rezultat s-a obținut și când s-au comparat cele două jumătăți ale corpusului mixt global, #1JCMG și #2JCMG.

Exemplificăm în continuare modul de citire al Tabelului 6.

Primele două coloane conțin corpusurile care se compară între ele.

Coloana 3 indică numărul de cuvinte investigate în comparații (care au îndeplinit condiția $Np^*(1-p^*) \geq 20$ în ambele texte care se compară).

Rezultatele din coloanele 4, 5 și 6 au fost obținute aplicând criteriul comparațiilor "cuvintelor ca atare".

Coloanele 4 și 5 arată câte cuvinte nu au trecut testul $T1$ de apartenență a probabilității la interval. Coloana 4 se referă la situația când intervalul fix $(a;b)$ este intervalul Δ "reprezentativ" din primul corpus al comparației, iar mulțimea de date *i.i.d.* supusă testului este mulțimea *i.i.d.* "reprezentativă" din al doilea corpus. Similar, în coloana 5: intervalul fix $(a;b)$ este intervalul Δ "reprezentativ" din al doilea corpus al comparației, iar mulțimea de date *i.i.d.* supusă testului este mulțimea *i.i.d.* "reprezentativă" din primul corpus.

Coloana 6 conține numărul de cuvinte care sunt rejectate de testul $T2$ de egalitate între probabilități.

Coloanele 7, 8 și 9 conțin același tip de informație specificat în coloanele 4, 5 și 6, cu diferența că de această dată se compară cuvintele care ocupă același rang în loc de cuvintele “ca atare”.

Când se compară domenii diferite, spre exemplu literar și științific, apar multe diferențe marcate de ambele teste $T1$ și $T2$ și de cele două criterii de comparație.

Rezultatele comparațiilor puntează unele diferențe între domeniile literar și științific. Testele nu au indicat diferențe când s-au comparat corpusuri organizate după aceeași regulă (jumătățile corpusului mixt global între ele sau jumătățile corpusului literar global între ele); reamintim că atât corpusul mixt global cât și cel literar global au fost obținute prin concatenarea aleatoare a cartilor respective.

3. Legea lui Zipf. Studiu experimental

Ierarhiile frecvențelor relative ale cuvintelor (prezentate în Cap. 1 și întărite de analiza de staționaritate din Cap. 2) au constituit o bază de plecare pentru studiul nostru experimental asupra legii lui Zipf. În lingvistică legea lui Zipf este una din cele mai cunoscute dependente rang – frecvență. (Aceste dependente rang – frecvență au fost observate de a lungul timpului și în diverse alte domenii: economie, fizică, biologie, demografie, etc. [11], [12].) Obiectivul acestui capitol a fost de a stabili dacă și în ce măsură (cu ce acuratețe) limba română scrisă satisface legea lui Zipf.

Fie un text (corpus) având o lungime de L_c cuvinte, dintre care N_c sunt distincte. Aceste N_c cuvinte se sortează într-o listă în ordine descrescătoare a numărului de apariții în textul natural. Se notează cu k rangul unui cuvânt în listă și cu $f(k)$ frecvența relativă a acestuia (numărul de apariții raportat la L_c): $f(1) \geq f(2) \geq \dots \geq f(N_c)$. (Altfel spus, $f(k)$ este de tipul p^* din capitolele precedente). Legea lui Zipf afirmă că produsul dintre rang și frecvența relativă este constant, [11] – [14].

$$k f(k) = A \quad (3)$$

Se observă că membrul stâng al ecuației (3) corespunde realității fiind vorba de măsurători efectuate pe texte naturale în timp ce membrul drept corespunde modelului teoretic presupus.

Este știut din considerații privind alte limbi naturale că legea Zipf, apreciată ca foarte simplă și foarte atractivă, funcționează cu aproximație pentru o plajă limitată de ranguri, anume nu prea mici și nu prea mari. Astfel un prim pas al studiului nostru teoretic și experimental a fost să reprezentăm grafic dependența rang – frecvență pe tot corpusul de care am dispus (corpusul mixt global, #CMG). Fig. 3 prezintă această dependență la scara logaritmică ($f(k)$ versus k). La o primă vedere am putea spune că mărimea A din (3) este aproximativ constantă pentru un interval de ranguri $k \in [k_{\min}; k_{\max}]$ unde $k_{\min} > 50$. Am limitat studiul la acele ranguri pentru care numărul de apariții ale cuvintelor a fost mai mare decât 50 pentru a beneficia de rezultatele anterioare privind studiul de staționaritate prezentat în Cap. 2. Aceasta face ca rangul k_{\max} să depindă de corpusul analizat.

Legea lui Zipf este descrisă în numeroase referințe dintre care în limba română menționăm în special [13] și [14]. Capitolul de față urmărește determinarea constantei legii atât pe corpusul de ansamblu, #CMG, cât și pe diverse texte naturale (grupate după autori sau pe subdomenii ale limbii). Se analizează și în ce măsură comportamentul real se abate de la cel teoretic.

3.1 Elemente teoretice

3.1.1 Determinarea parametrului legii Zipf prin minimizarea erorii practice

Presupunând valabilitatea legii Zipf pentru rangurile $k \in [k_{\min}; k_{\max}]$ ne-am propus să determinăm mărimea A din condiția de minimizare pe acest interval a următoarei funcții (suma patratelor erorilor):

$$g(A) = \sum_{k=k_{\min}}^{k_{\max}} \left[f(k) - \frac{A}{k} \right]^2 \quad (4)$$

Derivând funcția $g(A)$ și egalând cu 0 se obține valoarea mării A corespunzând minimului:

$$A = \frac{\sum_{k=k_{\min}}^{k_{\max}} \frac{f(k)}{k}}{\sum_{k=k_{\min}}^{k_{\max}} \frac{1}{k^2}} \quad (5)$$

Valorile k_{\min} și k_{\max} sunt la dispoziția experimentatorului. Pentru o evaluare a acurateții cu care limba naturală verifică legea lui Zipf definim următoarele tipuri de erori:

- ε , suma patratelor erorilor pe intervalul $k \in [k_{\min}; k_{\max}]$ și forma ei normată, ε_n :

$$\varepsilon = \sum_{k=k_{\min}}^{k_{\max}} \left[f(k) - \frac{A}{k} \right]^2 \quad \varepsilon_n = \varepsilon / \left[\sum_{k=k_{\min}}^{k_{\max}} \left(\frac{A}{k} \right)^2 \right] \quad (6)$$

- ❖ ε_{\max} , eroarea relativă maximă pe intervalul de optimizare $k \in [k_{\min}; k_{\max}]$:

$$\varepsilon_M = \varepsilon_r(k_M) = \max_k \varepsilon_r(k) \quad \varepsilon_r(k) = \left| f(k) - \frac{A}{k} \right| / \left(\frac{A}{k} \right) \quad (7)$$

3.1.2 Determinarea parametrului legii lui Zipf considerând cazul ideal

Dacă acceptăm legea lui Zipf ca fiind corectă pe întreg domeniul de ranguri $k \in [1; N_c]$, atunci valoarea constantei A se determină prin raționamentul descris în [13], [14]:

$$A = \frac{1}{c + \ln N_c} \quad (8)$$

unde c este constanta lui Euler, egala cu 0,577215 si $N_c > 50$.

Observam ca marimea A calculata cu relatia (8) nu depinde decât de numarul N_c de cuvinte distincte din textul analizat. Prin urmare sunt de asteptat unele diferente între evaluarile marimii A pe baza datelor experimentale cu relatia (5) si cazul ideal, pur teoretic, relatia (8).

3.1.3 Corolar al legii lui Zipf

Rezultatele experimentale cuprind si verificarea unui corolar al legii Zipf care se refera la determinarea cotei parti, l_s/L_c , pe care o acopera cele mai frecvente s cuvinte într-un text de lungime L_c , [13], [14].

$$\frac{l_s}{L_c} = \frac{c + \ln s}{c + \ln N_c} \quad (9)$$

Relatia (9) este valabila pentru un numar de cuvinte $s > 50$.

Observam ca valoarea raportului l_s/L_c nu depinde de marimea A . De aceea diferentele existente între diversele moduri de evaluare ale marimii A nu vor influenta acest raport. În consecinta ne asteptam la o buna verificare experimentală a acestui corolar.

3.2. Rezultate experimentale si concluzii

Analiza experimentală a legii lui Zipf a început cu corpusul global #CMG (vezi Fig. 3) si a continuat pentru comparatie cu o serie de texte naturale incluse în acesta (prezentate în Introducere). Rezultatele experimentale sunt concentrate în Tabelul 7. Pentru fiecare text analizat Tabelul 7 contine în coloanele 2 si 3 numarul total de cuvinte L_c si numarul cuvintelor distincte N_c . În toate textele analizate s-au investigat toate cuvintele cu numar de aparitii mai mare decât 50; acesta determina rangul k_{\max} corespunzator fiecarui text analizat (coloana 4). k_{\max} difera de la text la text; k_{\min} este ales întotdeauna 51. Pentru acest interval de ranguri, $k \in [k_{\min}; k_{\max}]$, s-a determinat cu relatia (5) marimea A cuprinsa în coloana 5. Coloanele 6 – 9 contin rezultatele numerice calculate cu relatiile (6) si (7) unde marimea A este cea din coloana 5 (determinata din textul natural respectiv). Coloana 9 contine rangul k_M pentru care s-a obtinut eroarea relativa maxima ε_M .

Ne-am pus problema si daca marimea $A = 0,0909$ determinata pentru corpusul mixt global #CMG, ar putea fi acceptata drept referinta pentru limba româna. De aceea coloanele 10 – 13 contin succesiv marimile din relatiile (6) si (7) unde $A = 0,0909$ pentru toate textele naturale analizate. Eroarea relativa maxima ε_M este însoțita de rangul corespunzator, k_M .

Tabel 7

Studiu experimental al legii lui Zipf în limba română scrisă

Text	L_c	N_c	k_{\max}	A $\times 10^2$	ε 10^6	ε_n $\times 10^2$	ε_M $\times 10^2$	k_M	ε	ε_n	ε_M	k_M
1	2	3	4	5	6	7	8	9	10	11	12	13
#CMG	8806433	202403	14543	9,09	0,36	0,22	9,81	286	0,36	0,22	9,81	286
#CLG	6255235	162124	10299	9,60	0,30	0,17	13,93	10136	0,81	0,50	15,43	149
#1JCLG	3127618	116247	5568	9,58	0,26	0,14	10,53	136	0,72	0,44	16,41	136
#2JCLG	3127617	116860	5529	9,74	0,29	0,16	10,02	122	1,11	0,69	17,86	122
#1	226420	26943	466	9,81	0,37	0,22	9,07	68	1,28	0,88	16,81	173
#2	121177	18457	260	10,15	0,15	0,09	8,21	256	1,95	1,48	20,83	256
#3	88827	13768	190	10,07	0,20	0,14	10,17	186	1,60	1,33	22,07	186
#2+#3	210004	25036	484	9,97	0,52	0,29	18,71	478	1,89	1,29	30,18	478
#4	130743	18223	274	10,71	0,26	0,14	8,38	53	4,47	3,35	25,85	110
#5	75698	10351	187	11,56	0,42	0,22	10,86	183	9,22	7,71	40,92	183
#6	197889	23206	399	10,34	0,15	0,08	7,03	121	2,85	1,99	21,73	121
#4+#5+#6	404330	33555	849	10,53	0,20	0,10	8,90	103	4,03	2,62	26,08	103
#7	644794	49434	1195	10,03	0,35	0,18	10,80	477	2,04	1,30	21,49	77

În Fig. 3 sunt prezentate pentru corpusul mixt global două traiectorii, una experimentală (cu 'o') și cea teoretică (cu '*') conform relației (3) cu parametrul $A = 0,0909$ din coloana 5, Tabelul 7. Se observă o bună concordanță a celor două curbe pentru $k \in [k_{\min}; k_{\max}]$.

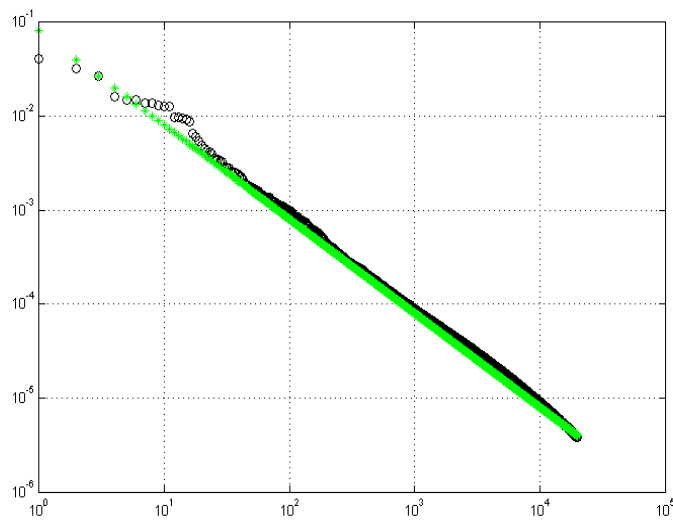


Figura 3. Dependenta rang – frecventa relativa de aparitie a cuvintelor în corpusul mixt global #CMG (scara logaritmica $f(k)$ versus k). Curba experimentală marcată cu ‘o’; curba teoretică, relația (3) pentru $A = 0,0909$, marcată cu ‘*’

În cazul ideal, pur teoretic, mărimea A poate fi determinată cu relația (8) pe baza coloanei 2 din Tabelul 7. Spre exemplu în corpusul global #CMG, unde au fost identificate $N_c = 202403$ cuvinte, $A = 0,0781$. În corpusul #CLG, pentru $N_c = 162124$ cuvinte distincte, aplicând relația (8) rezulta $A = 0,0795$.

Tabel 8

Valori teoretice, relația (9), și experimentale ale raportului l_s/L_c în corpusul literar global #CLG

λ	0,1%	0,05%	0,01%
s	104	189	911
l_s/L_c (experimental)	43,69%	49,64%	64,38%
l_s/L_c (teoretic)	41,53%	46,28%	58,78%

Tabelul 8 conține date despre cota parte acoperită de cuvintele pentru care $f(k) \geq \lambda$ unde $\lambda = 0,1\%; 0,05\%; 0,01\%$, în textul literar global #CLG. S-a folosit relația (9) unde $N_c = 162124$, iar numărul de cuvinte s corespunzător pragului λ este conținut în linia 2 a Tabelului. Se observă o concordanță destul de bună între valorile teoretice și cele experimentale.

Nota: Din cele 189 cuvinte din corpusul literar global #CLG din care au frecvență relativă mai mare decât 0,05%, doar 156 au îndeplinit condiția de Moivre – Laplace și au fost investigate cu control statistic aparând și în Tabelul 5.

Ca o remarcă finală legea lui Zipf poate fi considerată ca valabilă și pentru limba română pentru ranguri nu prea mici și nu prea mari, fapt susținut de Fig. 3 și datele din Tabelul 7.

4. Concluzii. Perspective

Unul din principalele rezultate obținute în cadrul acestei lucrări este de a aduce probe noi privind staționaritatea limbii române scrise, de această dată pe baza structurii de cuvinte. (Ipoteza de staționaritate este inclusă în presupunerea generală conform căreia limbile naturale sunt lanțuri Markov multiple ergodice). Analiza staționarității s-a făcut prin extinderea unei metode dezvoltate de autori pentru studiul structurii statistice de m -grame (litere, digrame, trigrame, tetragrame). În consecință s-au putut obține probabilitățile cuvintelor cu intervale de încredere statistică 95% “representative”. Aceste intervale pe care le-am numit “representative” au avut compatibilitate cu toate multimile de date *i.i.d.*

obținute prin esantionarea periodica a textului natural. Simultan au rezultat multimile de date *i.i.d.* "representative" pentru cuvântul investigat și textul natural analizat.

O alta contribuție constă în procedura de comparație matematică între texte naturale facilitată de intervalul "representativ" pentru probabilitate și de multimile de date *i.i.d.* "representative". Comparațiile făcute între corpusuri organizate în aceeași manieră (literar *versus* literar sau mixt *versus* mixt) au întărit ideea de staționaritate a limbii și au confirmat modelul matematic prezentat anterior prin intervale de încredere statistică 95% "representative" pentru probabilitățile cuvintelor. Au apărut unele diferențe între domeniile literar și științific.

Rezultatele experimentale dau un plus de semnificație frecvenței relative, p^* , marime de care orice experimentator este interesat. Acest plus de semnificație este datorat faptului că în toate situațiile analizate de noi (cuvânt sau text natural) am putut obține o estimatie a probabilității practic egală cu p^* , iar intervalul de încredere statistică asociat acestei estimatii a fost confirmat ca interval "representativ" pentru probabilitate.

Lucrarea conține totodată și confirmarea valabilității pentru limba română scrisă a legii lui Zipf (lege de tip rang - frecvență) și a unui corolar al acesteia de interes lingvistic.

Autorii doresc să multumească D-lui dr. ing. Dan TUFIS, membru corespondent al Academiei Române, pentru sprijinul științific acordat continuu în studiul limbii române scrise.

Autorii menționează, de asemenea, sugestiile utile primite din partea D-lui Prof. dr. ing. Alexandru Serbanescu de la Academia Tehnică Militară.

Referințe bibliografice

- [1] Shannon C. E., "Prediction and Entropy of Printed English", Bell Syst. Tech. J., Vol. 30, pp. 50-64, January 1951.
- [2] Adriana Vlad, and A. Mitrea 1997 "Estimating conditional probabilities and digram statistical structure in printed Romanian". In Tufis D, Andersen P. (Eds), *Recent Advances in Romanian Language Technology*, Bucharest, Ed. Academiei, ISBN 973-27-0626-0, pp. 57-72, <http://www.racai.ro/books/awde/vlad.html>.
- [3] Adriana Vlad, A. Mitrea, M. Mitrea, D. Popa, "Statistical methods for verifying the natural language stationarity based on the first approximation. Case study: Printed Romanian" în Vol. **VEXTAL'99** (Conferința Veneția per il trattamento automatico della lingue), Ed. Unipress, ISBN 88-8098-112-9, pp. 127-132, Nov. 22-24, 1999, Venetia; <http://byron.cgm.unive.it/events/papers/vlad.pdf>.
- [4] Adriana Vlad, A. Mitrea, M. Mitrea, "Verifying Printed Romanian Language Stationarity Based on the Digram Statistical Structure", Proceedings of the Romanian Academy, Series A, Vol. I, No. 2/2000, pp. 129-139.

-
- [5] Vlad Adriana, Mitrea A., Mitrea M., "Two frequency–rank laws for letters in printed Romanian", *Procesamiento del Lenguaje Natural*, Revista N^o 24, Septiembre de 2000, pp. 153-160, ISSN 1135-5948.
- [6]. Adriana Vlad, A. Mitrea, M. Mitrea, "The trigram statistical structure in printed Romanian", în **ROMJIST** (Romanian Journal of Information Science and Technology), Vol. 4, No. 3, 2001, pp. 353-372.
- [7]. Adriana Vlad, A. Mitrea, M. Mitrea, "A Corpus – based Analysis of how Accurately Printed Romanian Obeys Some Universal Laws", Capitolul 13 în *Cartea A Rainbow of Corpora: Corpus Linguistics and the Languages of the World*; A. Wilson, P. Rayson, and T. McEnery Editors, Lincom–Europa Publishing House, Munich, 2001, pp. 155-167; lucrarea a fost prezentata la *Corpus Linguistics 2001*, Aprilie 2001, Lancaster, Marea Britanie iar rezumatul este cuprins în *Proc. of CL2001*, pp. 600.
- [8]. Adriana Vlad, A. Mitrea si M. Mitrea, "Contributii privind structura statistica de tetragrame în limba româna scrisa", in *Proc A XXIX-a Sesiune de Comunicari Stiintifice cu Participare Internationala "Tehnologii Moderne în Secolul XXI"*, Academia Tehnica Militara, nov. 2001, Bucuresti, Sectiunea 9.1, pp. 60-65, ISBN 973-8290-27-9.
- [9]. J. Devore, 1987, *Probability and Statistics for Engineering and the Sciences*, 2nd ed., Brooks/Cole Publishing Company, Monterey, Ca.
- [10]. Adriana Vlad, B. Badea si M. Mitrea, 1999, *Metode Statistice în Prelucrarea Informatiei. Compendiu si Aplicatii*, Ed. Metropol, Bucuresti, ISBN 973-562-104-5.
- [11] Kanter I., Kessler D. A., "Markov Processes: Linguistics and Zipf's Law", *Physical Review Letters*, Vol. 74, No. 22, pp. 4559-4562, May 1995.
- [12] Günther R., Levitin L., Schapiro B., Wagner P., "Zipf's Law and the Effect of Ranking on Probability Distributions", *Intl. J. of Theoretical Physics*, Vol. 35, No. 2, pp. 395-417, 1996.
- [13] S. Marcus, Ed.Nicolau, S. Stati, *Introducere în lingvistica matematica*, Ed. Stiintifica, Bucuresti, 1966.
- [14] Dinu M., *Personalitatea limbii române*, Ed. Cartea Româneasca, Bucuresti, 1996.

Anexa 1. Test de apartenență a probabilității la un interval dat – T1

Fie $I = (a; b)$ un interval în care presupunem ca se afla probabilitatea p a unui eveniment urmarit. Dispunem de o multime de date experimentale $[x_1, x_2, \dots, x_N]$, date care satisfac modelul statistic *i.i.d.*. Ne interesează dacă datele experimentale $[x_1, x_2, \dots, x_N]$ confirmă ipoteza că probabilitatea p aparține intervalului $I = (a; b)$, pentru un prag de semnificație statistică, α , ales.

Procedura de test este următoarea:

Se formulează cele două ipoteze statistice, ipoteza nulă H_0 și respectiv ipoteza alternativă H_1 :

H_0 : p aparține intervalului $(a; b)$; $p \in (a; b)$;

H_1 : p este în afara intervalului $(a; b)$; $p \notin (a; b)$.

Se alege pragul de semnificație α (echivalent, nivelul de încredere statistică $1 - \alpha$). Se calculează estimatia $\hat{p} = m/N$, unde cu m s-a notat numărul de succese ale evenimentului în multimea de date $[x_1, x_2, \dots, x_N]$. Verificăm dacă estimatia \hat{p} se afla sau nu în zona de acceptare a datelor. Regiunea de acceptare a datelor este un interval $(c_1; c_2)$ care include $(a; b)$. Intervalul $(c_1; c_2)$ se determină conform relației (10), [3]-[8]:

$$\int_{c_1}^{c_2} \frac{1}{\sqrt{2\pi a(1-a)/N}} \exp\left(-\frac{(x-a)^2}{2a(1-a)/N}\right) dx = \int_{c_1}^{c_2} \frac{1}{\sqrt{2\pi b(1-b)/N}} \exp\left(-\frac{(x-b)^2}{2b(1-b)/N}\right) dx = 1 - \alpha \quad (10)$$

În relația (10) apar două funcții de densitate de probabilitate corespunzătoare legii normale: de medie a și dispersie $a(1-a)/N$ și respectiv de medie b și dispersie $b(1-b)/N$.

Ipoteza nulă H_0 va fi acceptată dacă și numai dacă estimatia \hat{p} aparține intervalului $(c_1; c_2)$. În caz contrar, $\hat{p} \notin (c_1; c_2)$, datele se resping ca fiind semnificative pentru pragul de semnificație α ales (se acceptă ipoteza H_1).

Ca în orice test statistic, pot să apară două tipuri de erori:

Eroarea de tipul (genul) I: Eroarea de a fi respinse date bune, adică să fie respinsă ipoteza H_0 când ea este corectă. Aceasta situație apare atunci când estimatia \hat{p} nu satisface testul, adică $\hat{p} \notin (c_1; c_2)$, și totuși probabilitatea adevărată p este în intervalul $(a; b)$. Probabilitatea acestui tip de eroare este mai mică decât α .

Eroarea de tipul (genul) II: Eroarea de a fi acceptate date false, adica sa fie acceptata H_0 când ea este, de fapt, falsa. Aceasta situatie apare atunci când estimatia \hat{p} satisface testul, $\hat{p} \in (c_1; c_2)$, si totusi probabilitatea adevarata p a evenimentului nu apartine intervalului $(a; b)$, $p \notin (a; b)$. Pentru α si N fixate, probabilitatea acestui tip de eroare depinde de valoarea adevarata necunoscuta p , si se calculeaza cu relatia:

$$\beta(p) = \int_{c_1}^{c_2} \frac{1}{\sqrt{2\pi p(1-p)/N}} \exp\left(-\frac{(x-p)^2}{2p(1-p)/N}\right) dx, \quad p \notin (a; b).$$

$\beta(p)$ este mare atunci când p este la stânga lui a (sau la dreapta lui b), dar foarte aproape de a (respectiv de b). Practic, deranjante sunt situatiile în care $p \leq (1-\delta) \cdot a$ sau $p \geq (1+\delta) \cdot b$, si totusi testul este trecut, adica $\hat{p} \in (c_1; c_2)$. Valoarea δ este determinata (prestabilita) de utilizator, în functie de cât de mult deranjeaza aceasta situatie.

În studiul nostru asupra stationaritatii limbii române acest test a fost absolut necesar, vezi Cap. 2. A trebuit sa stabilim daca probabilitatea p a unui anumit cuvânt este aceeași când dispunem de diverse multimii de date experimentale extrase dintr-un acelasi text natural (unde multimile sunt compatibile cu modelul statistic *i.i.d.*, dar nu sunt independente între ele). Testul a fost folosit si în comparatii între texte naturale.

Anexa 2. Test de egalitate între două probabilități – T2

Disponem de două mulțimi de date experimentale în model statistic *i.i.d.*, de volume N_1 , respectiv N_2 . Notând cu m_1 numărul de succese (aparitii) ale unui eveniment în prima mulțime de date experimentale, estimatia probabilității este $\hat{p}_1 = (m_1 / N_1)$. Similar, pentru a doua mulțime de date experimentale, estimatia probabilității este $\hat{p}_2 = (m_2 / N_2)$. Urmărim să stabilim dacă cele două estimatii \hat{p}_1 și \hat{p}_2 provin din aceeași probabilitate teoretică, respectiv $p_1 = p_2$.

Procedura de test este următoarea:

Se formulează cele două ipoteze statistice, ipoteza nulă H_0 și respectiv ipoteza alternativă H_1 :

H_0 : cele două probabilități teoretice sunt egale $p_1 = p_2$;

H_1 : cele două probabilități teoretice sunt diferite $p_1 \neq p_2$.

Se alege pragul de semnificație statistică α .

Se construiește o valoare de test z conform, [9], [10]:

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\frac{p_1(1-p_1)}{N_1} + \frac{p_2(1-p_2)}{N_2}}}.$$

Valoarea z depinde de datele experimentale prin estimatiile \hat{p}_1 și \hat{p}_2 . În condițiile în care ipoteza H_0 este adevărată z provine dintr-o variabilă aleatoare a cărei lege de repartiție este practic legea normală standard.

Întrucât p_1 și p_2 sunt necunoscute, se consideră $p_1 = p_2 \cong \frac{m_1 + m_2}{N_1 + N_2}$.

În aceste condiții valoarea de test z devine:

$$z = \sqrt{\frac{N_1 + N_2}{N_1 N_2}} \frac{m_1 N_2 - m_2 N_1}{\sqrt{(m_1 + m_2)(N_1 + N_2 - m_1 - m_2)}}. \quad (11)$$

Ipoteza nulă H_0 va fi acceptată (se va considera că probabilitățile sunt egale, $p_1 = p_2$) dacă și numai dacă $|z| \leq z_{\alpha/2}$ ($z_{\alpha/2}$ corespunde pragului de semnificație statistică α ales; am folosit $z_{\alpha/2} = 1,96$). În caz contrar se respinge ipoteza de egalitate a celor două probabilități pentru pragul de semnificație statistică α ales.

Această procedură de test a fost folosită când am comparat între ele diverse texte naturale.

Dezambiguizarea automata a cuvintelor din corpusuri paralele folosind echivalentii de traducere

Dan TUFIS
Institutul de Inteligenta Artificiala
Academia Româna

Rezumat

Corpusurile paralele constituie surse de cunostinte extrem de valoroase, traducerea unui text reprezentând o succesiune de decizii lingvistice pe care traducatorul le ia în vederea asigurarii unei transpuneri cât mai naturale si mai fidele a semnificatiei din textul sursa în textul tradus. Explicitatea si extragerea acestor cunostinte prin metode algoritmice, formalizarea si reutilizarea lor ulterioara constituie provocari ale inteligentei artificiale, subiecte de interes fierbinte în cercetarea actuala. Lucrarea prezinta o serie de contributii în aceasta directie, prezentând mai întâi o metoda originala de identificare a echivalentilor lexicali de traducere a cuvintelor dintr-un corpus paralel (extragând deci un dictionar multilingv) si apoi o metoda extrem de promitatoare pentru identificare automata a diferitelor sensuri ale cuvintelor polisemantice.

Motivatii

Evolutia stiintifica si tehnologica este o sursa permanenta de formare a noi termeni sau a noi sensuri specializate pentru cuvintele existente. În domeniul lexicografiei multilinguale, pastrarea în actualitate a dictionarelor bi- si multilingve fara a apela la tehnologiile informatice, cu precadere cele din sfera ingineriei lingvistice, este aproape imposibila. O serie de studii în domeniul traducerii automate au aratat ca principalele probleme în *acceptabilitatea* traducerilor automate si cu atât mai mult al celor implicând pre- sau post-editare umana, nu sunt legate de probleme de natura sintactica (topica, acorduri, structura frazala) ci ele se regasesc cu precadere în sfera lexicala, mai precis al semanticii lexicale. Evaluarea sistemelor existente de prelucrare a limbajului natural si mai ales a celor de traducere automata (cu variantele ce presupun interventia expertului uman) a condus la identificarea unor puncte sensibile, deficitare (pentru o interesanta trecere în revista a problemelor privind evaluarea sistemelor de prelucrare a limbajului natural si a sistemelor de traducere a se vedea <http://www.isi.edu/natural-language/mteval/>). De pilda, traducerea gresita a unui cuvânt sau al unei expresii într-o fraza perfecta din punct de vedere sintactic este perceptuta de imensa majoritate a consumatorilor de traduceri, cu

precadere de natura stiintifica, ca mult mai grava decât un dezacord gramatical sau vreo alta abatere de la norma gramaticii. S-a invocat pe buna dreptate ca dificultatea majora a prelucrării automate a limbajului este rezolvarea ambiguitatilor lexicale, a omonimiilor si a polisemiei ce apar în orice text (scris sau vorbit). Spre deosebire de oameni, care de multe ori nici nu constientizeaza aceste fenomene (ele sunt „obturate” fie de contextul textului, fie de cunostintele de „bun simt” ale fiecărei persoane), procesoarele artificiale de limbaj natural încearca rezolvarea ambiguitatilor printr-un proces inteligent de alegere, dintr-un spatiu al tuturor solutiilor posibile în raport cu o modelare a limbajului, a solutiei care respecta cel mai bine restrictiile modelului. Raportarea la modelul limbajului este esentiala întrucât dificultatea procesului de prelucrare este cu atât mai mare cu cât modelul este mai complex: spatiul de cautare a solutiilor poate creste exponential, iar procesul decizional poate deveni nedeterminist sau de complexitate neoperationala.

Rezolvarea algoritmica eficienta a omografieii a cunoscut spectaculoase progrese în ultimii 10-15 ani, dar identificarea automata a sensului pe care îl are un anumit cuvânt polisemantic într-un context dat este înca o problema nerezolvata satisfactor si, prin urmare, un subiect „fierbinte” de cercetare. Problema identificării sensului cu care este utilizat un cuvânt este vitala în traducerea automata, întrucât se cunoaste faptul ca de foarte multe ori un cuvânt polisemantic dintr-o limba se traduce într-o alta limba prin cuvinte diferite, în functie de sensul considerat. Este interesant de remarcat ca daca un cuvânt polisemantic din limba sursa se traduce printr-un singur cuvânt polisemantic în limba tinta, sau altfel spus toate sensurile cuvântului de tradus se regasesc în cuvântul reprezentând traducerea sa, necesitatea identificării sensului de utilizare al cuvântului sursa nu mai este obligatorie (cel puțin la nivelul fazei de transfer lexical) cu exceptia situatiei în care diferitele sensuri ale cuvântului tinta se realizeaza lingvistic prin structuri de subcategorizare distincte.

În aceasta lucrare vom prezenta în prima parte o metoda de extragere automata a echivalentilor de traducere si vom descrie apoi procedura de discriminare a sensurilor cuvintelor din corpusuri paralele pe baza echivalentilor de traducere.

2. Echivalenti de traducere

1.1. Notiuni preliminare

O pereche de texte în doua limbi diferite L_A si L_B , astfel încât unul reprezinta traducerea celuilalt constituie cea ce se numeste un *bitext*. Un bitext suficient de mare, constituie un corpus paralel. L_A si L_B se numesc echivalenti de traducere. Notiunea de echivalenta de traducere se poate rafina la niveluri subtextuale, de pilda la nivelul paragrafului, al propozitiei sau chiar la nivel lexical, al cuvântului sau al expresiei. În continuare elementul de aliniere lexicala îl vom numi, generic, *atom lexical* sau simplu *atom*. Un bitext în care echivalentii de traducere sunt explicitati se numeste un bitext *aliniat*. Cea mai mica unitate textuala la nivelul careia se realizeaza alinierea defineste *granularitatea*

echivalentilor de traducere. Echivalentii lexicali de traducere (obiectul nostru de interes în aceasta lucrare) depind evident de bitextul din care sunt extrasi iar procesul de extragere a lor devine echivalent cu extragerea unui dictionar bilingv, specific unui anumit domeniu⁸². Extragerea unui dictionar de echivalenti de traducere dintr-un bitext este în fond un proces de explicitare a dictionarului mental folosit de translatorul (sau translatorii) textului original.

Presupozitia fundamentala în încercarea de a alinia corpusurile paralele este ca *aceiasi* semnificatie este exprimata în doua sau mai multe limbi. Definirea identitatii de înțeles între doua sau mai multe reprezentari ale (presupus) aceluiași lucru este o binecunoscuta problema filozofica care ramâne deschisa chiar în domenii mult mai precise decât cel al limbii (de pilda în ingineria software). Prin urmare, notiunea de echivalent de traducere este un concept vag, si pentru operationalizarea sa în domenii ca traducerea automata, terminologie, managementul multilingual al documentelor si altele asemenea avem nevoie de o definitie precisa în termeni direct cuantificabili. Una dintre cele mai larg acceptate definitii a echivalentei de traducere este cea folosita în (Melamed, 2001): „the translation equivalence defines a (symmetric) relation that holds between two different language texts such that expressions appearing in corresponding parts of the two texts are reciprocal translations. These expressions are called *translation equivalents*”.

Majoritatea abordarilor moderne în extragerea automata a echivalentilor de traducere⁸³, sprijinite de forta de calcul din ce în ce mai mare a calculatoarelor, utilizeaza metode statistice si pot fi clasificate în doua mari categorii:

- paradigma „*presupune si testeaza*” (Gale, Church, 1991), (Smadja et al., 1996) etc., se bazeaza pe generarea unei multimi de potentiali echivalenti de traducere (spatiul ipotezelor) din care se selecteaza ulterior, pe baza unor teste de independenta statistica, echivalentii de traducere. Selectarea fiecarui echivalent de traducere se face independent de echivalentii extrasi anterior (procesul poate fi considerat ca fiind unul de optimizare locala).
- paradigma „*modelului de limba*” (Brown et al., 1993), (Kupiec, 1993), (Hiemstra, 1997) etc. presupune construirea unui model statistic al bitextului, model ai carui parametri se estimeaza, prin metode de optimizare globala. În aceasta abordare un candidat supus estimarii nu mai este o pereche de atomi lexicali ci o multime de perechi, numita *assignare* (Brown et al., 1993).

Exista sustinatori si critici ai ambelor abordari si o discutie a avantajelor si dezavantajelor lor este prezentata în (Hiemstra, 1997). În esenta, paradigma „*presupune si testeaza*” este mult mai eficienta din punct de vedere computational deoarece presupune investigarea unui spatiu al solutiilor proportional cu N^2 , unde N este maximul dintre

⁸² Posibilitatea de a genera automat dictionare bilingve în domenii specializate, coroborata cu performantele tot mai bune ale programelor de clasificare automata a textelor, deschide noi perspective traducerii automate si în general prelucrării multilinguale a textelor.

⁸³ În continuare, daca nu vom specifica altminteri, prin „echivalenti de traducere” vom înțelege implicit „echivalenti lexicali de traducere”

numerele de articole lexicale distincte din cele doua parti ale bitextului, dar echivalentii de traducere cu numar mic de aparitii sunt de obicei pierduti. Paradigma „*modelului de limba*” este extrem de costisitoare din punct de vedere computational întrucât spatiul de cautare al solutiilor este teoretic proportional cu $N!$, în schimb având potentialitatea identificarii corecte chiar si echivalentilor de traducere cu o singura aparitie în bitext (*hapax-legomena*). În (Brown et al., 1993) sunt prezentati o serie de algoritmi foarte eficienti, bazati pe o serie de ipoteze simplificatoare dar rationale, ce permit o ignorarea unor mari regiuni din spatiul de cautare, regiuni în care este improbabil sa existe solutii acceptabile.

Metoda descrisa aici poate fi încadrata în categoria abordarilor de tip „*presupune si testeaza*”. Algoritmul genereaza mai întâi o lista de candidati si apoi succesiv, alege din aceasta lista perechile cu cele mai mari scoruri de co-ocurenta în regiuni corespondente ale bitextului. Dupa cum se va vedea în continuare, acest algoritm nu are nevoie de un dictionar bilingv initial, dar daca acesta exista, utilizarea sa poate spori substantial viteza si acuratetea prelucrării.

2.2. Ipoteza corespondentei lexicale 1:1

În general, un cuvânt dintr-un segment ce apare într-o parte a bitextului se traduce în segmentul corespunzator din cea de a doua parte a bitextului tot printr-un singur cuvânt. Daca acest lucru s-ar întâmpla întotdeauna, problema alinierii lexicale a unui bitext ar fi substantial mai simpla decât în realitate. Din pacate ipoteza „*cuvânt la cuvânt*” nu este adevarata în foarte multe cazuri, astfel încât adoptarea ei ca premisa de calcul nu pare foarte promitatoare. Dificultatea poate fi însa ocolita prin considerarea ipotezei conform careia **un articol lexical** dintr-o limba se traduce în cealalta tot printr-un **singur articol lexical**. Asa cum am aratat în sectiunea precedenta, un articol lexical este reprezentat fie de un cuvânt, fie de o secventa de cuvinte (sintagma, compus, expresie). Aceasta formulare, cunoscuta sub numele de „*ipoteza corespondentei lexicale 1:1*”, adoptata ca premisa computationala, simplifica mult problema tinta a alinierii lexicale a unui bitext, dar introduce probleme noi si anume definirea si respectiv recunoasterea automata a articolelor lexicale. Din fericire aceste probleme sunt reductibile la contexte monolingve si au solutii simple si foarte eficiente. Un program capabil sa realizeze recunoasterea automata a articolelor lexicale se numeste *segmentator lexical*. Un segmentator lexical este în general independent de limba, iar functionarea sa este controlata prin resurse specifice (dictionare continând cuvinte, secvente de cuvinte sau expresii regulate definite peste un vocabular limitat). În (Tufis&Barbu, 2001b) este discutata structura resurselor necesare segmentării lexicale a textelor în limba româna cu ajutorul segmentatorului *MtSeg*, dezvoltat la Universitatea Aix-en-Provence în cadrul proiectului european „Multext”.

Adoptarea „*ipotezei corespondentei lexicale 1:1*” reduce dramatic complexitatea problemei extragerii echivalentilor lexicali, indiferent de paradigma în care este abordata rezolvarea (a se vedea (Tufis&Barbu, 2001b, 2002) pentru detalii). Trebuie mentionat însa ca o segmentare lexicala perfecta (din punctul de vedere al utilitatii ei într-un context multilingv) este practic imposibila din cauza incompletitudinii inerente a oricarui dictionar

frazal. În (Tufis 2001b, Tufis&Barbu2002) se arata cum poate fi surmontata aceasta incompletitudine a resurselor necesare segmentarii lexicale.

2.3. Etape de preprocesare

2.3.1 Alinierea frazal?

Înainte de extragerea propriu-zisa a echivalentilor de traducere, un corpus paralel este supus unor prelucrari preliminare, de aducere a bitextului într-o forma normalizata. După ce fiecare parte a bitextului a fost supusa segmentarii lexicale, urmeaza etapa de aliniere la nivelul propozitiei a corpusului paralel. Pentru acest scop, am utilizat o varianta puțin modificata a algoritmului prezentat și documentat (Gale&Church, 1993). În (Tufis&Barbu, 2001b) este descris procesul de aliniere la nivel de fraza și furnizate exemple și statistici pentru diferite perechi de limbi prezente în corpusul paralel multilingv „1984”, conținând traduceri în șase limbi ale romanului omonim al lui George Orwell. Acolo aratam ca, în marea majoritate a cazurilor, traducerile din limba engleza s-au realizat în celelalte limbi pastrând corespondenta de 1:1 la nivelul frazei⁸⁴, cu alte cuvinte, aproape întotdeauna o fraza din textul englezesc a fost tradusa ca o singura fraza în celelalte limbi reprezentate în corpusul paralel. Algoritmul de aliniere la nivelul frazei poate depista și acele cazuri în care traducerea s-a realizat fara pastrarea corespondentei 1:1. Astfel, au fost cazuri în care doua fraze sursa au fost traduse printr-o singura fraza, sau invers, când o fraza din limba engleza a fost tradusa prin 2 sau chiar 3 fraze în celelalte limbi. În cele ce urmeaza, indiferent de tipul de aliniere (1:1, 2:1, 1:2 etc.) vom numi portiunile aliniate la nivelul frazal, *unitati de traducere* (UT).

Ratiunea acestei etape de prelucrare consta în intuitia comuna ca elementele lexicale aflate în relatie de echivalenta de traducere se regasesc în frazele ce se constituie în unitati de traducere. Pe de alta parte, procesul de aliniere la nivelul frazei este mult mai simplu, pentru ca în general indiferent de perechile de limbi considerate într-un bitext ordinea frazelor dintr-o limba este pastrata în cealalta limba. Aceasta ipoteza, operationalizata de un algoritm de optimizare dinamica de genul celui descris în (Gale&Church, 1993), permite printre altele și identificarea portiunilor netraduse într-una din limbi (alinieri de tipul N:0 sau 0:M).

O alta ipoteza simplificatoare pentru procesul identificarii echivalentilor lexicali de traducere se bazeaza pe observatia ca în marea majoritate a traducerilor, categoriile gramaticale din limba sursa se conserva în limba tinta (Melamed, 2001). Cu alte cuvinte, un verb se traduce de obicei printr-un verb, un substantiv printr-un substantiv s.a.m.d. Melamed a numit o astfel de pereche de traducere, pereche de tip V, distingând-o de perechile de tip P, în care atomii lexicali în cele doua limbi au categorii gramaticale diferite. Melamed, distinge și o a treia categorie de perechi de traducere, tipul I, perechile de traducere incomplete, rezultate ca urmare a unei segmentari lexicale partiale și a

⁸⁴ *Notiunea de fraza este luata aici în sensul ei larg, al unei propozitii sau fraze (enunt terminat cu un semn de punctuatie din categoria celor finale: punct, punct și virgula, doua puncte, semnul exclamarii, semnul întrebării, trei puncte).*

utilizării „ipotezei de aliniere lexicală 1:1”. Considerațiile lui Melamed referitoare la distribuția celor trei tipuri de traduceri lexicale sunt foarte bine confirmate de experimentele noastre, în ciuda faptului că textul nostru este un text literar în timp ce textul sau este un text politic (dezbaterile din Parlamentul Canadian) conținând traduceri literale, mult mai puțin afectate de personalitatea literară a traducătorului. Ceea ce este demn de remarcat este că perechile de tip P nu conțin categorii gramaticale arbitrare, și că de la o pereche de limbi la alta, se pot identifica regularități în alternanța categoriilor gramaticale la traducere (de ex. participiu-adjectiv, gerunziu-substantiv, gerunziu-adjectiv). Astfel de regularități pot fi abstractizate prin expresii regulate, efectul net fiind că multe din perechile de tip P pot fi asimilate (algoritmice) perechilor de tip V. Prin urmare, necesitatea identificării rapide și precise a categoriei gramaticale (și eventual al altor trăsături morfologice sau lexicale) pentru atomii lexicali dintr-un bitext impune o altă prelucrare preliminară, respectiv etichetarea morfo-lexicală, prelucrare pe care o prezentăm în secțiunea următoare.

2.3.2 Etichetarea morfo-lexicală? și lematizarea

Etichetarea morfo-lexicală este procesul prin care fiecărui articol lexical dintr-un text arbitrar i se atribuie un cod morfo-lexical unic dintr-o mulțime specifică articolului lexical respectiv, numită clasă sau de ambiguitate. Codul morfo-lexical reprezintă o reprezentare compactă, și de obicei standardizată, a proprietăților morfologice și lexicale ce caracterizează apariția unui atom lexical într-un text. Clasa de ambiguitate a unui atom lexical reprezintă mulțimea tuturor interpretărilor posibile în orice context legal al atomului respectiv. De exemplu cuvântul "*urâți*" are cel puțin 8 interpretări posibile putând fi substantiv, adjectiv sau verb. Lema sa poate fi una dintre "*urât*" (substantiv sau adjectiv), "*a urâți*" sau "*a urî*" (verb).

urâți	urâți	Vmnp	(inf.: A <i>urâți</i> înseamnă a face să devină urât)
urâți	urâți	Vmis3s	(ind., perf.simplu, sing., pers. 3: El <i>urâți</i> totul în viața ei)
urâți	urâți	Vmm-2s	(imp., sing: Prietene, nu <i>urâți</i> singurul lucru frumos din viața lui!)
urâți	urî	Vmip2p	(ind., prez., pl., pers. 2: De pomană îi <i>urâți</i> pe ei, ceilalți sunt de vină)
urâți	urî	Vmsp2p	(subj., prez., pl., pers. 2: Voi ar trebui să <i>urâți</i> tot ce e împotriva vieții)
urâți	urî	Vmm-2p	(imp., sing: Nu-i <i>urâți</i> pe apărătorii planetei!)
urâți	urât	Afpmp-n	(adj., masc. pl., neart. : Doi câini <i>urâți</i> și răi păzeau intrarea.)
urâți	urât	Ncmp-n	(subs. com., masc. pl., neart.: Niște <i>urâți</i> m-au băgat în sperieți.)

Asadar, clasa de ambiguitate a cuvântului "*urâți*" este mulțimea (Vmnp, Vmis3s, Vmm-2s, Vmm-2p, Vmip2p, Vmsp2p, Afpmp-n, Ncmp-n), iar etichetarea morfo-lexicală a acestui cuvânt înseamnă a alege, în funcție de contextul apariției sale, unul și numai unul

dintre cele 8 coduri reprezentând interpretarea contextuală a cuvântului. În cercetările anterioare am dezvoltat o metoda statistica de etichetare morfo-lexicală (Tufis, 1999), numita etichetarea cu doua niveluri si modele de limba combinate (TT-CLAM: tiered-tagging with combined language models), bazata pe programul TnT al lui Thorsten Brants (Brants, 2000) de prelucrare a modelelor markov cu legaturi ascunse de ordin 2 (3-gram HMM), program ce poate fi descarcat de la adresa www.coli.uni-sb.de/~thorsten/tnt/. Abordarea TT-CLAM a aratat ca texte arbitrare în limba română pot fi etichetate morfo-lexical în mod corect în peste 98.5% din cazuri si ca atunci când de interes este numai categoria gramaticală, procentul de etichetare corectă depășește 99.5%. Metoda TT-CLAM s-a dovedit independentă de limba, rezultate mai bune decât în alte abordări fiind raportate în literatura de specialitate pentru limbi foarte diferite de limba română: limba maghiară (Varadi, 2002, Oravecz *et al.*, 2000, Tufis *et al.*, 2000) limba germană (Hinrics&Truskina, 2002), Slovene (Erjavec, 2002).

Lematizarea este procesul prin care o formă flexionată a unui articol lexical (cuvânt sau expresie) este redusă la forma normală de dicționar. Lematizarea se poate realiza fie printr-un proces de analiză morfologică fie prin căutarea într-o bază de date lexicale, conținând cuvinte în formă flexionată însoțite de analiză lor morfologică și de forma lema. Lematizarea se realizează în acest caz prin identificarea în baza de date a lemei pentru care forma flexionată și analiză morfo-lexicală sunt identice cu cele din textul de lematizat, care desigur a fost în prealabil etichetat. Pentru limba română, noi am experimentat cu ambele metode și datorită vitezei mult superioare, am optat pentru varianta a doua.

În figura de mai jos este exemplificat rezultatul prelucrărilor preliminare discutate în această secțiune (segmentare lexicală, aliniere frazală, etichetare morfo-lexicală și lematizare) pentru începutul bitextului Englez-Român din corpusul multilingv „1984”. Prima linie arată că în limba română, fraza cu identificatorul Oro.1.2.2.1, reprezintă traducerea a două fraze din textul englezesc, respectiv a celor cu identificatorii Oen.1.1.1.1 și Oen.1.1.1.2 (avem deci o aliniere de tip 1:2). Linii următoare, specifice pentru fiecare articol lexical din fiecare limba tipul sau (TOK, LSPLIT, DATE, ABR etc.), forma ortografică, lema, codul morfo-lexical și categoria gramaticală (ultimele 3 separate prin caracterul „\”).

<link targets="Oro.1.2.2.1; Oen.1.1.1.1 Oen.1.1.1.2">

(<S FROM="Oro.1.2.2.1">		(<S FROM="Oen.1.1.1.1">
LSPLIT	Într-	Întru\Spsay\S
TOK	o	un\Tifsr\T
TOK	zi	zi\Ncfsrn\N
TOK	senina	senin\Afpfsrn\A
...		...
...		</S>
...		<S FROM="Oen.1.1.1.2">
...		...
</S>		</S>

))

Figura 1: Bitext preprocesat pentru extracția echivalentilor lexicali de traducere

O descriere a principiilor de codificare morfo-lexicală, în conformitate cu recomandările EAGLES poate fi găsită în Erjavec and Ide (1998). Codificarea specifică pentru limba română, conformă cu standardul respectiv este pe larg descrisă în (Tufis *et al.*, 1997).

2.4. Un prim algoritm de extragere automată a echivalentilor lexicali de traducere

Există, așa cum am văzut mai sus, mai multe ipoteze simplificatoare care permit ținerea sub control a complexității problemei extragerii automate a echivalentilor de traducere. Nici una dintre aceste ipoteze nu este valabilă în cazul general, dar situațiile în care ele nu sunt adevărate sunt suficient de rare astfel încât adoptarea lor nu alterează valoarea rezultatelor. Trebuie subliniat faptul că ipotezele simplificatoare folosite de noi, discutate anterior și rezumate în continuare, în general nu afectează precizia (corectitudinea) dictionarelor bilingve extrase și completitudinea lor. Altfel spus, o serie de perechi corecte (echivalenți de traducere reali), deși prezente în bitext, pot să nu sunt găsite. Precizia și completitudine (în limba engleză acești termeni sunt *precision* și *recall*) se definesc în mod standard astfel:

PREC=(număr de echivalenți corect extrasi)/(număr total de echivalenți extrasi)

COMP=(număr de echivalenți corect extrasi)/(număr total de echivalenți existenți în bitext)

Mai trebuie precizat și faptul că ipotezele simplificatoare enumerate mai jos nu împiedică recuperarea ulterioară a echivalentilor negasiți din cauza adoptării acestor ipoteze de lucru. În (Tufis, 2000) sunt discutate metode de recuperare a unor echivalenți de traducere ce nu respectă ipoteza „echivalenței lexicale 1:1”.

- ipoteza „echivalenței lexicale 1:1”; ea stă la baza majorității abordărilor cunoscute: (Kay & Röscheisen, 1993), (Brew & McKelvie, 1996), (Hiemstra, 1997), (Tiedemann, 1998), (Ahrenberg et al., 2000), (Melamed, 2001), etc. Așa cum am arătat mai devreme, un articol lexical identificat corespunzător de un segmentator lexical adecvat diminuează considerabil efectul contrazicerii acestei ipoteze;
- un articol lexical polisemantic ce apare de mai multe ori în aceeași unitate de traducere este folosit cu același înțeles; această presupozitie este explicit utilizată de (Melamed, 2001) și implicit de toți cercetătorii amintiți mai sus;

- un articol lexical dintr-o parte a unitatii de traducere UT poate fi aliniat unui articol lexical în cealalta parte a UT doar daca cele doua articole au categorii gramaticale compatibile; în majoritatea cazurilor compatibilitatea categoriilor gramaticale se reduce la identitate, dar cum am specificat anterior, este posibil sa se defineasca corespondente compatibile (de pilda, verbele la participiu si gerunziu din limba engleza sunt destul de frecvent traduse în limba româna ca adjective sau substantive, si reciproc).
- Desi ordinea cuvintelor nu este un invariant al traducerii, ea nu este nici arbitrara; când doua sau mai multe perechi de articole lexicales candideaza la statutul de echivalenti de traducere, iar alte criterii de evaluare nu permit departajarea lor, atunci este preferata perechea continând articolele cele mai apropiate în pozitile lor relative. Aceasta euristica este, de asemenea, folosita de [Ahrenberg et al., 2000].

Pe baza bitextului preprocesat asa cum s-a prezentat în sectiunea precedenta, primul pas al algoritmului este de a delimita spatiul de cautare al solutiilor. Acest lucru se realizeaza prin constructia unei liste a tuturor candidatilor posibili (în conformitate cu ipotezele de lucru amintite mai sus). Aceasta lista, pe care o notam cu TECL (Translation Equivalence Candidates List) contine la rândul ei o multime de sub-liste (câte una pentru fiecare categorie gramaticala luata în considerare). Fiecare sublista contine perechi de forma $\langle token_S, token_T \rangle$ unde $token_S$ si $token_T$ sunt articole lexicales de categorii gramaticale compatibile si care au aparut în partile corespunzatoare ale aceleiasi unitati de traducere. Fie TU^j cea de a j^{th} unitate de traducere (translation unit). Prin colectarea tuturor articolelor lexicales aparținând aceleiasi categorii gramaticale POS_k (pastrând ordinea lor relativa si eliminând duplicatele) se construiesc pentru fiecare TU^j multimele ordonate $L_{POS_k}^{S_j}$ si $L_{POS_k}^{T_j}$. Pentru fiecare POS_i , fie $TU_{POS_i}^j$ produsul cartezian $L_{POS_i}^{S_j} \otimes L_{POS_i}^{T_j}$. Atunci, definim lista de corespondente în unitatea de traducere TU^j ca fiind CTU^j (correspondences in the j^{th} translation unit):

$$CTU^j = \bigcup_{i=1}^{no.of.pos} TU_{POS_i}^j$$

Cu aceste notatii, si presupunând ca bitextul de intrare contine n unitati de aliniere, atunci TECL se defineste astfel:

$$TECL = \bigcup_{j=1}^n CTU^j$$

TECL contine desigur foarte mult „zgomot” si cele mai multe perechi candidate (TEC=Translation Equivalence Candidate) sunt extrem de improbabile. Pentru a elimina cât mai multe din perechile TEC improbabile, TECL este filtrata pe baza unor functii scor ce supun fiecare TEC la o analiza a ipotezei statistice de independenta a asocierii articolelor lexicales. Pentru a prezenta functiile scor pe care le-am utilizat în experimentele noastre, vom mai defini o serie de notatii:

- $TEC = \langle T_S, T_T \rangle \in TECL$, un potential echivalent de traducere definit ca perechea formata din articolul lexical sursa T_S si posibila sa traducere T_T în limba tinta;
- n_{11} = numarul de ocurente ale $\langle T_S, T_T \rangle$ din TECL;
- n_{12} = numarul de perechi $\langle T_S, \neg T_T \rangle$ din TECL în care T_S a fost asociat cu un articol lexical diferit de T_T ;
- n_{21} = numarul de perechi $\langle \neg T_S, T_T \rangle$ din TECL în care T_T a fost asociat cu un articol lexical diferit de T_S ;
- n_{22} = numarul de perechi $\langle \neg T_S, \neg T_T \rangle$ din TECL ce nu contin nici pe T_S si nici pe T_T ;
- n_{1*} = numarul de perechi $\langle T_S, * \rangle$ din TECL în care apare T_S indiferent cu cine este asociat;
- n_{*1} = numarul de perechi $\langle *, T_T \rangle$ din TECL în care apare T_T indiferent cu cine este asociat;
- n_{2*} = numarul de perechi $\langle \neg T_S, * \rangle$ din TECL în care T_S nu apare;
- n_{*2} = numarul de perechi $\langle *, \neg T_T \rangle$ din TECL în care T_T nu apare;
- n_{**} = numarul de perechi $\langle *, * \rangle$ din TECL;

Tabela de contingenta din figura de mai jos ilustreaza aceste notatii:

	T_T	$\neg T_T$	
T_S	n_{11}	n_{12}	n_{1*}
$\neg T_S$	n_{21}	n_{22}	n_{2*}
	n_{*1}	n_{*2}	n_{**}
	$n_{1*} = n_{11} + n_{12}, n_{2*} = n_{21} + n_{22}$		
	$n_{*1} = n_{11} + n_{21}, n_{*2} = n_{12} + n_{22}$		
	$n_{**} = \sum_{j=1}^2 \sum_{i=1}^2 n_{ij}$		

Figura 2: Tabela de contingenta pentru un potential echivalent de traducere $\langle T_S, T_T \rangle$

Pentru ordonarea potentialilor echivalenti de traducere în vederea filtrarii (eliminarea candidatilor cei mai puțin plauzibili) am realizat experimente folosind 4 functii de calcul al scorului de echivalenta: MI (*informatia mutuala*), DICE, LL (log likelihood), and χ^2 (chi-patrat). Folosind notatiile de mai sus, aceste functii-scor se definesc în felul urmator:

$$(1) \quad MI(T_T, T_S) = \log_2 \frac{n_{**} * n_{11}}{n_{1*} * n_{*1}},$$

$$(2) \quad \text{DICE}(T_T, T_S) = \frac{2n_{11}}{n_{1*} * n_{*1}},$$

$$(3) \quad \text{LL}(T_T, T_S) = 2 * \sum_{j=1}^2 \sum_{i=1}^2 n_{ij} * \log \frac{n_{ij} * n_{**}}{n_{i*} * n_{*j}} \text{ si}$$

$$(4) \quad \chi^2(T_T, T_S) = n_{**} \sum_{j=1}^2 \sum_{i=1}^2 \frac{\left(n_{ij} - \frac{n_{i*} * n_{*j}}{n_{**}} \right)^2}{n_{i*} * n_{*j}}$$

Figura 3: Functii-scor pentru evaluarea unui potential echivalent de traducere $\langle T_S T_T \rangle$

O formula mai simpla de calcul pentru $\chi^2(T_T, T_S)$ este urmatoarea:

$$(4') \quad \chi^2(T_T, T_S) = \frac{n_{**}(n_{11} * n_{22} - n_{12} * n_{21})^2}{(n_{11} + n_{12}) * (n_{11} + n_{21}) * (n_{21} + n_{22}) * (n_{21} + n_{22})}$$

Filtrarea potentialilor echivalenti de traducere se face în raport cu un prag numeric impus scorului calculat cu una dintre functiile de mai sus. Toate perechile ce obtin un scor mai mare decât pragul ales sunt considerate plauzibile si vor fi supuse unor prelucrari suplimentare iar celelalte sunt eliminate. Orice metoda de filtrare statistica va elimina multi echivalenti falsi de traducere, dar pe lânga acestia si un numar de perechi corecte. Alegerea pragului de scor s-a facut avânt ca obiectiv minimizarea numarului de perechi corecte dar eliminate în mod gresit si a numarului de perechi incorecte acceptate ca urmare a scorului superior pragului de selectie. Dupa mai multe experimente, cele mai bune rezultate s-au obtinut folosind functia de scor LL cu limita pragului de acceptanta egala cu 9.

Într-o prima varianta, algoritmul nostru de extragere a echivalentilor de traducere, având unele asemanari cu algoritmul iterativ prezentat în (Ahrenberg et al. 2000), implementa o strategie de selectie indiferenta la locul si pozitia în corpus a articolelor lexicale aparând în perechea TEC analizata la un anumit moment. O diferenta majora fata de algoritmul descris în (Ahrenberg et al. 2000) este ca în programul nostru calculul diferitelor probabilitati (mai exact al estimatilor de probabilitate) si al scorurilor (testul t) devine nenecesar, conducând la o viteza de prelucrare cu cel puțin un ordin de marime mai mare. Pornind de la lista filtrata a potentialilor echivalenti de traducere, algoritmul selecteaza în mod iterativ cei mai plauzibili candidati (vezi mai jos) si apoi îi sterge din lista initiala. Algoritmul se opreste dupa un numar prestabilit de iteratii sau mai devreme în cazul în care lista candidatilor s-a golit sau daca nici un candidat nu mai indeplineste conditia de selectie.

În iteratia k a algoritmului se construiesc o matrice de contingenta (TBLk) pentru fiecare categorie gramaticala (POS) având dimensiunile $S_m * T_n$ unde S_m si T_n reprezinta numarul de articole lexicale din limba sursa respectiv tinta care mai exista în lista de

candidati la pasul k (Figura 4). Liniile si coloanele tabelului sunt indexate cu articolele lexicale (având aceeași categorie gramaticală) din limba sursă respectiv limba țintă. Fiecare celulă (i,j) a matricii reprezintă numărul de ocurențe în lista de candidați a perechii $\langle T_{Si}, T_{Tj} \rangle$.

	T_{T1}	...	T_{Tn}	
T_{S1}	n_{11}	...	n_{1n}	n_{1*}
...				...
	
T_{Sm}	n_{m1}	...	n_{mn}	n_{m*}
	n_{*1}	...	n_{*n}	n_{**}

$$n_{ij} = \text{occ}(T_{Si}, T_{Tj}); n_{i*} = \sum_{j=1}^n n_{ij}; n_{*j} = \sum_{i=1}^m n_{ij}; n_{**} = \sum_{j=1}^n \left(\sum_{i=1}^m n_{ij} \right).$$

Figura 4: Matricea de contingenta la pasul k

Condiția de selecție la pasul k a multimii de echivalenți de traducere este exprimată de relația (5):

$$(5) \quad TP^k = \left\{ \langle T_{Si}, T_{Tj} \rangle \mid \forall p, q (n_{ij} \geq n_{iq}) \wedge (n_{ij} \geq n_{pj}) \right\}$$

Condiția de mai sus constituie esența algoritmului iterativ (numit în (Tufis&Barbu, 2002) algoritmul BASE) și ea spune că pentru a selecta perechea $\langle T_{Si}, T_{Tj} \rangle$ drept echivalent de traducere, numărul de asocieri ale lui T_{Si} cu T_{Tj} trebuie să fie mai mare sau cel puțin egal decât numărul de asocieri ale lui T_{Si} cu orice alt T_{Tp} ($p \neq j$) și simultan numărul de asocieri ale lui T_{Tj} cu T_{Si} trebuie să fie mai mare sau cel puțin egal decât numărul de asocieri ale lui T_{Tj} cu orice alt T_{Sq} ($q \neq i$). Toate perechile selectate în TP^k sunt eliminate din lista de candidați (ceea ce în matricea de contingenta pentru pasul $k+1$ implică punerea pe 0 a contoarelor de ocurență pentru perechile selectate anterior). Dacă T_{Si} este tradus în mai multe moduri (fie pentru că are sensuri ce se lexicalizează diferit în limba țintă, fie pentru că în limba țintă se folosesc diferiți sinonimi pentru T_{Tj}) restul traducerilor sale va fi extras în iterațiile următoare. Algoritmul discutat este schitat în figura 5:

```

procedure BASE(bitext, step; dictionary) is
  k=1;
  TP(0)={};
  TECL(k)=build-cand(bitext);
  for each POS in TECL do
    loop

```

```

TECL(k)=update(TP(k-1),TECL(k))
TBL(k)=build_TEC_table(TECL(k));
TP(k)=select(TBL(k)); ## relatia (5) ##
add(dictionary, TP(k));
k=k+1;
until {(TECL(k-1) is empty)or(TP(k-1) is empty)or(k >
step)}
endfor
return dictionary
end

```

Figura 5: Algoritmul iterativ de extragere a echivalentilor de traducere

2.5. Un algoritm îmbunătătit de extragere automată a echivalentilor lexicali de traducere

Una dintre principalele deficiente ale algoritmului BASE este vulnerabilitatea la ceea ce (Melamed, 2001) numeste *asociatii indirecte*. Dacă $\langle T_{Si}, T_{Tj} \rangle$ are un scor de coocurență ridicat iar T_{Tj} apare (dintr-un motiv sau altul) de mai multe ori împreună cu T_{Tk} , s-ar putea ca și perechea $\langle T_{Si}, T_{Tk} \rangle$ să primească un scor ridicat. Deși, așa cum observă și Melamed, în general, asociatiile indirecte au un scor mai mic decât cele directe (corecte), ele pot obține totuși scoruri mai mari decât multe alte perechi corecte ce n-au legătură cu T_{Si} iar acest lucru nu numai că generează echivalenți de traducere greșiți, dar va elimina din competiție și echivalenți corecți. Prin urmare asociatiile indirecte afectează atât precizia cât și completitudinea procesului. Pentru a slăbi această sensibilitate în implementarea algoritmului BASE a fost nevoie de stabilirea unei limite inferioare de ocurență pentru fiecare articol lexical luat în considerare. Această limită, conduce inevitabil la eliminarea din spațiul de căutare a soluțiilor a mai mult de 50% dintre echivalenții de traducere⁸⁵. Deficiența algoritmului BASE se explică prin faptul că scorurile de coocurență sunt calculate în mod global fără a verifica dacă atomii lexicali ai unei perechi evaluate sunt sau nu prezenți în unitățile de traducere prelucrate.

Pentru diminuarea influenței asociatiilor indirecte fără a mai impune un prag de ocurență, algoritmul BASE a fost modificat astfel încât ierarhizarea și alegerea celor mai probabili echivalenți de traducere se realizează în contextul local al fiecărei unități de traducere (deși scorurile lor se calculează tot la nivelul întregului bitext). Cu această modificare, noul algoritm (BETA) se apropie de algoritmul „competitive linking” al lui Melamed (Melamed, 2001). Candidații proveniți din unitatea de traducere curentă sunt analizați prin prisma scorului lor de coocurență și cel cu scorul cel mai mare este selectat. În baza ipotezei corespondenței lexicale 1:1, dintre candidații rămași sunt eliminați toți aceia care conțin unul din articolele lexicale ale perechii câștigătoare. Dintre candidații care rămân după această filtrare, se alege din nou cel cu scorul cel mai bun și iar se elimină candidații conținând unul dintre articolele lexicale din perechea selectată. Procesul se

⁸⁵ Pierderea unui așa mare număr de echivalenți de traducere desigur nu surprinde, întrucât una din legile distribuționale ale lui Zipf (celebra lege „rang-frecvență”) prognosticează acest lucru.

repetă până când nici un echivalent de traducere nu mai poate fi extras din unitatea de traducere curentă, caz în care algoritmul trece la prelucrarea următoarei unități de traducere.

Eliminarea pragului de ocurență a îmbunătățit substanțial completitudinea și calitatea dictionarelor de traducere (o detaliată comparație a performanțelor și o analiză cantitativă și calitativă a dictionarelor extrase automat este furnizată în (Tufis&Barbu, 2002)) dar a ridicat problema decelării între candidații cu una sau două apariții, pentru care scorul de coocurență este statistic nesemnificativ. În acest caz, criteriul frecvenței a fost înlocuit cu o combinație între un scor de similaritate ortografică și un scor de proximitate relativă. Funcția de similaritate ortografică folosită de noi, $COGN(T_S, T_T)$, este o variantă a funcției **XXDICE** descrisă în (Brew&McKelvie, 1996). Astfel, dacă T_S este un sir de m caractere $\alpha_1\alpha_2 \dots \alpha_m$ și T_T un sir de n caractere $\beta_1\beta_2 \dots \beta_n$ se construiesc două noi siruri T'_S și T'_T prin inserarea în T_S și T_T a unui număr minim de caractere speciale astfel încât în final sirurile T'_S și T'_T au aceeași lungime p ($\max(m, n) \leq p < m+n$) și un număr maxim de caractere pozitional identice. Fie α_i un caracter din T'_S și β_i un caracter din T'_T care se potrivesc și sunt puse în corespondență. Fie $\delta(\alpha_i)$ numărul de caractere speciale consecutive ce preced imediat caracterul α_i și $\delta(\beta_i)$ numărul de caractere speciale ce preced imediat caracterul β_i . Fie q numărul de caractere care se potrivesc în cele două siruri. Cu aceste notații, măsura de similaritate $COGN(T_S, T_T)$ se definește astfel:

$$(6) \quad COGN(T_S, T_T) = \begin{cases} \frac{\sum_{i=1}^q \frac{2}{1 + |\delta(\alpha_i) - \delta(\beta_i)|}}{m+n} & \text{if } q > 2 \\ 0 & \text{if } q \leq 2 \end{cases}$$

Limita de relevanță a scorului de similaritate a fost empiric găsită a fi 0.42. Această valoare este dependentă într-o oarecare măsură de pereche de limbi considerată în procesul de extragere a echivalențelor de traducere. Implementarea efectivă a testului de similaritate include și o serie de normalizări ale sirurilor testate (eliminarea unor afixe, reducerea consoanelor duble, ignorarea distincției create de diacritice etc.) normalizări care depind de morfologia fiecărei limbi în parte.

Cel de al doilea criteriu de evaluare a plauzabilității unui candidat este scorul de proximitate, $DIST(T_S, T_T)$ definit după cum urmează:

Dacă ($\langle T_S, T_T \rangle \in L_{posk}^{S_j} \otimes L_{posk}^{T_j}$) & (T_S este al n -lea element în $L_{posk}^{S_j}$) & (T_T este al m -lea element în $L_{posk}^{T_j}$) atunci $DIST(T_S, T_T) = |n-m|$

Filtrul $COGN(T_S, T_T)$ este mult mai semnificativ din punct de vedere lingvistic⁸⁶ decât $DIST(T_S, T_T)$, astfel încât scorul de similaritate are precedență asupra celui de proximitate. Funcția $DIST(T_S, T_T)$ este invocată doar atunci când $COGN(T_S, T_T) = 0$ (deci

⁸⁶ Motivația se bazează pe intuiția conform căreia dacă în două fraze (în limbi diferite) ce reprezintă una traducerea celeilalte, apar cuvintele asemănătoare din punct de vedere ortografic, atunci este foarte rezonabil a presupune că ele au și aceeași semnificație, adică sunt cogneți.

când atomii lexicali nu prezintă similaritate ortografică și perechea $\langle T_s, T_t \rangle$ nu reprezintă o pereche singulară în corpus (hapax-legomena), sau când mai multe perechi candidate au obținut același scor de similaritate.

Algoritmul BETA este schitat mai jos:

```

procedure BETA(bitext;dictionary) is:
  dictionary={};
  TECL(k)=build_cand(bitext);
  for each POS in TECL do
    for each  $TU_{POS}^i$  in TECL do
      finish=false;
      loop
        best_cand = get_the_highest_scored_pairs( $TU_{POS}^i$ );
        conflicting_cand=select_conflicts(best_cand);
        non_conflicting_cand = best_cand\conflicting_cand;
        best_cand=conflicting_cand;
        if cardinal(best_cand)=0 then finish=true;
        else
          if cardinal(best_cand)>1 then
            best_cand=filtered(best_cand);
          endif;
        best_pairs = non_conflicting_cand + best_cand
        add(dictionary,best_pairs);
         $TU_{POS}^i$  =
        remove_pairs_containing_tokens_in_best_pairs( $TU_{POS}^i$ );
        endif;
      until {( $TU_{POS}^i$ ={}) or (finish=true)}
    endfor
  endfor
  return dictionary
end

procedure filtered(best_cand) is:
  result = get_best_COGN_score(best_cand);
  if (cardinal(result)=0)&(non-hapax(best_cand)) then
    result = get_best_DIST_score(best_cand);
  else if cardinal(result)>1
    result = get_best_DIST_score(best_cand);
  endif
  endif
  return result;
end

```

Din corpusul paralel multilingv „1984” am extras 6 bitexte conținând textul în limba engleză și traducerea în una din cele 6 limbi amintite. Fiecare bitext a fost prelucrat conform celor prezentate în acest capitol și au fost extrase 6 dicționare bilingve, din care s-a obținut și un dicționar multilingv în 7 limbi (cele 6 plus engleza). În (Tufis&Barbu, 2002) este furnizată o analiză contrastivă cu alte sisteme de acest tip a vitezei de prelucrare. Timpul mediu de extragere a unui dicționar bilingv din corpusul paralel multilingv „1984”

(circa 110.000 de cuvinte în fiecare limba) este 3 minute. Esantioane ale acestor dictionare pot fi consultate la adresa: <http://www.racai.ro/~tufis/BilingualLexicons/AutomaticallyExtractedBilingualLexicons.html>.

3. Dezambiguizarea sensurilor lexicale folosind echivalentele de traducere

3.1. Ambiguitatea limbajului natural

Este binecunoscut faptul ca una dintre cele mai dificile probleme în prelucrarea automată a limbajului natural este ambiguitatea sa inerentă. Ambiguitatea se manifestă la toate nivelurile tradiționale ale analizei de limbaj: nivelul fonetic și/sau lexical, sintactic, semantic sau discursiv. Ambiguitatea prezintă pe fiecare nivel generează exploziv ambiguități pe nivelurile următoare. De pildă, omofonia sau omografia prezintă pe primul nivel la nivelul unui sau al mai multor cuvinte va produce secvențe lexicale diferite (combinația tuturor interpretărilor posibile la acest nivel) pentru intrarea fazei de analiză sintactică. Fiecare secvență poate conduce, din pricina unor ambiguități de natură structurală, la interpretări sintactice multiple, după cum o serie de secvențe lexicale vor putea fi abandonate pe motivul contrazicerii unor restricții postulate de modelul sintactic al limbii prelucrate. Fiecare dintre interpretările sintactice posibile, poate la rândul ei să conducă la multiple interpretări semantice, în virtutea multiplelor sensuri pe care le poate avea fiecare element frazal al unei analize sintactice. Desigur, interpretarea semantică poate elimina unele structuri sintactice generate în faza anterioară pe baza încălcării unor restricții semantice (valabile în orice univers de discurs sau specifice unor domenii discursive de interes). În sfârșit, analiza de discurs, în care contextul interpretativ transcende limita propoziției, ambiguitățile rămase se presupune a putea fi rezolvate prin utilizarea restricțiilor pragmatice motivate fie de principii generale ale dialogului (coeziune, coerență), fie de natura bine precizată a unui univers de discurs (modelată prin cunoștințe extra-lingvistice despre entitățile universului de discurs). De pildă, în (Cristea&Dima, 2001) rezolvarea anaforelor, proces tipic analizei de discurs, este modelată în termenii identificării cailor de accesibilitate a entităților menționate în discurs („vene ale discursului”), care la rândul lor sunt formal definite pe baza principiilor generale ale coeziunii și coerenței unui text.

Rezultă din cele spuse până aici că identificarea și rezolvarea timpurie, la fiecare nivel de prelucrare, a ambiguităților este un imperativ al oricărui demers computațional privind prelucrarea limbajului natural. Și cum cuvântul (sau mai exact spus, atomul lexical) este elementul primar în prelucrarea limbajului o mare parte a eforturilor de cercetare este îndreptată spre nivelul lexical al prelucrarilor. Metodele de etichetare morfo-lexicală (tagging), printre care etichetarea cu două niveluri și modele de limbă combinate - amintite în capitolul 2, permit rezolvarea cu mare acuratețe a ambiguităților categoriale și intracategoriale. De pildă cuvântul *vin* poate fi atât substantiv cât și verb (ambiguitate

categoriala), iar ca verb, el contine ambiguitatea intracategoriala de persoana, numar si mod („indicativ + persoana I + numar singular”, „conjunctiv + persoana I + numar singular” sau „indicativ + persoana III + numar plural”). Un program de etichetare morfo-lexicala „instruit” corect pentru limba româna este capabil sa rezolve, în contextul aparitiei sale, astfel de ambiguitati morfo-lexicale.

Curentul lexicalist predominant în modelarea sintactica a limbajului natural presupune precizarea în descrierea de dictionar a fiecarui cuvânt a proprietatilor si restrictiilor sale distributionale sau colocationale relevante pentru analiza sintactica. Pe baza acestor descrieri lexicalizate si a contextului local, multe din potentialele ambiguitati structurale pot fi eliminate, înainte unei costisitoare analize sintactice, prin tehnici cunoscute sub numele de analiza sintactica partiala (*partial parsing* sau *shallow parsing*).

Un cuvânt omograf, chiar dupa ce a fost corect clasificat din punctul de vedere al categoriei sale gramaticale si al proprietatilor sale distributionale sau colocationale, poate ramâne ambiguu din punct de vedere semantic. Identificarea sensului cu care este utilizat cuvântul polisemantic într-un context dat este desigur de mare interes. Exista însa diferite grade de rafinare a notiunii de sens, iar natura aplicatiei pentru care identificarea sensului este necesara poate impune o acceptie a notiunii de sens diferita de cea utilizata într-un dictionar explicativ. Sa luam, de pilda, problema traducerii automate. Întrucât în imensa majoritate a cazurilor rezultatul traducerii este destinat uzului uman, ceea ce este important este ca în textul tradus sa nu apara ambiguitati suplimentare fata de cele din textul sursa. Cu alte cuvinte, daca o analiza algoritmica evidentiaza în limba sursa o serie de ambiguitati si pornind de la premiza ca textul este admisibil pentru vorbitorii nativi ai limbii textului sursa, de cele mai multe ori este nenaturala o traducere ce încearca sa evite total ambiguitatea identificata. La nivel lexical, aceasta revine la a spune ca daca diferitele sensuri ale unui cuvânt din limba sursa nu se lexicalizeaza prin cuvinte diferite în limba tinta, este neproductiva o încercare a diferentierii sensului contextual, atâta timp cât indiferent care ar fi el, traducerea cuvântului respectiv în limba tinta este aceeasi. Sau cu alte cuvinte, sensurile unui cuvânt din limba sursa ce se regasesc împreuna într-un cuvânt al limbii tinta nu necesita obligatoriu diferentierea pentru traducere. De exemplu, cuvântul englezesc „bottle” are în Wordnet1.5 (Fellbaum, 1998) doua sensuri (ca substantiv) anume de vas de sticla sau plastic cilindric cu un gât îngust si fara mâner, respectiv cantitatea de substanta continuta într-un astfel de vas. Ambele sensuri se regasesc în cuvântul românesc „sticla” (care însa include si alte sensuri lexicalizate în engleza prin cuvântul „glass”). În acest caz, a încerca eliminarea ambiguitatii la traducerea textului „He drank only a bottle of beer” în limba româna, de pilda prin utilizarea unei parafraze de genul „El bause doar continutul unei sticle de bere”, este nenecesara. Orice vorbitor al limbii române va gasi traducerea „El bause doar o sticla de bere” mult mai naturala si desigur nu va avea dificultati a în înțelege despre ce este vorba.

Acelasi gen de consideratii se pot face si în raport cu ambiguitatile sintactice pure. Celebrul exemplu „I saw the Statue of Libery flying over New York” contine cel puțin 4 ambiguitati, dar daca de pilda rezolvarea omografului *saw* (am vazut / tai cu fierastraul)

este esentiala în traducere, rezolvarea ambiguitatii structurale poate fi lasata în sarcina mintii celui ce citește textul: „Am vazut Statuia Libertatii zburând deasupra New York-ului”, caci daca cititorul englez nu are dificultati în a înțelege cine si cum zbura, e plauzibil ca nici cititorul român (de exemplu) ne le va avea. Aceasta nu înseamna ca nu exista ambiguitati structurale a caror nerezolvare prealabila sa nu conduca la traduceri hazlii sau chiar incompreensibile. Ideea este ca metodele formale de analiza a limbajului, modelabile algoritmic, expliciteaza de multe ori ambiguitati greu de constientizat de omul obisnuit, iar luarea în considerare a factorului uman poate simplifica mult prelucrarile automate. Reconsiderarea conceptului de traducere automata în acceptiunea clasica (MT) în favoarea unor concepte mai realiste de tipul HAMT (human assisted machine translation) sau MAHT (machine assisted human translation) a relevat faptul ca, în numeroase ocazii, posteditarea umana a unui text tradus automat introduce ambiguitati care, desi nu sunt sezizabile usor la lectura, pot fi totusi puse în evidenta de algoritmi de analiza.

Cercetarile moderne în domeniul dezambiguizarii automate, în context, a sensurilor cuvintelor sunt motivate si de alte aplicatii informatice, cum ar fi clasificarea dupa continut a volumelor mari de texte, regasirea mai precisa a documentelor electronice, rezumarea automata a textelor, extragerea de cunostinte din texte, crearea de ontologii. Aceasta directie de cercetare, identificata în literatura engleza prin acronimul WSD (Word Sense Disambiguation) constituie de câtiva ani obiectul unor conferinte specializate si chiar a unei competitii de evaluare (SENSEVAL, ajunsă la a treia editie) a solutiilor propuse de specialisti din întreaga lume.

Primii care au sugerat ideea ca, pentru obiectivele WSD, sensurile ce trebuie differentiate sunt cele care se lexicalizeaza într-o alta limba prin cuvinte diferite au fost Resnik and Yarowsky (1997). Intuitiv, se poate presupune ca, daca un cuvânt din limba sursa se traduce în limba tinta în mai multe feluri si aceste traduceri nu sunt sinonimice, atunci trebuie sa existe o motivatie conceptuala. Analizând un numar suficient de mare de limbi si de texte, e plauzibil, afirmau cei doi specialisti, sa identificam diferentierile lexice semnificative care delimiteaza sensurile unui cuvânt. Aceste sensuri sunt numite de cei doi „*sensuri tari*”. Inabilitatea de a identifica corect sensurile tari este principala sursa a erorilor inacceptabile în orice aplicatie multilinguala. Utilizarea textelor paralele pentru WSD (Gale *et al.*, 1993), (Dagan *et al.*, 1991), (Dagan and Itai, 1994), în scopul identificarii proprietatilor semantice a lexemelor si a relatiilor dintre ele (Dyvik, 1998) a folosit implicit sau explicit notiunea de „sens tare”. Mai recent, pe baza echivalentilor de traducere extrasi din corpusul „1984” prin procedura noastra, descrisa în capitolul precedent, Ide (1999) a aratat ca diferentele de traducere în 5 limbi (din 4 familii diferite) pot constitui un criteriu extrem de eficace în identificarea sensurilor tari în limba de pornire (în acest caz, engleza). Resnik and Yarowsky (2000) au folosit în schimb traducerile unor propozitii izolate în limba engleza efectuate de vorbitori nativi ai limbilor tinta, dar în mare concluziile studiului lor au fost aceleasi cu ale lui Ide. În ambele studii amintite referinta pentru limba engleza a fost WordNet (Miller *et al.*, 1990) si desi rezultatele lor sunt promitatoare, mai ales pentru sensurile tari, ele se bazeaza pe o multime prestabilita de sensuri. Date fiind divergentele semnificative între distinctiile de sens realizate în dictionarele (monolingve)

existente, precum și inexistenta unui acord general asupra gradului de rafinare a descrierilor de sens în practica lexicografică internațională, raportarea la un inventar prestabilit de sensuri, cel puțin din perspectiva prelucrării automate a limbajului, nu pare a fi o soluție optimă. În continuare, vom prezenta o abordare alternativă, detaliată în (Ide *et al.* 2001, Ide *et al.* 2002).

3.2. Discriminarea automată a sensurilor lexicale: metodologia

Metoda pe care o vom descrie este menită să identifice sensurile distincte cu care unul sau mai multe cuvinte apar într-un text dat. Întrucât este foarte improbabil ca într-un text omogen, chiar foarte lung (de pildă un roman), un cuvânt să fie folosit în toate sensurile sale, metoda desigur va identifica, prin analiza textuală descrisă în continuare, doar acel sens sau acele sensuri cu care este folosit cuvântul respectiv în textul prelucrat. La limita prin prelucrarea unor texte foarte diferite este posibil teoretic să fie identificate toate sensurile atestate ale unui anumit cuvânt.

Din punct de vedere metodologic, studiul nostru s-a bazat pe corpusul paralel multilingv „1984” și pe dicționarul multilingv extras din acest corpus. Cele 7 limbi ale experimentului nostru fac parte din patru familii: germanică (engleză), romanică (română), slavică (bulgară, cehă și slovenă) și ugro-finică (estoniană, maghiară). Deși corpusul conține un text beletristic, textul orwelian ca și traducerea sa în celelalte limbi nu sunt foarte stilizate și, ca atare, oferă un esantion rezonabil de limbă modernă, comună. Mai mult, traducerea textului original, efectuate de traducători avizați (unii dintre ei fiind apreciați scriitori), reflectă riguros originalul: pentru mai mult de 95% din textul englezesc o frază sursă este tradusă în celelalte limbi tot ca o singură frază. Tipurile de aliniere frazale existente în corpusul „1984” sunt prezentate în tabelul de mai jos și discutate în (Tufis&Barbu, 2001b):

Estoniana-Engleza			Maghiara-Engleza			Romanian-Engleza		
Tip	Nr.	Proc	Tip	Nr.	Proc	Tip	Nr.	Proc
3-1	2	0.030321%	7-0	1	0.014997%	3-1	3	0.046656%
2-2	3	0.045482%	4-1	1	0.014997%	2-4	1	0.015552%
2-1	60	0.909642%	3-1	7	0.104979%	2-3	3	0.046656%
1-3	1	0.015161%	3-0	1	0.014997%	2-2	2	0.031104%
1-2	100	1.516070%	2-1	108	1.619676%	2-1	85	1.321928%
1-1	6426	97.422681%	1-6	1	0.014997%	2-0	1	0.015552%
1-0	1	0.015161%	1-5	1	0.014997%	1-5	1	0.015552%
0-2	1	0.015161%	1-2	46	0.689862%	1-3	14	0.217729%
0-1	2	0.030321%	1-1	6479	97.165573%	1-2	259	4.027994%
			0-4	1	0.014997%	1-1	6047	94.043551%
			0-2	3	0.044991%	0-3	2	0.031104%
			0-1	19	0.284943%	0-2	2	0.031104%
						0-1	10	0.155521%

Bulgară- Engleză			Cehă- Engleză			Slovena- Engleza		
2-2	2	0.030017%	4-1	1	0.015029%	3-3	1	0.014970%
2-1	23	0.345190%	3-1	2	0.030057%	2-1	48	0.718563%
1-2	72	1.080594%	2-1	109	1.638112%	1-5	1	0.014970%
1-1	6558	98.424134%	1-3	2	0.030057%	1-2	53	0.793413%
0-1	8	0.120066%	1-2	81	1.217313%	1-1	6572	98.383234%
			1-1	6438	96.753832%	1-0	2	0.029940%
			0-1	21	0.315600%	0-1	3	0.044910%

Figura 6: Distribuția tipurilor de aliniere frazala în corpusul paralel ”1984”

Alinierea de tipul N:M reprezintă situațiile în care M fraze din limba engleză au fost traduse cu N fraze în limba respectivă. Un caz particular îl reprezintă situațiile de omisiune în traducere (0:M) sau de inserare de text fără corespondent în original (N:0).

3.3. Experimentul initial

Textul original „1984” conține 7.069 leme diferite, iar dicționarul multilingv extras prin metoda descrisă în prima parte a acestei lucrări conține 1.233 de intrări. Aceste intrări au fost reținute respectând condiția ca un articol lexical din limba engleză să aibă traduceri (eventual multiple) în cât mai multe limbi țintă. Condiția impusă dicționarului multilingv este foarte restrictivă, având în vedere că majoritatea dicționarelor bilingve extrase automat conțin între 6000 și 7000 de intrări. Intrări tipice (parțiale) în dicționarul multilingv sunt ilustrate în figura 7. O informație suplimentară, ce nu apare în exemplificarea din figura 7 este multimea tuturor unităților de traducere din corpusul paralel în care cuvântul englezesc a fost tradus prin echivalenții săi listate în dicționar. Dintre aceste intrări, au fost selectate 845 pentru care s-au găsit una sau mai multe traduceri în toate limbile. Dintre acestea, s-a ales o multime de 33 de substantive, acoperind toate gamele de frecvență și ambiguitate, cu care s-a realizat experimentul ale cărui rezultate au fost validate de experți umani (Ide, *et al.*, 2001).

Engleză	Categorie	Bulgară	Cehă	Estoniană	Maghiară	Română	Slovenă
...
finally	R	îredr?	nakonec konečně	lõpuks viimaks	végül	în_cele_di n_urmă până_la_ur mă	končen nazadnje
...
wealth	N	áfármnnâi	bohatství	jõukus	jólét	avuție	blaginja

	ăăraî		rikkus	gazdagság	bogăție	bogastvo
--	-------	--	--------	-----------	---------	----------

Figura 7: Exemple de echivalenții de traducere identificate în corpusul paralel „1984”

Pentru fiecare substantiv din acest esantion au fost extrase toate frazele englezești în care apare, împreună cu toate frazele corespunzătoare din celelalte limbi și pentru fiecare ocurență a sa a fost construit un vector binar reprezentând toate traduceri posibile ale cuvântului respectiv. O valoare 1 în poziția n a acestui vector semnifică faptul că acea ocurență a fost tradus prin cuvântul ce reprezintă a n -a traducere posibilă. O valoare 0 semnifică faptul că a n -a traducere posibilă nu a fost folosită. De pildă pentru substantivul „wealth” (vezi figura 7) au fost depistate 11 traduceri posibile (2 în bulgară, estoniană maghiară română și slovenă, 1 în cehă). Un vector asociat oricărei ocurențe a lui *wealth* va avea prin urmare 11 poziții. Astfel, dacă a m -a apariție în textul original al romanului „1984” a cuvântului *wealth* are atașat vectorul 10101010101 acest lucru semnifică faptul că în varianta bulgărească el a fost tradus cu *ăăarnnăî*, în cea cehă cu *bohaství*, în cea estoniană cu *rikkus*, în cea maghiară cu *gazdagság*, în cea română cu *bogăție* iar în cea slovenă cu *bogastvo*. Vectorii astfel definiți au fost prelucrați cu un algoritm de clasificare de tip aglomerativ (Stolcke, 1996), clasele rezultate fiind considerate a reprezenta sensuri distincte în care cuvântul curent a fost folosit de-a lungul romanului. Clasele produse de algoritm au fost comparate cu clasele rezultate prin dezambiguizarea manuală efectuată, independent de 2 vorbitori nativi ai limbii engleze. Dezambiguizarea manuală a fost efectuată utilizând numerotarea sensurilor din WordNet 1.6.

Pentru a putea compara rezultatele produse de dezambiguizatorii umani (numiți în continuare adnotatori) cu cele produse de algoritmul nostru, datele au fost normalizate în felul următor: pentru fiecare adnotator și pentru algoritm fiecare din cele 33 de cuvinte a fost reprezentat printr-un vector binar de lungime $n(n-1)/2$, unde n este numărul de ocurențe ale cuvântului în tot corpusul. Pozițiile în vector reprezintă o asignare de tip “DA/NU” indicând dacă ocurența respectivă a fost clasificată la fel de către adnotatori, respectiv algoritm. Rezultatele acestui prim experiment sunt rezumate în tabelul din figura 8 indicând procentul de acord între clasificările propuse de algoritm și cele ale fiecărui adnotator, acordul dintre cei doi adnotatori și acordul dintre toți cei trei clasificatori.

Algoritm/Adnotator 1	66.7%
Algoritm /Adnotator 2	63.6%
Adnotator 1/Adnotator 2	76.3%
Algoritm /Adnotator 1/ Adnotator 2	53.4%

Figura 8: Concordanța între diferite clasificări

3.4. Cel de-al doilea experiment

Rezultatele primului experiment au aratat ca metoda discriminarii sensurilor folosind echivalentii de traducere este foarte competitiva, acuratetea procesului fiind comparabila (si uneori superioara) cu performantele obtinute de alti cercetatori ce au folosit ca referinta acelasi dictionar (Wordnet). Mai mult, diferentele de acord asupra clasificarii dintre cei 2 adnotatori pe de o parte si dintre fiecare adnotator si algoritmul pe de alta parte este de numai 10-13%, ceea ce din nou este foarte competitiv în raport cu scorurile obtinute în alte experimente.

Pentru a valida aceste rezultate empirice, în cea de a doua faza a experimentului a fost luat în considerare un numar dublu de substantive (76) dintre cele „dificile”, adica cu grad de ambiguitate mare, atât din clasa celor abstracte cât si a celor concrete (de exemplu, „thought”, „stuff”, „meaning”, „feeling” respectiv „hand”, „boot”, „glass”, „girl” etc.). Am ales acele substantive care au aparut cel puțin de 10 ori în corpus (pentru a elimina efectul de „insuficienta a datelor”) si în plus care au cel puțin 5 traduceri în cele 6 limbi tinta. Restrictia de 10 aparitii a aparut din pragul înalt de confidenta pe care l-am impus procesului de extragere a echivalentilor de traducere:

$$LL(T_T, T_S) = 2 * \sum_{j=1}^2 \sum_{i=1}^2 n_{ij} * \log \frac{n_{ij} * n_{**}}{n_{i*} * n_{*j}} \geq 18$$

În plus, pentru adnotarea manuala au fost cooptati înca doi vorbitori nativi ai limbii engleze, astfel încât fiecare dintre ocurențele cuvintelor selectate a fost etichetata, în mod independent, de 5 clasificatori: 4 adnotatori si algoritmul discutat aici. În tabela din figura 9 sunt rezumate datele si rezultatele de acord între cei 4 adnotatori:

Nr. de cuvânte	76
Nr. ocurențe	2399
Număr mediu de ocurențe-cuvânt	32
Nr. de sensuri găsite de adnotatorul 1	241
Nr. de sensuri găsite de adnotatorul 2	280
Nr. de sensuri găsite de adnotatorul 3	213
Nr. de sensuri găsite de adnotatorul 4	232
Nr. de sensuri găsite împreună de toți adnotatorii	345
Numărul mediu de sensuri pe cuvânt	4.53
Procent de acord între adnotatori	
Full agreement (4/4)	54.27
75% agreement (3/4)	28.13
50% agreement (2/4)	16.92
No agreement	0.66

Figura 9: Datele experimentului si acordul între 4 adnotatori umani independenti

Rezultatele produse de algoritmul de clasificare si clasificarile realizate de adnotatori prin asignarea sensurilor din Wordnet1.6 au fost de data aceasta normalizate în mod diferit, prin ignorarea etichetei puse de adnotatori si considerând doar clasele rezultând din aceasta etichetare. Pentru a clarifica acest aspect sa urmarim modul în care doi dintre adnotatori au dezambiguizat cele 7 ocurente ale cuvântului “youth”:

Ocurența nr.	1	2	3	4	5	6	7
Adnotatorul 1	3	1	6	3	6	3	1
Adnotatorul 2	2	1	4	2	6	2	1

Figura 10: Acordul de clasificare pentru cuvântul ” youth” între 2 adnotatori umani independenti

Acordul între cei doi adnotatori este doar de 43% (doar ocurențele 2, 5 si 7 au asignate sensuri consensuale); totusi, ambii adnotatori au clasificat ocurențele 1, 4, and 6 ca având acelasi sens, desi primul le-a etichetat cu sensul 3 din Wordnet, în timp ce al doilea le-a etichetat cu sensul 2. Daca însa ignoram eticheta clasificarea celor 3 ocurențe este consistenta, în sensul ca ambii adnotatori au decis ca ele au acelasi sens. Acordul de clasificare se dubleaza în acest caz⁸⁷, iar datele sunt mult mai usor de comparat cu rezultatele produse de algoritmul.

În acest al doilea experiment am luat în considerare determinarea momentului optim de oprire a clasificării aglomerative. În primul experiment, am folosit o distanta minima predefinita, pentru determinarea numarului de sensuri între care se realizeaza discriminarea. Aceasta solutie nu tinea însa cont de proprietatile individuale ale cuvintelor (numarul maxim de sensuri, prescrise de Wordnet, frecventa de aparitie a cuvântului, numarul mediu de traduceri pe care le-a primit cuvântul în corpus). Noul algoritmul de clasificare a fost modificat astfel încât sa-si calculeze un numar optim de clase⁸⁸, optimalitatea fiind judecata în raport cu numarul mediu de clase identificate de adnotatori. Drept criteriu de oprire am folosit distanta minima dintre clasele existente la fiecare pas de aglomerare. La un pas de aglomerare, clasele cu cea mai mica distanta relativa sunt reunite într-o clasa mai mare. Procesul începe cu fiecare ocurența într-o clasa distincta si se opreste

⁸⁷ Singurul dezacord ramas consta în faptul ca Adnotatorul 1 considera ocurențele 3 si 5 ca având acelasi sens, în timp ce Adnotatorul 2 atribuie un sens diferit ocurenței 3—în fapt, realizând o discriminare mai fina între sensurile celor doua ocurențe.

⁸⁸ În principiu, limita superioara a numarului de sensuri pe care îl poate avea un cuvânt englezesc într-un text este data de numarul de sensuri listate în Wordnet; dupa cum era de asteptat însa nu exista în corpusul nostru nici un exemplu în care vreun cuvânt polisemantic sa fi aparut cu toate sensurile din WordNet.

când distanțele relative între clasele existente este „suficient” de mare. Distanța dintre două clase se calculează pe baza vectorilor caracteristici (centrozii) ai celor două clase (evident depinzând de cuvânt, de numărul de ocurențe și de numărul de sensuri ale cuvântului clasificat):

$$\text{dist}(v_1, v_2) = \sqrt{\sum_{i=1}^n (v_1(i) - v_2(i))^2}$$

Cele mai bune rezultate în discriminarea automată au fost obținute pe cale experimentală, impunând drept criteriu de oprire a algoritmului condiția:

$$\frac{\text{mindist}(k) - \text{mindist}(k+1)}{\text{mindist}(k+1)} < 0.12$$

în care $\text{mindist}(k)$ reprezintă distanța minimă între clasele existente la pasul k de aglomerare.

Pentru medierea opiniilor adnotatorilor am definit o adnotare de referință reprezentând clasificarea majoritară între cei 4. În cazul egalității de voturi, adnotatorul care a fost în cele mai multe cazuri de aceeași opinie cu majoritatea a impus clasa. Folosind această clasificare mediata și raportând-o la clasificarea produsă de algoritm pentru cele 76 de cuvinte, am analizat diferențele de clasificare, considerate ca fiind erori. Marea majoritate a erorilor de clasificare pentru cele 2399 de ocurențe au apărut în cazul cuvintelor pentru care distribuția sensurilor este foarte inegală; ca urmare am adăugat algoritmului o fază suplimentară de postprocesare, în care clasele cu un număr mult mai mic de ocurențe decât clasa cu cele mai multe ocurențe au fost încorporate în ultima. Raportul minim între numărul de ocurențe al celei mai mari clase și numărul de ocurențe din clasele potențial absorbabile în cea dintâi a fost ales empiric ca fiind 10^{89} . Motivatia acestei euristici, constă în constatarea făcută de mai mulți cercetători în domeniul lingvisticii corpusului (fapt sugerat chiar de Zipf cu peste 50 de ani în urmă) ca utilizarea frecvența a unui cuvânt într-un text omogen tinde să-i pastreze sensul.

Cu această nouă euristica încorporată, algoritmul de clasificare a atins cifra de 74,6% acord cu clasificarea mediata. În (Ide *et al.* 2002) sunt prezentate alte variante ale algoritmului care au condus prin evaluarea empirică la versiunea sa finală. Clasele produse de fiecare pereche de clasificatori (om sau mașină) au fost evaluate printr-un algoritm ce calculează alinierea claselor astfel încât intersecția lor să fie maximală. Diferențele dintre două clase astfel alinate au fost considerate dezacorduri de clasificare. Scorul de acord a fost calculat ca fiind raportul dintre suma numărului de ocurențe comune pentru fiecare clasă aliniată și numărul total al ocurențelor cuvântului respectiv. În tabela din figura 10 este exemplificat modul de calcul al acordului dintre clasificarea produsă de algoritm și clasificarea mediata a adnotatorilor pentru cuvântul *movement*. Acesta a apărut în text de 40 de ori. Atât algoritmul cât și cei patru adnotatori au identificat 4 sensuri distincte în care

⁸⁹ 10 este o valoare precaută; experimente viitoare, de mai mare anvergură ar putea furniza argumente pentru coborârea acestui prag.

acest cuvânt a fost utilizat. Așa cum se vede din figura 10, cea mai numeroasă clasă (clasa 1) conține în clasificarea mediată 28 dintre cele 40 de apariții ale cuvântului *movement*, în timp ce clasa corespunzătoare creată de algoritm conține doar 25 de ocurențe. Dintre acestea, 24 sunt comune cu cele din clasa 1 a clasificării mediate. În conformitate cu definiția anterioară a scorului de acord, pentru acest exemplu algoritmul a produs rezultatul corect în 85% din cazuri.

CLASA	1	2	3	4	Σ
Clasificare mediată	28	6	3	3	40
Clasificare algoritmică	25	7	6	2	40
Intersecție	24	6	3	1	34
Precizie	85%				

Figura 11: Clasificarea mediată și cea produsă de algoritm pentru cuvântul *movement*

3.5. Rezultate

Rezultatele obținute cu ultima variantă a clasificatorului în cel de-al doilea experiment sunt sintetizate în tabelul din figura 12. Tabelul indică procentul de acord între diverse clasificări: 1, 2, 3, 4, reprezintă clasificările realizate de adnotatorii umani, M reprezintă clasificarea mediată a clasificatorilor umani, A reprezintă clasificarea produsă de algoritm, iar B este referința de bază (baseline) care presupune toate ocurențele unui cuvânt ca având același sens.

	1	2	3	4	M	A
B	71.1	65.1	76.3	74.1	75.5	81.5
1		78.1	75.6	83.1	88.6	74.4
2			71.3	75.9	82.5	66.9
3				77.3	82.1	77.1
4					90.4	75.9
M						77.3

Figura 12: Acorduri între diverse clasificări

Tabela arată că acordul între adnotatorii umani comparat cu cel dintre algoritm și adnotatorii umani (cu excepția unuia dintre ei (4), pe care îl suspectăm că a văzut clasificările celorlalți trei și în consecință și-a revizuit unele decizii) nu diferă substanțial. Acest lucru demonstrează (cel puțin în raport cu datele experimentului nostru) că dezambiguirea automată este comparabilă ca acuratețe cu cea efectuată de adnotatori

umani. Diferența fundamentală constă în faptul că programul a terminat în circa 2 minute clasificarea pentru care adnotatorilor le-au trebuit între 4 și 5 săptămâni.

Experimentul descris a evaluat dezambiguizarea automată a cuvintelor englezești pornind de la traducerea lor în celelalte 6 limbi. Aceasta direcționare a fost impusă doar de disponibilitatea pentru limba engleză a textului dezambiguizat de experți umani (vorbitori nativi ai limbii engleze). Întrucât algoritmul de clasificare nu depinde în nici un mod de limba pentru care se realizează dezambiguizarea (limba țintă) și nici de limbile martor în raport cu care se face acest proces, rezultă că exact același procedeu descris până aici poate fi folosit pentru dezambiguizarea cuvintelor românești folosind echivalenții lor de traducere în engleză, bulgară, cehă, estoniană, maghiară și slovenă, ori pentru dezambiguizarea cuvintelor bulgarești pe baza echivalenților lor de traducere în celelalte 6 limbi. Întrucât sensul este (în principiu) un invariant al traducerii, nu pare a se justifica și pentru celelalte limbi efortul de adnotare umană făcut pentru limba engleză. Este rațional să presupunem că rezultate similare (raporturi relative) s-ar obține indiferent de limba țintă și de limbile martor.

Să mai menționăm și faptul că există o anumită corelație (factorul Spearman - 0.51) între numărul de sensuri în Wordnet ale unui cuvânt și nivelul de acord între diferitele clasificări ale ocurențelor sale. Cele mai scăzute scoruri de acord au fost obținute pentru "line" (29 sensuri), "step" (10), position (15), "place" (17) și "corner" (11). Acorduri perfecte s-au obținut pentru majoritatea cuvintelor cu mai puțin de 5 sensuri, ca de exemplu "hair" (5), "morning" (4), "sister" (4), "tree" (2), and "waist" (2) care toate au fost considerate, atât de adnotatori cât și de algoritm, a fi fost folosite cu un singur sens în tot textul. Pe de altă parte, gradul de acord pentru câteva cuvinte cu mai puțin de 5 sensuri ("rubbish" (2), "rhyme" (2), "destruction" (3) și "belief" (3)) a fost semnificativ mai mic decât media pentru toate perechile de clasificări (adnotator-adnotator, adnotator-algoritm). Concluzia a fost că pentru unele cuvinte, distincțiile de sens sunt atât de fine în Wordnet, încât chiar vorbitorii nativi (și cu atât mai mult algoritmul de clasificare) nu pot face diferențieri sistematice de sens ale diferitelor ocurențe ale acestor cuvinte. O astfel de hiperdiferențiere a sensurilor este în imensa majoritate a cazurilor irelevantă pentru aplicațiile de prelucrare a limbajului natural.

4. Concluzii

Rezultatele experimentelor noastre arată că acuratețea discriminării sensurilor pe baza echivalenților de traducere extrasi din corpusuri paralele este comparabilă cu cea produsă de adnotatori umani. Întrucât abordarea noastră este complet automatizată ea poate fi folosită la crearea de volume mari de texte, având discriminate sensurile cuvintelor polisemantice. Utilizarea experților umani este prohibitivă sub aspectul costului și al timpului de realizare a unei asemenea sarcini, iar procentajul suplimentar de acuratețe, presupus de activitatea umană, este prea mic pentru a justifica procedurile manuale.

Metoda pe care am descris-o în aceasta lucrare nu etichetează clasele de ocurențe ale unui cuvânt cu un număr de sens ales dintr-un inventar prescris de sensuri iar majoritatea aplicațiilor de prelucrare a limbajului natural (de pildă clasificarea textelor, regăsirea informațiilor, rezumarea automată etc.) nici nu au nevoie de această informație suplimentară; pentru aceste tipuri de aplicații este suficient a decide ca două sau mai multe ocurențe ale unui cuvânt sunt folosite în același sens sau nu. O etichetare convențională a sensurilor identificate pentru un anumit cuvânt ar putea să se bazeze pe frecvența sensurilor respective (sensul 1 corespunzând clasei cu cele mai multe ocurențe). Evident o astfel de etichetare depinde de registrul lingvistic al textului pe baza căruia se identifică sensurile distincte.

O direcție foarte promițătoare (Tufis, 2002b), (Tufis&Cristea, 2002a,b) o constituie utilizarea metodologiei prezentată aici în construcția și validarea ontologiilor multilingve de tip EuroWordNet. Folosind echivalenții de traducere și clasificarea ocurențelor echivalente din punct de vedere al sensului se poate verifica dacă proiecția interlinguală a două sau mai multe dicționare semantice este corectă. Aceasta presupune ca sensurile cuvintelor extrase ca echivalenți de traducere ai cuvintelor englezești dezambiguizate să fie puse în corespondență cu același concept interlingual aparținând Indexului Interlingual (ILI - vezi Tufis&Cristea, 2002b – în acest volum). În cazul contrar (echivalenții de traducere sunt puși în corespondență cu concepte interlinguale diferite) este vorba de o eroare propriu-zisă de proiecție conceptuală într-unul sau mai multe dintre dicționarele semantice aliniate ori conceptele interlinguale sunt atât de apropiate semantic încât se poate propune unificarea lor într-un concept mai general cu lexicalizare în mai multe limbi. Aceasta este esența conceptului de „soft-clustering” definit în comunitatea EuroWordNet. Fata de identificarea prin metode statistice a conceptelor interlinguale foarte apropiate semantic, analiza prin metoda echivalenților de traducere și a discriminării sensurilor a proiecțiilor sensurilor făcute de lexicografi profesioniști peste o multitudine de sensuri conceptualizate în ILI este mult mai robustă. Experimentele preliminare discutate în (Tufis, 2002b) au arătat că în diferite limbi pentru care se realizează o ontologie lexicală multilingvă (bulgară, cehă, greacă, română, sârbă, turcă) există dificultăți identice de proiecție conceptuală a sensurilor unor cuvinte din limbile considerate. Faptul că aceleași concepte interlinguale creează același tip de dificultate în proiecția sensurilor unor cuvinte aparținând unor limbi foarte diferite indică cu claritate că acele concepte trebuie generalizate.

Un alt aspect care merita subliniat este că metodologia prezentată aici, corelată cu existența a tot mai multor dicționare semantice de tip Wordnet, ce aderă la principiul EuroWordNet de aliniere la Indexul Interlingual, va permite dezvoltarea de corpusuri adnotate semantic (de tipul SemCor) pentru orice limbă. Tranzitivitatea relațiilor de tip „EQ-SYN” folosite în proiecția sinseturilor unui wordnet monolingv peste ILI, corelată cu echivalența de traducere (relație tot între sensuri) extrasă dintr-un corpus paralel, în care textul dintr-una din limbi este adnotat semantic, permite importul adnotărilor în toate celelalte limbi. Deoarece limba din care se importa adnotarea semantică nu este relevantă pentru această procedură, rezultă că eforturile depuse de-a lungul timpului în crearea celor

câteva corpusuri cu adnotare semantică pentru limbile „mari” pot fi valorificate pentru orice altă limbă în care există (sau se creează) traduceri ale textelor din corpusurile adnotate. Mai mult, se poate imagina crearea unui consorțiu multilingv care să aleagă un corpus paralel în cât mai multe limbi cu scopul de a-l adnota semantic. Prin adnotarea independentă, în fiecare limbă, a unor porțiuni distincte din corpusul paralel, folosind o metodologie de genul celei prezentate în această lucrare (și desigur având un dicționar semantic multilingv de tip EuroWordNet) adnotările secțiunilor monolingve vor putea fi importate în secțiunile corespunzătoare ale tuturor celorlalte texte monolingve, în final putându-se obține adnotarea semantică, consistentă, a întregului text din fiecare limbă a corpusului paralel.

Mulumiri

Rezultatele prezentate în această lucrare sunt rodul mai multor proiecte internaționale de cercetare desfășurate la Institutul de Inteligență Artificială, alături de colegii Ana Maria Barbu, Eduard Barbu, Radu Ion, Catalin Mititelu, Octavian Popescu. De asemenea, colaborarea cu Nancy Ide de la Universitatea Vassar din Poughkeepsie, SUA, și cu Tomaz Erjavec de la Institutul „Jozef Stefan” din Ljubljana, Slovenia, parteneri în proiectele amintite, a fost și este extrem de productivă. Tuturor le aduc aici cuvenitele mulumiri.

Referințe bibliografice

- Ahrenberg, L., M. Andersson, M. Merkel. (2000). "A knowledge-lite approach to word alignment", in (Véronis, 2000: 97-116)
- Brants, T. (2000) "TnT – A Statistical Part-of-Speech Tagger", in *Proceedings of the Sixth Applied Natural Language Processing Conference, ANLP-2000*, April 29 – May 3, 2000, Seattle, WA
- Brew, C., McKelvie, D. (1996). "Word-pair extraction for lexicography", <http://www.ltg.ed.ac.uk/~chrisbr/papers/nemplap96>
- Brown, P., Della Pietra, S. A., Della Pietra, V. J., Mercer, R. L. (1993). "The mathematics of statistical machine translation: parameter estimation" in *Computational Linguistics* 19(2): 263-311.
- Carletta, J. (1996). Assessing Agreement on Classification Tasks: The Kappa Statistic. *Computational Linguistics*, 22/2, 249:254.
- Cristea, D., Dima, G. E. (2001). „An Integrating Framework for Anaphora Resolution”, *Journal on Information Science and Technology, Romanian Academy Publishing House*, Bucharest, vol. 4, no. 3, 273:292.

- Dagan, I., Itai, A. (1994). Word sense disambiguation using a second language monolingual corpus. *Computational Linguistics*, 20/4, 563:596.
- Dagan, I., Itai, A., Schwall, U. (1991). Two languages are more informative than one. *Proceedings of the 29th Annual Meeting of the ACL*, 18-21 Berkeley, California, 130:137.
- Dunning, T. (1993), "Accurate Methods for the Statistics of Surprise and Coincidence" in *Computational Linguistics* 19(1), 61:74.
- Dyvik, H. (1998). Translations as Semantic Mirrors. *Proceedings of Workshop Multilinguality in the Lexicon II, ECAI 98*, Brighton, UK, 24:44.
- Erjavec T., Ide, N. (1998). "The Multext-East corpus". In *Proceedings of First International Conference on Language Resources and Evaluation*, Granada, Spain, 971:974
- Erjavec. T. (2002) "An Experiment in Automatic Bi-Lingual Lexicon Construction from a Parallel Corpus", *Proceedings of the 7th TELRI International Seminar on Corpus Linguistics*, Dubrovnik, Croatia (forthcoming).
- Fellbaum C. (1998) *Wordnet: An Electronic Lexical Database*, MIT Press, 423p.
- Gale, W.A., K.W. Church, (1991). "Identifying word correspondences in parallel texts". In *Fourth DARPA Workshop on Speech and Natural Language*, 152:157
- Gale, W.A., K.W. Church, (1993). "A Program for Aligning Sentences in Bilingual Corpora". In *Computational Linguistics*, 19(1), 75:102
- Gale, W. A., Church, K. W., Yarowsky, D. (1993). A method for disambiguating word senses in a large corpus. *Computers and the Humanities*, 26, 415:439.
- Hinrics, H., Trushkina, J. (2002): "Forging Agreement: Morphological Disambiguation of Noun Phrases", *Proceedings of the Workshop on Treebanks and Linguistic Theories 2002*, Sozopol, Bulgaria, 1:18.
- Ide, N. (1999). Cross-lingual sense determination: Can it work? *Computers and the Humanities*, 34/1-2, 223:234.
- Ide, N., Erjavec, T., and Tufis, D. (2001). Automatic sense tagging using parallel corpora. *Proceedings of the Sixth Natural Language Processing Pacific Rim Symposium*, Tokyo, 83:89.
- Ide, N., Erjavec, T., Tufis, D. (2002): "Sense Discrimination with Parallel Corpora" in *Proceedings of the SIGLEX Workshop on Word Sense Disambiguation: Recent Successes and Future Directions*. ACL2002, July Philadelphia, 56:60
- Hiemstra, D. (1997). "Deriving a bilingual lexicon for cross language information retrieval". In *Proceedings of Gronics*, 21:26
- Kay, M., Röscheisen, M. (1993). "Text-Translation Alignment". In *Computational Linguistics*, 19/1, 121:142

- Kupiec, J. (1993). "An algorithm for finding noun phrase correspondences in bilingual corpora". In *Proceedings of the 31st Annual Meeting of the Association of Computational Linguistics*, 17:22
- Miller, G. A., Beckwith, R. T. Fellbaum, C. D., Gross, D. and Miller, K. J. (1990). WordNet: An on-line lexical database. *International Journal of Lexicography*, 3/4, 235:244.
- Melamed, D. (2001). "*Empirical Methods for Exploiting Parallel Texts*", MIT Press, 373p.
- Resnik, P. and Yarowsky, D. (2000). Distinguishing systems and distinguishing senses: New evaluation methods for word sense disambiguation. *Journal of Natural Language Engineering*, 5(2), 113:133.
- Resnik, P. and Yarowsky, D. (1997). A perspective on word sense disambiguation methods and their evaluation. *ACL-SIGLEX Workshop Tagging Text with Lexical Semantics: Why, What, and How?* Washington, D.C., 79:86.
- Smadja, F. (1993). "Retrieving Collocations from Text: Xtract". In *Computational Linguistics*, 19/1, 142:177
- Smadja, F., K.R. McKeown, and V. Hatzivassiloglou (1996). "Translating collocations for bilingual lexicons: A statistical approach". *Computational Linguistics*, 22/1, 1:38
- Stolcke, A. (1996). Cluster 2.9. <http://www.icsi.berkeley.edu/ftp/global/pub/ai/stolcke/software/cluster-2.9.tar.Z>.
- Tiedemann, J. (1998). "Extraction of Translation Equivalents from Parallel Corpora", In *Proceedings of the 11th Nordic Conference on Computational Linguistics*, Center for Sprogteknologi, Copenhagen, 1998, <http://stp.ling.uu.se/~joerg/>
- Tufis, D., Barbu, A.M., Patrascu, V., Rotariu, G., Popescu, C. (1997) "Corpora and Corpus-Based Morpho-Lexical Processing" in D. Tufis, P. Andersen (eds.) "Recent Advances in Romanian Language Technology", Editura Academiei, 35:56
- Tufis, D., Mason, O., (1998). "Tagging Romanian Texts: a Case Study for QTAG, a Language Independent Probabilistic Tagger". In *Proceedings of the First International Conference on Language Resources & Evaluation (LREC)*, Granada, Spain, 589:596.
- Tufis, D. (1999). "Tiered Tagging and Combined Classifiers" In F. Jelinek, E. Nöth (eds) *Text, Speech and Dialogue, Lecture Notes in Artificial Intelligence 1692*, Springer, 29:33
- Tufis, D., Ide, N. Erjavec, T. (1998). "Standardized Specifications, Development and Assessment of Large Morpho-Lexical Resources for Six Central and Eastern European Languages". In *Proceedings of First International Conference on Language Resources and Evaluation (LREC)*, Granada, Spain, 233:240

-
- Tufis, D. (2000) "Using a Large Set of Eagles-compliant Morpho-Syntactic Descriptors as a Tagset for Probabilistic Tagging. In *Proceedings of the Second International Conference on Language Resources and Evaluation (LREC)*, Athens, Greece, 1105:1112
- Tufis, D., Dienes, P., Oravecz, C., Váradi T. (2000). "Principled Hidden Tagset Design for Tiered Tagging of Hungarian" *Proceedings of the Second International Conference on Language Resources and Evaluation (LREC)*, Athens, Greece, 1421:1426
- Tufis, D., Barbu, A.M., (2001a). "Extracting multilingual lexicons from parallel corpora" in *Proceedings of the ACH/ALLC 2001*, New York University, 42:46
- Tufis, D., Barbu, A.M. (2001b). "Automatic Learning of Translation Equivalents" in "Romanian Journal on Information Science and Technology", Romanian Academy, vol.4, no. 3-4, 325:351.
- Tufis, D., Barbu, A.M., (2002). „Revealing translators knowledge: statistical methods in constructing practical translation lexicons for language and speech processing”, in *International Journal of Speech Technology*. Kluwer Academic Publishers, no.5, 199:209.
- Tufis, D.(2001a). "Building an ontology from a large Romanian dictionary of synonyms by importing Wordnet relations", RACAI Research report, June, 68 pp.
- Tufis, D. (2001b). "Partial translations recovery in a 1:1 word-alignment approach", RACAI Research report, June, 32pp.
- Tufis, D. (2002a). "A cheap and fast way to build useful translation lexicons" in *Proceedings of the 19th International Conference on Computational Linguistics, COLING2002*, Taipei, China, 246:251
- Tufis, D. (2002b). "Interlingual alignment of parallel semantic lexicons by means of automatically extracted translation equivalents", *Proceedings of the 7th TELRI International Seminar on Corpus Linguistics*, Dubrovnik, Croatia (forthcoming)
- Tufis, D., Cristea, D. (2002a). "Methodological issues in building the Romanian Wordnet and consistency checks in Balkanet", In *Proceedings of LREC2002 Workshop on Wordnet Structures and Standardisation*, Las Palmas, Spain, 35-41
- Tufis, D., Cristea, D. (2002b). "RO-BALKANET - ontologie lexicalizata, în context multilingv, pentru limba româna", în acest volum.
- Varadi, T.(2002) The Hungarian National Corpus, *Proceedings of LREC2002*, Las Palmas, Spain, 385:389
- Véronis, J. (ed), (2000). *Parallel Text Processing*. Text, Speech and Language Technology Series, Kluwer Academic Publishers Vol. 13, 2000

- Yarowsky, D., Florian. R. (1999). Taking the load off the conference chairs: towards a digital paper-routing assistant. *Proceedings of the Joint SIGDAT Conference on Empirical Methods in NLP and Very Large Corpora*, 220:230.
- Zipf, G.K., (1936). "The Psycho-biology of Language: an Introduction to Dynamic Philology". Routledge, London, UK

Referentialitate si cursivitate în relatie cu structura de discurs

Dan CRISTEA

Universitatea "A.I. Cuza" Iasi, Facultatea de Informatica

Academia Româna, Institutul de Informatica Teoretica - Filiala Iasi

1. Introducere

În ultimii 25 de ani s-a studiat enorm pentru a se înțelege ce anume face dintr-un text (considerat o secvență de propoziții sintactice corecte) să fie un discurs, adică de ce un discurs e coerent și ce elemente îi atribuie coeziune. Dintre teoriile computationale ale discursului, trei au avut o influență covârșitoare asupra dezvoltărilor ultimilor ani din acest domeniu: teoria structurilor retorice, teoria starilor atentionale și teoria centrelor.

Dezvoltată inițial din perspectiva generării textelor, teoria structurilor retorice (*rhetorical structure theory*, de aici încolo RST), a fost elaborată de Mann și Thompson între 1986 și 1987 ca o teorie a organizării textelor [Mann, Thompson, 1988, Hovy, 1988, Scott, de Souza, 1990]. Ea caracterizează structura de discurs în termeni de relații ce leagă părți componente ale textului. Unitatea elementară de discurs în RST este, de regulă, o propoziție ce, la nivel semantic, formulează o predicatie. O structură de discurs este descrisă de o **schema**. Ea grupează o secvență de unități, sau de unități și scheme, sau o secvență de scheme. Dintr-un anumit punct de vedere o schema poate fi asemuită cu o regulă a unei gramatici, ea relevând structura de constituenți a unui compus. O schema constă dintr-o **relație** (27 în RST) care leagă două sau mai multe întinderi de text, fiecare dintre ele având, la rândul lor, o structură (constituenții schemei). Un discurs este fie o unitate, care este o întindere de text elementară, fără structură, fie o schema (un text mai lung decât o singură unitate și care manifestă o structură). Relațiile pot fi de două tipuri: **hipotactice** – dacă argumentele sunt constituenți neegali ca importanță și **paratactice** (sau echinucleare) – dacă constituenții pe care-i agregă sunt egali ca importanță. Între constituenții uniți de relațiile hipotactice există întotdeauna unul singur mai important, numit **nucleu**, ceilalți fiind numiți **sateliți**. La relațiile paratactice, prin convenție, se consideră ca toți constituenții sunt nucleari. Satelitul este, în general, mai susceptibil de a fi schimbat sau eliminat complet decât nucleul, fără ca, prin aceasta, înțelesul discursului să se modifice. Dimpotrivă înlocuirea sau stergerea nucleului este o opțiune mult mai drastică, care poate duce la denaturări ale înțelesului. Relațiile hipotactice sunt în general cele intenționale, în care o întindere de text comunică un scop și o altă exprimă un subscop ce

completeaza, dezvolta etc. scopul principal. Pe de alta parte, relatiile paratactice sunt, în general, de natura informatională, simetrice, neputându-se stabili dacă și care componentă predomină.

În RST accentul este pus pe performanța retorică: prin ce mijloace un scriitor (sau vorbitor) reușește să convingă un cititor (ascultător) de intențiile pe care le are de comunicat. Ca produs secundar al liniei principale de investigație în RST, multe eforturi de cercetare care au succedat elaborarea teoriei s-au concentrat asupra îmbunătățirii setului de relații propus inițial de teorie. Într-adevăr pare extrem de convenabil, inclusiv din punct de vedere computațional, să vedem discursul reprezentat ca un arbore, în care nodurile terminale să reconstituie, în secvența lor, textul. Cu toate acestea RST nu aduce nici o lumină în privința vreunei legături care ar exista între structura și referențialitate. RST este deci o teorie asupra structurii globale a discursului.

Teoria stărilor atenționale (*attentional state theory*, AST) [Grosz, Sidner, 1986] reprezintă o dezvoltare a liniei de cercetare în discurs dominată de Barbara Grosz și Candace Sidner asupra manierei în care focarul ori centrul de discurs (*focus* în engleză) se modifică pe parcursul derulării textului și a recunoașterii intențiilor comunicate de discurs [Grosz, 1981, Sidner, 1983]. Grosz și Sidner nu cred că varietatea atât de mare a intențiilor ce pot fi comunicate de un discurs poate fi condensată într-un număr fix de sabloane retorice exprimate ca relații, ca în RST sau tentative similare acestora. Teoria se dorește a fi un model formal, care se distanțează de detaliile ce ar putea fi asociate participanților la discurs. Realizând proiecții corespunzătoare utilizatorului de limbaj, însoțite de detalii specifice, ea s-ar putea regăsi atât în construcția unui sistem automat cât și într-o teorie psihologică, ambele consumatoare de limbaj natural. Deși recunoaște însemnătatea mesajului transmis de un discurs, teoria nu adresează problema înțeleșului discursului și a manierei în care acesta poate fi dedus din elementele constitutive ale textului. Ea este, primordial, o teorie a structurii discursului, prin aceasta plasându-se la baza oricărei tentative de a aborda problema construirii sensului.

Conform lui Grosz și Sidner intențiile joacă rolul principal în explicarea structurii discursului, în timp ce dinamica atenției joacă rolul principal în explicarea interpretării discursului. Structura discursului are trei componente distincte, dar strâns corelate:

- o **structură lingvistică**, care face ca una sau mai multe propoziții, exprimări (*utterance* în engleză) să fie agregate într-un segment de discurs iar limitele dintre segmente să fie indicate de expresii lingvistice, intonație, schimbări ale timpului și aspectelor verbale. Segmentul de discurs are însă o definiție recursivă: un segment poate îngloba alte segmente, acestea pe altele s.a.m.d.;
- o **structură intențională**, care face să vedem discursul ca având un scop global (scopul discursului – SD), care este scopul fundamental al vorbitorului/scriitorului la emiterea discursului și fiecare segment al său are un scop al segmentului (scopul segmentului de discurs – SSD) care este un subscop al scopului segmentului din care face parte. Dintr-un punct de vedere intuitiv, SSD specifică cum contribuie respectivul subsegment la

realizarea scopului segmentului din care face el parte. Teoria admite ca nu exista o lista finita de scopuri ale discursului, care sa faca posibila o comparatie cu lista categoriilor gramaticale, de exemplu. Conform teoriei, doua relatii structurale sunt suficiente pentru a compune structura discursului: **relatia de dominare** (daca SSD_1 domina SSD_2 atunci SSD_2 contribuie la SSD_1 , sau SSD_2 este intentionata sa satisfaca partial SSD_1) si **relatia de satisfacere-precedenta** (SSD_1 satisface-precede SSD_2 daca SSD_1 trebuie satisfacut înainte de SSD_2);

- o **stare attentionala**, prin care se asociaza fiecarui segment al discursului un spatiu al entitatilor aflate în centrul atentiei. Starea attentionala este o proprietate a discursului iar nu a participantilor la discurs. Ea reprezinta o trasatura dinamica a discursului, pastrând **obiecte, proprietati si relatii** ce sunt importante la fiecare moment al parcurgerii discursului. Starea attentionala e modelata printr-un **set de spatii ale centrelor atentiei**, în timp ce schimbarile ce pot avea loc în starea attentionala sunt restrictionate de un set de reguli de tranzitie care arata conditiile de adaugare si stergere a spatiilor. Colectia tuturor spatiilor centrelor de atentie ce sunt disponibile în fiecare moment al interpretarii unui discurs formeaza o structura a atentiei ce are dinamica unei **stive** si care ar fi capabila sa explice procesele implicate în procesarea discursului, inclusiv accesibilitatea referentiala: domeniul în care trebuie cautate entitatile de discurs referite în segmentul corespunzator starii attentionale aflate în vârful stivei este cel al starilor aflate în stiva.

Structura recursiva a segmentului de discurs din AST permite si aici acceptarea unei reprezentari arborescente, în cadrul careia cele doua relatii între segmente, de dominare si de satisfacere-precedenta, nu sunt altceva decât relatiile topologice normale pe orice structura de arbore: cea dintre parinte si orice fiu al sau si, respectiv, cea de ordine dintre frati. AST se constituie într-o teorie globala asupra structurii si a coeziunii discursului.

Cercetatori precum Moser si Moore [1996] sau Marcu [1999] pun în evidenta similaritati semnificative între AST si RST, inclusiv în ceea ce priveste maniera de reprezentare prin arbori a structurii de discurs, ceea ce permite combinarea puterii de reprezentare, mai fine în RST, datorita proliferarii relatiilor, cu implicatiile pe care structura le poate avea asupra referentialitatii, puse în evidenta de AST. Utilizând structura de segmente si stiva, ca mecanism de prelucrare, AST propune o maniera de a rezolva accesibilitatea referintelor anaforice printr-o transparenta pe verticala, de sus în jos, de-a lungul starilor attentionale ce se afla la un moment dat în stiva. Reprezentarea prin segmente din AST are însa o slabiciune: modelul stiva nu poate reflecta relatia de dominare atunci când scopul dominat corespunde unui segment care apare în text înaintea celui care domina [Ide, Cristea, 2000]. Sa remarcam ca defectul este unul de granularitate pentru ca identificarea segmentului dominat ce precede pe cel dominator cu însusi segmentul dominator elimina problema. AST nu e, asadar, capabila sa reprezinte segmente având o

granularitate oricât de fina: coborând de la o granularitate grosiera la una fina, exista o limita dincolo de care ne putem astepta la grave contradictii.

Teoria centrelor (*centering*, CT) [Grosz *et al.*, 1995, Brennan *et al.*, 1987] furnizeaza explicatii convingatoare asupra contextelor ce permit utilizarea pronumelor pentru realizarea referintelor si asupra ce anume face un discurs sa fie coerent. CT nu se aplica însa dincolo de limitele unui segment (vazut în acceptiunea din AST). Avem de a face, asadar cu o teorie locala asupra coeziunii si coerenței. Desi nu este definita riguros în teorie, în toate exemplele autorilor unitatea elementara a structurii lingvistice este fraza (*utterance*, exprimare). Abordari ulterioare întrevad posibilitatea de a considera o segmentare mai fina, la nivel de propozitie (v. [Kameyama, 1998] de exemplu). Noi vom considera drept unitate a structurii de discurs acelasi tip de întindere lexicala ca si în cazul RST, adica acea întindere ce la nivel sintactic este o propozitie iar la nivel semantic – o predicatie. Fiecare unitate de discurs u_n ce intra în compozitia unui segment este caracterizata de o lista de **centre anticipatoare** (*forward-looking*) notata $C_f(u_n)$. Centrele listei $C_f(u_n)$ sunt entitati semantice ce corespund, la nivelul textului, expresiilor referentiale cuprinse în unitatea u_n . Spunem ca o expresie referentiala **realizeaza** un centru. Elementele acestei liste sunt ordonate pentru a reflecta importanta relativa în u_n . Criteriile de ordonare a elementelor listei C_f sunt, în forma originara a teoriei, de natura sintactica, desi alte abordari le diferentiaza în functie de limba (v. de exemplu [Walker *et al.*, 1994] pentru japoneza, [deEugenio, 1990, de Eugenio, 1998] pentru italiana, sau [Strube, Hahn, 1996] pentru germana). Pentru limba engleza, autorii CT dau urmatorul criteriu: subiect > obiect-direct > obiect-indirect > complemente > adjuncti. Elementele listei $C_f(u_n)$ sunt acele entitati despre care se vorbeste în unitatea u_n si deci despre care e cel mai probabil ca se va continua sa se vorbeasca si în unitatea urmatoare, u_{n+1} , daca aceasta apartine aceluiasi segment ca si u_n . Cel mai bine plasat element al listei $C_f(u_n)$ se numeste **centru principal** si se noteaza $C_p(u_n)$. Fiecarei unitati îi este asociat un unic **centru retroactiv** (*backward-looking*), notat $C_b(u_n)$. Prin conventie, centrul retroactiv al primei unitati a segmentului este considerat centrul principal, în timp ce, pentru toate celelalte unitati ale segmentului, el este cel mai bine plasat element al listei C_f a unitatii precedente care este de asemenea realizat si în unitatea curenta.

Teoria face o clasificare a tranzitiilor posibile între unitati consecutive, din punctul de vedere al invariantei ori nu a centrelor retroactive si al identificarii ori nu a lor cu centrele principale. Astfel, cu exceptia cazului în care între unitati succesive ale aceluiasi segment nu exista centre comune, urmatoarele patru tipuri de tranzitii sunt posibile:

CONTINUARE (*continuing*, CON): $C_b(u_{n+1}) = C_b(u_n)$ si $C_p(u_{n+1}) = C_p(u_n)$, corespunzând situatiei în care atât în u_n cât si în u_{n+1} se vorbeste despre aceeasi entitate si este de asteptat ca si în unitatea urmatoare sa se vorbeasca despre ea.

RETINERE (*retaining*, RET): $C_b(u_{n+1}) = C_b(u_n)$ dar $C_p(u_{n+1}) \neq C_p(u_n)$, a carui interpretare este ca, desi atât în u_n cât si în u_{n+1} se vorbeste despre

aceeasi entitate, este de asteptat ca în unitatea urmatoare sa se vorbeasca despre o alta.

SCHIMBARE LINA (*smooth-shifting*, SSH): $C_b(u_{n+1}) \neq C_b(u_n)$ dar $C_b(u_{n+1}) = C_p(u_{n+1})$, cu semnificatia ca în u_n si în u_{n+1} nu se vorbește despre aceeași entitate si este de asteptat ca în unitatea urmatoare sa se vorbeasca despre entitatea despre care s-a vorbit ultima oara.

SCHIMBARE ABRUPTA (*abrupt-shifting*, ASH): $C_b(u_{n+1}) \neq C_b(u_n)$ si $C_b(u_{n+1}) \neq C_p(u_{n+1})$, cu semnificatia ca în u_n si în u_{n+1} nu se vorbește despre aceeași entitate si este de asteptat ca în unitatea urmatoare sa se vorbeasca despre o alta entitate decât ultima mentionata.

Nucleul CT este concentrat în doua reguli, prima enunțând o constrângere asupra formei de realizare a centrelor prin pronume, iar cea de a doua formulând preferințe asupra secvențelor de tranziții ale centrelor. Regula a doua, cea care se refera la coerența, formulează presupunerea ca anumite secvențe produc o încărcare inferențială în ascultător mai mare decât altele:

Regula 2: Secvențele de continuari sunt preferabile secvențelor de retineri, care sunt preferabile secvențelor de schimbări line, iar acestea sunt preferabile secvențelor de schimbări bruste: CON > RET > SSH > ASH.

Daca ne abținem de a penaliza CT, ca teorie locală, asadar aplicabilă la întinderea unui segment, pe motivul fragilității noțiunii de segment, care are o definiție recursivă (un segment este constituit din alte segmente), slabiciune moștenită de la AST, atunci apare naturală tentativa de a largi aplicabilitatea CT la întregul discurs, într-o manieră recursivă, pe chiar această structură de segment, definită, ea însăși, recursiv. Teoria nervurilor propune o astfel de generalizare.

Teoria nervurilor (*veins theory*, VT) [Cristea *et al.*, 1998], preluând de la RST diferențierea dată de nuclearitate între argumentele relațiilor retorice dar ignorând, ca și în AST, numele acestora, releva o structură "ascunsă" în arborele de discurs, numită **nervura**. Fără a nega structura lingvistică a segmentelor de discurs, cât și pe cea intențională a relațiilor dintre scopurile comunicate de segmente și care, prin echivalarea de care am amintit ([Moser, Moore, 1996, Marcu, 1999]), poate fi recuperată din structura de arbore proprie analizelor RST, VT corectează defectul de accesibilitate al AST înlocuind modelul accesibilității în stivă cu accesibilitatea de-a lungul nervurilor arborelui de discurs și explicând naturalitatea unor referințe la distanță realizate prin mijloace de evocare foarte economice (pronume) [Fox, 1987]. Concluziile VT sunt, de asemenea, stabile la granularitate. În felul acesta VT se constituie într-o teorie globală a coeziunii discursului. VT generalizează totodată partea din CT relativă la încărcarea inferențială (regula a doua), extinzând concluziile ei la întregul discurs, prin această VT constituindu-se și într-o teorie globală a coerenței.

În secțiunea urmatoare sunt prezentate argumente lingvistice în favoarea teoriei. Secțiunea 3 prezintă definițiile teoriei, secțiunea 4 enunță conjectura VT relativă la

referentialitate, iar secțiunea 5 – conjectura VT referitoare la coerența. Secțiunea 6 descrie rezultate experimentale în sprijinul presupuzițiilor VT, secțiunea 7 prezintă o proprietate de granularitate, iar ultima secțiune este dedicată concluziilor și prezentării unor aplicații ale VT.

2. Intuițiile VT

Notiunea de nervură s-a născut sintetizând observațiile asupra modului în care se aliniaza referințele pe o reprezentare arborescentă a discursului. Considerând organizarea ierarhică dată de structura de arbore și principiul compozitionalității, care permite ca unități de discurs aflate la distanță să fie frați sub aceeași relație, aceste observații au fost următoarele (pentru simplificarea exprimării vom spune că "o unitate A referă o unitate B" și vom înțelege "o expresie referențială aparținând unei unități A referă o entitate de discurs introdusă de (sau referită dintr-o) unitate B"; de asemenea vom nota cu u_1, u_2, u_3 – unități de discurs iar cu R, R_1, R_2 – relații. Atunci când apar ca argumente ale unei relații, unitățile de discurs vor purta un indice ridicat " sau s , cu semnificația de nucleu sau satelit):

- un satelit sau un nucleu poate referi un frate nuclear aflat la stânga: în combinații $u_1^n R u_2^s$, sau $u_1^n R u_2^n$, u_2 poate referi u_1 ;

Ex. 1

1. *Ion a plecat de acasă fără umbrelă*
2. *desi dimineața r aflate la radio că va ploua.*

Subiectul vid (notat r) din unitatea 2, un satelit al unității 1, referă entitatea [Ion]⁹⁰ introdusă de expresia referențială *Ion* din prima unitate.

- un nucleu poate referi un satelit al sau aflat la stânga: în combinații $u_1^s R u_2^n$, u_2 poate referi u_1 . Astfel, în exemplul:

Ex. 2

1. *Ion i-a dat Mariei o floare.*
2. *Pentru că r s-a simțit frustrată,*
3. *sotia lui s-a separat.*

unitatea 2 este un satelit al unității 3. Pe cine desemnează pronumele vid (notat r) din 2, pe [Maria] sau pe [sotia lui Ion]? Într-o interpretare incrementală a textului, la sfârșitul recepționării celei de a doua unități avem tendința de a asocia, prea timpuriu, subiectul vid [Mariei]. După citirea unității 3 are loc însă o reconsiderare a legării $r \rightarrow$ [Maria] și o

⁹⁰ Vom nota prin [text] entitatea de discurs introdusă/referită de expresia referențială text.

identificare a expresiei referentiale *sotia lui* cu subiectul vid din 2, ambele indicând entitatea [**sotia lui Ion**].

- un satelit drept al unui nucleu u nu e accesibil dintr-un alt frate drept, nuclear sau satelit, al lui u : în combinații $(u_1^n R_1 u_2^s)^n R_2 u_3^n$ sau $(u_1^n R_1 u_2^s)^n R_2 u_3^s$, u_3 poate referi u_1 dar nu u_2 .

Ex. 3

1. *Ion i-a marturisit Mariei ca o iubeste.*
2. *El n-a fost niciodata casatorit*
3. *si a trait pâna la 40 de ani lângă mama sa.*
4. *Ea, dimpotriva, a fost maritata de doua ori.*

Secventa 2-3-4 ofera o explicatie la 1. Secventa 2-3 se afla într-o relatie de CONTRAST (o relatie paratactica) fata de 4, iar 3 aduce o completare la 2. Structura este deci urmatoarea: $u_1^n R_1 ((u_2^n R_2 u_3^s)^n R_3 u_4^n)^s$ în care R_3 este relatia CONTRAST. Pentru cei mai multi cititori, *ea* din unitatea 4 trebuie sa fie [**Maria**], iar nu [**mama lui Ion**], desi [**mama lui Ion**] este entitatea cea mai recent referita, din pozitia unitatii 4, cu care pronumele feminin se potriveste în numar si gen. Motivul preferarii Mariei în locul mamei este acela ca cititorul recunoaste unitatea 4 ca fiind într-o relatie de CONTRAST cu unitatea 2 (relatie pusa în evidenta prin *dimpotriva*), ceea ce face ca cele doua unitati sa fie percepute ca fiind adiacente. Apropierea lor nu este însa una liniara, ci ierarhica, pe structura. Unitatea 3 este închisa la referinta din unitatea 4.

- un nucleu blocheaza accesibilitatea dintr-un satelit drept spre un satelit stâng: în combinații $(u_1^s R_1 u_2^n)^n R_2 u_3^s$, u_3 poate referi u_2 dar nu u_1 .

Ex. 4

1. *Înca înainte cu un an de terminarea mandatului sau de presedinte al firmei*
2. *dl. W. Ross începuse masinatiile pentru falimentarea acesteia.*
- *3. *De altfel, circulau vorbe ca l-ar fi obtinut fraudulos.*

În acest exemplu 1 si 3 sunt sateliti ai lui 2 (1 este o circumstantiala a lui 2, în timp ce unitatea 3 da o explicatie la purtarea necinstita a lui Ross). Referinta $l=[\mathbf{mandatul de presedinte al firmei al lui Ross}]$ se deduce cu dificultate, ceea ce face ca întregul discurs sa fie defectuos. Dimpotriva, în urmatoarea varianta, discursul câstiga în cursivitate:

Ex. 5

1. *Dl. W. Ross începuse masinatiile pentru falimentarea firmei al carei presedinte era*
2. *înca înainte cu un an de terminarea mandatului sau.*

3. *De altfel, circulau vorbe ca l-ar fi obtinut fraudulos.*

În Ex. 5 unitatea 2 este un satelit al lui 1, iar 3 – un satelit al lui 2 (aici *de altfel* anunta o paranteza la informatia asupra mandatului de presedinte). Referinta $l=[\text{mandatul de presedinte al firmei al lui Ross}]$ poate fi recuperata acum fara dificultate.

Motivatia acceptarii Ex. 5 si rejectarii Ex. 4, consta nu în departarea liniara mai mare a anaforului de antecedent în Ex. 4 decât în Ex. 5, ci în faptul ca în Ex. 4, spre deosebire de Ex. 5, accesul anafor-antecedent se face dinspre un satelit catre un alt satelit, între ei interpunându-se un nucleu. Sa remarcam ca Ex. 4 este reparat daca se elimina aceasta referinta:

Ex. 6

1. *Înca înainte cu un an de terminarea mandatului sau de presedinte al firmei*
2. *dl. W. Ross începuse masinatiile pentru falimentarea acesteia.*
3. *De altfel, circulau vorbe ca el ar fi fraudulat alegerile.*

3. Definitiiile teoriei

Intuitia fundamentala care sta la baza dezvoltarilor unificatoare asupra structurii de discurs si accesibilitatii în VT este ca distinctia specifica RST dintre nuclee si sateliti constrânge plaja de referenti asupra carora pot fi rezolvati anaforii; cu alte cuvinte, distinctia nucleu-satelit, corelata cu o structura de discurs, induce pentru fiecare unitate de discurs un domeniu de accesibilitate referentiala imediata pentru anaforii pe care-i contine. Mai precis, pentru fiecare anafor x aparținând unei unitati de discurs u , VT avanseaza ipoteza ca x poate fi rezolvata cu usurinta examinând doar un subset al multimii entitatilor de discurs care preced u . Daca antecedentul lui x este plasat într-o unitate de discurs aflata în afara domeniului lui u atunci legatura anafor-antecedent este refacuta cu greutate, sau pentru realizarea ei e nevoie de mijloace referentiale tari, cum sunt, de exemplu, numele proprii.

Mai mult decât atât, aceeași corelatie nuclearitate-structura, aplicata întregului discurs, permite generalizarea CT dincolo de granitele unui segment, ceea ce face posibila aplicarea concluziilor CT asupra coerentei la întregul discurs.

VT se bazeaza, în mare masura, pe aceleasi elemente ale structurii de discurs ca si RST:

- unitatile de baza ale discursului sunt întinderi de text (în engleza – *text span*) ce nu se intersecteaza. Dupa cum am precizat mai sus, noi le vom asimila cu propozitii, la nivel semantic fiecare continând o predicatie (careia îi corespunde o reprezentare evenimentiala sau situationala);

- structura unui discurs este reprezentată ca un arbore. Spre deosebire de RST, dar fără a reduce generalitatea, în VT vom considera arborii de discurs ca fiind binari (fiecare nod are exact doi descendenți) (pentru argumentație, v. [Marcu, 2000] și [Cristea, Webber, 1997]);
- principiul secvențialității [Cristea, Webber, 1997]: secvența de noduri de pe frontiera terminală a arborelui corespunde secvenței de unități de discurs ce compune textul⁹¹;
- principiul compoziționalității [Marcu, 2000]: o relație ce se aplică între două întinderi de text se aplică, de asemenea, și între subîntinderile nucleare ale întinderilor aflate în relație;
- la fel ca în RST, nuclearitatea nodurilor este importantă, nodurile fiind clasificate în nuclee (cele mai importante) și sateliți (cele mai puțin importante);
- nodurile terminale ale arborelui reprezintă unitățile de discurs, în timp ce nodurile neterminale reprezintă relații retorice între întinderi adiacente de text. Spre deosebire de RST, în VT nu interesează numele relațiilor, ceea ce contează fiind topologia arborelui, nuclearitatea nodurilor și etichetarea nodurilor terminale;
- între fiii fiecărui nod intermediar al arborelui există cel puțin un nod nuclear. Nodul radacina, prin convenție, e considerat satelit.

În vizualizarea arborilor vom reprezenta nodurile neterminale prin dreptunghiuri fără nume, pe cele terminale – prin ovaluri etichetate, iar nodurile nucleare vor fi subliniate (v. Figura 1). În definițiile ce urmează vom folosi următoarele convenții de notare:

- $mark(\alpha)$ este o funcție care întoarce șirul α în care fiecare simbol este marcat (de exemplu, este poziționat între paranteze);
- $unmark(\alpha)$ este funcția inversă lui $mark()$, ce îndepărtează toate marcasele atasate simbolurilor din expresia α (ex. $unmark(mark(\alpha)) = \alpha$);
- $simpl(\alpha)$ este funcția care elimină toate simbolurile marcate din expresia argumentului α (ex. $simpl(mark(\alpha)) = r$, șirul vid și $simpl(\alpha \cdot mark(\beta) \cdot \gamma) = \alpha \cdot \gamma$);
- $seq(\alpha, \beta)$ este o funcție de secvențiere, care întoarce acea permutare a concatenării simbolurilor din α și β dată de citirea de la stânga la dreapta a nodurilor corespunzătoare simbolurilor din α și β pe frontiera terminală a

⁹¹ Unitățile de discurs întrerupte nuanțează acest principiu. Astfel într-un discurs precum următorul: *O dată,¹ când treceau unul pe lângă altul pe coridor,² ea îi aruncase o privire piezișă¹ care parcă-l străpunsese³ și pentru o clipă fusese cuprins de o groază oarbă.⁴* (G. Orwell, 1984), unitatea 1 este întreruptă de unitatea 2.

arborelui. Funcția menține marcajele asupra simbolurilor, dacă acestea există, $seq(r, \beta) = \beta$, și $seq(\alpha, seq(\beta)) = seq(seq(\alpha), \beta) = seq(\alpha, \beta)$;

- $H(n)$ și $V(n)$ reprezintă expresiile *head* și nervura (în engleză – *vein*) ale unui nod n ;
- $pref(u, \alpha)$ reține prefixul expresiei simbolice α până la simbolul u inclusiv, o etichetă de nod terminal.

Teoria nervurilor calculează două expresii pe care le atasează fiecărui nod al structurii.

3.1 Expresia *head* a unui nod al arborelui

Intenția expresiei *head* a unui nod al arborelui de discurs este de a pune în evidență secvența celor mai importante unități de discurs din întinderea de text acoperită de nod. Ea este o secvență de etichete de unități, după cum urmează:

Definiii

1. Expresia *head* a unui nod terminal este eticheta sa;
2. Expresia *head* a unui nod neterminal este dată de concatenarea, în ordinea apariției lor în arbore de la stânga la dreapta, a expresiilor *head* ale descendenților săi nucleari.

Definițiile expresiilor *head* sugerează un proces de calcul care se propagă de jos în sus în arborele de discurs. Cele mai importante unități de discurs sunt proiectate în sus până în primul nod satelit întâlnit.

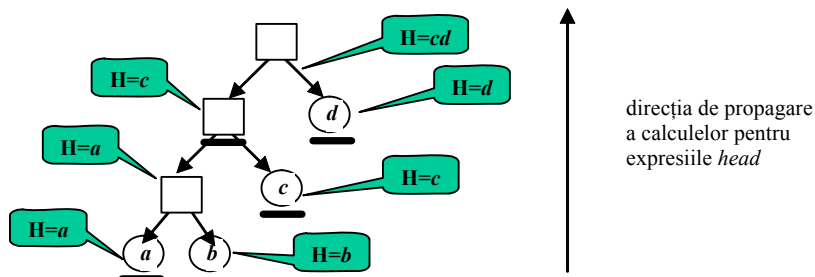


Figura 1: Calculul expresiilor *head*

3.2 Expresia nervurii unui nod al arborelui

Expresia **nervurii** unui nod intenționează să surprindă **secvența unităților de discurs care sunt semnificative pentru a sintetiza**⁹², **în contextul întregului text, întinderea de text** (în engleză – *text span*) **acoperita de nod**. Pentru orice nod al structurii, expresia nervurii este formată din cele mai importante unități din întinderea acoperită de nod, împreună, eventual, cu alte unități din afara acestei întinderi.

Definițiile care urmează, datorită recursivității lor, vor face posibilă considerarea contextului dat de totalitatea textului din exprimarea "a înțelege, în contextul întregului text, întinderea s" marginit la întinderea de text acoperită de nodul părinte al celui corespunzător întinderii s. Cu alte cuvinte, la fiecare nivel al structurii, cu excepția rădăcinii, adică acolo unde există două noduri fii sub un nod părinte, cu întinderile celor două noduri fii însumând întinderea nodului părinte, expresia nervurii a părintelui conține deja informația care permite înțelegerea/ rezumarea întinderii acoperite de el în contextul global. Coborârea pentru înțelegerea/ rezumarea subîntinderii acoperite de nodul curent al definiției (unul dintre cele două noduri fii) înseamnă adăugarea și/ sau ștergerea unei secvențe noi/ subsecvențe la/ din secvența de etichete contribuită de nervura părintelui, în funcție de polaritatea și poziția specifică a întinderii corespunzătoare nodului fiu curent în întinderea nodului părinte. În continuare, întinderea întregului text, o constantă pentru orice subîntindere, va fi numită **contextul total**. În figurile 2-6, nodurile curente – cele vizate de definițiile curente de nervură – apar în gri. Ele sunt notate simultan cu un dreptunghi și un oval pentru a sugera că pot fi atât noduri interioare (neterminale), cât și noduri terminale.

Definiții

1. Expresia nervurii rădăcinii este egală cu expresia sa *head*.

Expresia nervurii nodului rădăcină, conform intenției generale a nervurii unui nod, ar trebui să fie formată din cele mai semnificative unități de discurs necesare înțelegerii/ rezumării întinderii acoperite de nod (în cazul de față – întregul text) în contextul total. Cum contextul este aici egal cu textul în totalitatea lui, el poate fi lăsat la o parte în descriere, ceea ce ne lasă cu definiția expresiei *head* a nodului rădăcină.

2. Pentru fiecare nod nuclear, al cărui părinte are nervură v:

⁹² Prin sinteză, sau rezumatul, unei întinderi de text se înțelege un text mai scurt care redă ideea principală a textului supus sintezei. Indiferent dacă este realizat prin parafrază sau prin punerea cap la cap a unor subsecvențe ale întinderii originale [Mani, 2001], orice rezumat trebuie să fie comprehensibil, adică trebuie să poată fi înțeles prin el însuși (printre altele, de exemplu, rezumatul trebuie să continue toate elementele care să permită rezolvarea anafurilor). Adesea însă, atunci când întinderea este decupată dintr-un context mai larg, pentru ca rezumatul să fie comprehensibil, el trebuie să continue și elemente din afara întinderii și care aparțin contextului. Avem de a face, în acest caz, cu o sinteză a unei întinderi de text **în contextul** unei întinderi mai vaste. Sa mai observăm că, în multe privințe, "a sintetiza" e analog cu "a înțelege", pentru că ceea ce ne rămâne după lectura unui text este o sinteză a lui.

a. dacă nodul nu are un frate nenuclear în stânga, atunci expresia nervurii este v (v. Figura 2);

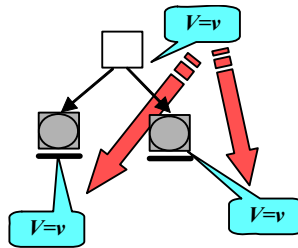


Figura 2: Expresia nervurii unui nod nuclear fara frate satelit în stânga

Definiția exprimă faptul că secvența de unități necesară înțelegerii/ rezumării, în contextul total, a unei întinderi nucleare de text ce are ca frate în structură o altă întindere nucleară necesită aceeași secvență de unități ca și cea necesară înțelegerii/rezumării, în contextul total, a reuniunii celor două întinderi. Cu alte cuvinte, o întindere nucleară ce este frate, în structură, întinderii nucleare curente este esențială înțelegerii/rezumării întinderii curente.

b. dacă nodul are un frate nenuclear în stânga de *head* h , atunci expresia nervurii lui este $seq(mark(h), v)$ (v. Figura 3);

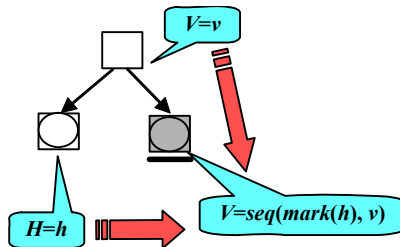


Figura 3: Expresia nervurii unui nod nuclear având un frate satelit în stânga

Secvența de unități necesară înțelegerii/rezumării, în contextul total, a unei întinderi nucleare de text ce are ca frate stâng în structură o întindere nenuclăară necesită, suplimentar față de secvența necesară înțelegerii în contextul total a întinderii acoperită de

nodul parinte (comunicata de expresia nervura a nodului parinte) si secventa *head* a întinderii frate stângi (adica cele mai importante unitati din întinderea stânga). Considerarea, în expresia nervurii întinderii nucleare curente, a expresiei *head* a întinderii nenuclare frate stângi, corespunde, prin prisma definitiei 2a, cu atribuirea întinderii stângi a calitatii de a se comporta ca un nucleu. Marcarea contributiei satelitului frate stâng prin functia *mark()* face însa aceasta revizuire a nuclearitatii lui, una cu valoare temporara, dupa cum se va dovedi mai jos, în definitia 3b.

3. Pentru fiecare nod nenuclear de *head* h , al carui parinte are nervura v :
 - a. daca nodul este descendentul stâng al parintelui sau, atunci expresia nervurii este $seq(h, v)$;

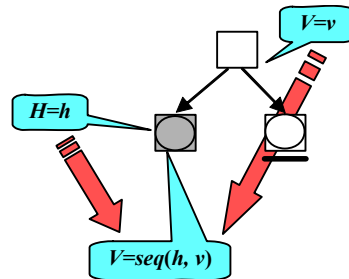


Figura 4: Expresia nervurii unui nod satelit stâng

Definitia exprima faptul ca pentru a înțelege/rezuma, în contextul total, o întindere nenuclara de text ce este descendent stâng, în structura, nodului parinte, la secventa de unitati necesara înțelegerii/rezumarii contextului total (contribuita de expresia nervura a parintelui) trebuie adaugate cele mai importante unitati din întinderea proprie (contribuite de expresia *head* proprie). Sa observam ca în expresia nervurii nodului parinte, care mosteneste expresii *head* ale nodurilor superioare, nu poate razbate influenta unui fiu satelit al sau, deci numai includerea *head*-ului fiului satelit, direct în expresia nervurii sale poate completa aceasta influenta.

- b. daca nodul este descendentul drept al parintelui sau, atunci expresia nervurii lui este $seq(h, simpl(v))$.

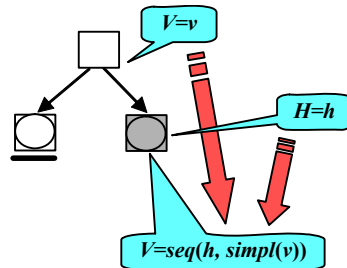


Figura 5: Expresia nervurii unui nod satelit drept

Pentru a înțelege, în contextul total, o întindere nenucleară de text ce este descendent pe dreapta al nodului parinte, la secvența de unități necesară înțelegerii/rezumării contextului total (contribuită de expresia nervurii a parintelui) și din care s-au sters unitățile marcate trebuie adăugate cele mai importante unități din întinderea proprie (contribuite de expresia *head* proprie). În acest fel, dacă expresia nervurii a nodului parinte nu conține unități marcate (prin contribuția definiției 2b), atunci expresia nervurii a unui satelit drept nu diferă de expresia nervurii a aceluiași satelit ce ar fi fost poziționat pe stânga (conform definiției 3a). Dacă însă nervurii parintelui conține unități marcate, atunci acestea dispar din expresia nervurii satelitului drept. Cum, conform definiției 2b, unitățile marcate pot fi contribuite doar de un satelit stâng, frate al celui mai apropiat ascendent nuclear al întinderii curente, urmează ca definiția curentă exprima o proprietate de blocare a accesibilității dinspre un satelit plasat în dreapta unui nucleu către un satelit plasat în stânga sa (v. Figura 6).

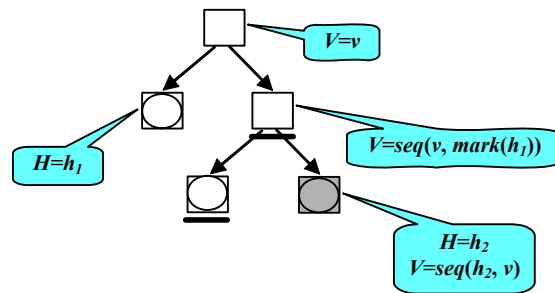


Figura 6: Simplificări în calculul expresiei nervurii a unui satelit drept:

$$V = \text{seq}(h_2, \text{simpl}(\text{seq}(v, \text{mark}(h_1)))) = \text{seq}(h_2, \text{seq}(v)) = \text{seq}(h_2, v)$$

Daca semnificatia expresiei nervurii unui nod oarecare din structura este particularizata la un nod terminal, obtinem: **expresia nervurii unei unitati de discurs reprezinta secventa unitatilor de discurs care sunt semnificative pentru a intelege/rezuma, în contextul întregului text, însasi unitatea de discurs în cauza.** Printre altele, aceasta înseamna ca expresia nervurii unei unitati de discurs este suficienta pentru a interpreta toate referintele anaforice continute în unitate.

4. Relatia dintre structura de discurs si referentialitate

Ipoteza pe care o avansam este ca exista doua tipuri de procese de rezolutie anaforica: **evocative** (sau **immediate**) si **post-evocative** (sau **inferentiale**). Procesele evocative, cele mai frecvente, sunt rapide si pot fi realizate prin orice mijloace de evocare referentiala, inclusiv cele fragile (de tipul subiectelor vide si pronomelor). Ele dau textului fluanta si-l fac coeziv. Cele post-evocative sunt mult mai putin frecvente decât cele evocative, necesita o încarcare inferentiala mai mare pentru procesarea lor si utilizeaza mijloace referentiale tari (nume proprii, substantive comune articulate).

Vom asocia spatiul de cautare al proceselor evocative unui **domeniu de accesibilitate referentiala evocativa** sau **imediata** (*domain of evocative accessibility – dea*) pe baza definitiei nervurii si al urmatoarelor observatii:

- **natura semantica a relatiei anaforice** [Halliday, Hassan, 1976]: o relatie anaforica are doi termeni: anaforul si antecedentul. Anaforul este reprezentat de o expresie referentiala a carei natura este textuala. Natura semantica a relatiei anaforice trebuie înțeleasa ca rasfrângându-se asupra antecedentului care nu trebuie identificat cu o anumita expresie referentiala ce precede în text anaforul ci cu o reprezentare a acesteia într-un plan semantic în asa fel încât semnificatia anaforului se construiește din cea a antecedentului însusi iar nu a semnificatiei lui. În cazul particular al unui lant co-referential acest lucru înseamna ca antecedentul este "realizat" repetat în text în aceeasi entitate de discurs. Expresiile co-referentiale "ancoreaza" în diverse pozitii ale textului entitatea de discurs.
- **dinamica incrementală a interpretării discursului**: un discurs este un text în procesul citirii ori ascultării lui de catre un subiect (om sau masina). Când citirea/ascultarea unui text s-a terminat discursul este încheiat si ceea ce ramâne este o reprezentare a lui în memoria subiectului. De asemenea, la un moment dat pe parcursul interpretării unui text, anumite elemente ale discursului pot fi plasate privilegiat în sfera atentiei [Grosz, Sidner, 1986, Sidner, 1983, Walker, 1996], iar trecerea de la o unitate de discurs la urmatoarea poate produce schimbari în structura memorata ce configureaza sfera attentionala.

- **natura cognitiva comuna a anaforei si a cataforei:** dintr-un punct de vedere cognitiv, toate referintele anaforice se fac dinspre expresii referentiale (entitati textuale) catre entitati ale discursului (entitati semantice) deja introduse de discursul trecut. Acest lucru înseamna ca într-o limba în care textul se noteaza de la stânga spre dreapta nu exista referinte anaforice spre dreapta. Distinctia dintre anafora si catafora, devine, în aceasta viziune care încearca sa reconstituie procesele cognitive ce stau la baza înțelegerii textelor (cu sau fara scopul simulării lor pe masina), inutila. În aceeași maniera în care, într-o anafora, un antecedent este o entitate de discurs propusa de o expresie referentiala ce precede anaforul si pe care anaforul o refera apoi, pronumele ce precede un nume într-o catafora propune o reprezentare, mai saraca, pe care numele o refera si o completeaza în acelasi timp [Cristea, Dima, 2001]. Acest lucru atribuie interpretării discursului o unica direccionalitate, care corespunde axei timpului lecturii, si care este cea a desfășurării liniare a textului (pentru limbile europene, de exemplu, de la stânga la dreapta). Relatia de referentialitate trebuie deci sa se proiecteze pe aceasta axa, dinspre entitati "noi" catre entitati "vechi", mereu catre înapoi pe axa timpului lecturii.

Ex. 7

1. Pentru ca ϕn -a vrut sa-si lase tata singur,
2. Ion a renuntat la concediu.

Expresia referentiala vida de pe pozitia de subiect a unitatii de discurs 1 propune o entitate de discurs caracterizata cel mult de o descriere [**type human**] (contribuita, cel mai probabil, de surse de cunoastere de natura pragmatica: cineva care nu poate sa-si lase tatal singur trebuie sa fie o persoana). Apoi, substantivul propriu *Ion*, din unitatea 2, refera entitatea construita precedent si o completeaza pâna la o reprezentare: [**type human, name Ion**].

Corelarea definitiei nervurii cu observatiile de mai sus, conduce la definirea domeniului de accesibilitate referentiala evocativa ca fiind format din toate unitatile de discurs care preced unitatea în care se gaseste expresia referentiala (si din care au fost îndepartate eventualele marcate, ce îndeplineau un rol de memorie temporara):

$$dea(u) = pref(u, unmark(V(u))).$$

Definitia *dea* formalizeaza prima conjectura a VT (sau a coeziunii), care pune în legatura accesibilitatea referentiala imediata de structura de discurs: antecedentii expresiilor referentiale dintr-o unitate de discurs *u* se gasesc, cu precadere, printre entitatile de discurs ancorate în unitatile ce preced pe *u*, inclusiv *u*, în expresia nervurii acesteia.

Paul Cornea [1998] vorbeste despre recodificarea sensului si memorizarea. El pune în evidenta trei tipuri de memorie, ce apar, de altfel, la mai multi cercetatori

[Kinntsch, Van Dijk, 1975, Schank, Abelson, 1977, Walker, 1996]: memoria imediata, memoria de scurta durata (de termen scurt – MST) si cea de lunga durata (de termen lung – MLT). Memoria imediata este un sistem de stocaj senzorial al informatiilor, retinerea urmelor din ultima jumătate de secunda. MST conserva câteva secunde informatia. Lungimea acestei memorii pare a fi de 7 ± 2 semne (cuvinte, cifre, litere – functie de context, v. si [Miller, 1956]; alti cercetatori apreciaza acest “empan” mijlociu la 13-15 cuvinte, la un lector lent fiind de 8 cuvinte, la unul rapid – de 16-20, de ex. [Richadeau, 1969] – citat în [Cornea, 1998] p. 166).

Constructia structurii de discurs se face dinamic, în actul lecturii. Sa ignoram un posibil proces de multi-interpretare ce poate duce la sintetizarea simultana a mai multor constructii alternative din care sa se selecteze, în urma unui proces de dezambiguizare, una sau mai multe structuri arborescente finale. Arborele însusi poate fi considerat rezumat în diverse grade, conform capacitatii de memorare a subiectului. Daca unitatea curenta este u_n , sa notam AR_n arborele de structura rezumat, la momentul prelucrării unitatii u_n . Nervura acesteia, culeasa pe AR_n , este $V(u_n)$, iar domeniul ei de accesibilitate imediata $dea(u_n)$. Noi credem ca MST poate fi considerata o fereastră de lungime 7 ± 2 semne în directa legatura cu $dea(u_n)$: fie 7 ± 2 unitati din aceasta secventa, fie tot atâtea structuri evenimentiale – ca reprezentari ale unitatilor de discurs, fie înca numai simboluri (cuvinte etc.) culese din acest sir de unitati. Tranzitarea la urmatoarea unitate, u_{n+1} , înseamna înlocuirea memoriei de scurta durata $dea(u_n)$ cu $dea(u_{n+1})$. Acest lucru duce uneori la o simpla prelungire a domeniului de accesibilitate precedent, alteori la o alterare a lui prin stergerea unor unitati si adaugarea altora, de fiecare data domeniul încheindu-se cu unitatea curenta. MST este asadar o proiectie a unui sir de unitati de discurs (sau de microstructuri suportate de unitati) decupate din structura dinamica curenta. Modificarile ce apar în sirul MST reflecta schimbarile de focalizare, în parcurgerea discursului. Componenta acestui sir este influentata de uitare (deci de un proces de abstractizare) si de modificarea de interes curenta în parcurgerea discursului. Când interesul s-a mutat pe o alta axa, componenta nervurii si, de aici, a domeniului de accesibilitate imediata sunt si ele actualizate. Includerea sau excluderea din MST a unor unitati de discurs în ritmul citirii, pentru ca *dea* evolueaza eliminând unele unitati si “redesteptând” altele “uite”, amintesc de procesele de “chemare” în sfera atentiei ale memoriei *cash* a lui Walker [Walker, 1996]. Pe de alta parte, structura memorata (rezumata) a discursului este pastrata în MLT si folosita pentru aducerea în prim plan a unitatilor de interes curent ce au fost temporar retrogradate de o comutare a atentiei într-o alta directie. Procesele evocative se desfasoara asadar în memoria de scurta durata. Pe de alta parte, procesele post-evocative sunt procese de rezolutie anaforica de natura inferentiala, ce presupun un anumit efort de regasire a unei entitati de discurs într-o zona a memoriei de lunga durata sau evoca entitati ale cunoasterii generice din sfera culturala a subiectului. Noi credem ca aceste procese se dezvoltă tot pe structura de discurs dezvoltata deja, iesind din *dea*, când rezolutia a esuat acolo.

Dintr-un punct de vedere ce se concentreaza asupra relatiei dintre referintele anaforice si structura de discurs, celor doua tipuri de procese anaforice le corespund **referinte evocative**, respectiv **post-evocative** (sau **inferentiale**). Diferenta dintre ele este

ca primele apar când lantul retroactiv al unitatilor ce ancoreaza expresii aflate în relatii referentiale intersecteaza domeniul de accesibilitate referentiala imediata al unitatii anaforului în cel puțin înca un punct decât unitatea anaforului, pe când în cazul referintelor post-evocative nu exista aceasta intersectie dubla. În [Cristea *et al.*, 2000, Cristea, 2000] referintele evocative sunt, mai departe, detaliate în **directe** si **indirecte**.

În referintele directe a doua unitate de intersectie este cea mai recenta liniar unitate ce ancoreaza aceeasi entitate de discurs ca si anaforul (în cazul relatiei de co-referinta) sau o entitate corelata functional cu aceasta (în cazul unei relatii de referinta functionala). În referintele indirecte intersectia *dea* cu lantul co/func-referential se realizeaza într-o unitate mai departata decât cea mai recenta liniar de unitatea anaforului. În referintele inferentiale lantul retroactiv al legaturilor anaforice al anaforului nu intersecteaza *dea* (în Figura 7 lantul legaturilor anaforice este reprezentat punctat, iar *dea* printr-o linie groasa).

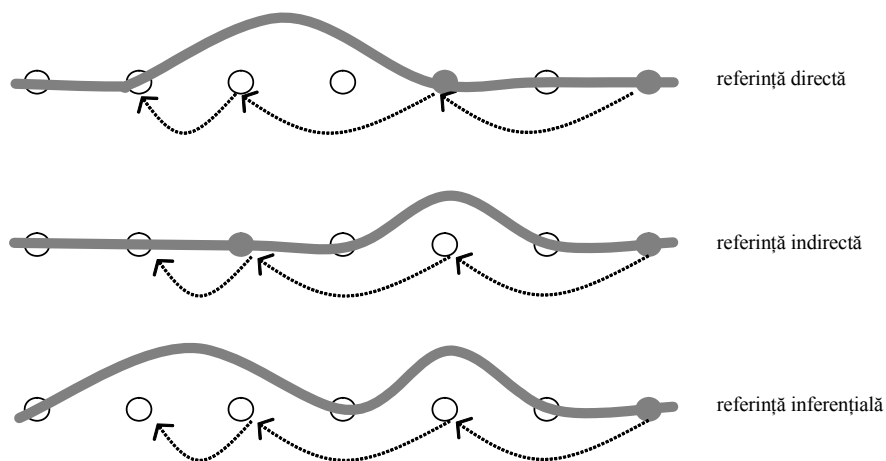


Figura 7: Referinte evocative si post-evocative

O categorie particulara de referinte post-evocative sunt **referintele pragmatice** (ce pot fi numite si **pseudo-referinte**). În acest tip de referinte participa expresii referentiale care pot fi interpretate fara un antecedent pentru ca interpretarea lor se bazeaza pe cunostinte exterioare textului, ce vin din cunoasterea comuna asupra lumii, deci din pragmatica. Desi exista cel puțin înca o expresie referentiala în text ce realizeaza aceeasi entitate de discurs, expresiile co-referentiale pot sa nu aiba, în mod necesar, o reprezentare unica, fara ca prin acesta înțelegerea textului sa sufere.

Recunoasterea antecedentului se datoreaza, în toate cazurile, unor procese de *pattern-matching* îmbogățite cu euristici, în care intervin structura de caracteristici morfo/sintactico/semantice ce definesc anaforul și structurile de caracteristici ce definesc entitățile de discurs deja introduse [Cristea, Dima, 2001, Cristea *et al.*, 2002a].

5. Relatia dintre structura de discurs si cursivitate

5.1 Linii de argumentatie

Expresiile nervura ale unitatilor ce compun un discurs arata tot atâtea moduri diferite în care poate fi citit acel discurs. Fiecare în parte da o rezumare a discursului prin prisma unitatii de discurs curente. Atunci când interesul este orientat catre un anumit episod al povestirii, putem sari peste pasaje întregi pentru a ne concentra asupra manierei în care elementul de interes se leaga cu ansamblul discursului. În acelasi fel, putem avea în vedere o alta pista și atunci lectura focalizeaza un alt fir de interes. Acest nou fir poate sa aiba elemente în comun cu primul dar poate, de asemenea, sa incorporeze și altele noi. Fiecare fir în parte poate pune în evidenta anumite particularitati, legate însa strâns de linia principala a discursului. Toate aceste sub-discursuri sunt coerente și nu exista referinte anaforice pentru a caror interpretare sa avem nevoie de fragmente aflate în afara rezumatului însusi. Acest lucru înseamna ca traseele referentiale ale rezumatului contin suficiente elemente care sa duca la recuperarea înțelesului anaforilor.

Sa luam urmatorul text:

Ex. 8

1. *Piton primise-n taina porunca de la Hera sa-l pîndeasca pe Apolo,*
2. *cînd va trece prin munte,*
3. *si sa-i rapuna viata.*
4. *Hera-l ura pe fiul cel nou nascut al Letei,*
5. *pentru ca sotul sau, prea puternicul Zeus, tinea mai mult la dînsul decit la fiii ei: Hefaistos si Ares.*
6. *Cînd a ajuns Apolo în muntele Parnas,*
7. *dihania uriasa s-a avîntat spre dînsul,*
8. *dornica sa-l ucida.*
9. *Dar zeul si-a întins arcul.*
10. *A tras prima sageata.*
11. *Erau doar patru zile de cînd vazuse lumea,*
12. *si întiia lui sageata a si nimerit monstrul.*

Alexandru Mitru - *Legendele Olimpului*, Editura Tineretului, 1966

Structura de discurs a acestui text este cea din Figura 8. Tabela 1 da expresiile nervura si domeniile de referentialitate evocativa ale nodurilor terminale. În coloana $dea(u)$ au fost, totodata, marcate în aldine domenii de referentialitate imediata maximale vis-r-vis de relatia de incluziune (cele mai lungi trasee dea). Astfel $dea(1) \subseteq dea(2) \not\subseteq dea(3) \subseteq dea(4) \subseteq dea(5) \not\subseteq dea(6)$ s.a.m.d. Vom numi aceste secvente care întrerup lanturi de incluziuni **linii de argumentatie** (la), în cazul nostru: 1 2, 1 3 4 5, 1 3 6 7, 1 3 7 8 si 1 3 7 9 10 11 12. Daca $la(u_1)$ precede imediat $la(u_2)$, atunci în $la(u_2)$ se regasesc domeniile tuturor unitatilor dintre u_1+1 si u_2 . În particular, în $la(u_2)$ se regasesc unitatile ce preced imediat unitatea u , pentru orice u între u_1+1 si u_2 , în domeniul lor de accesibilitate imediata (adica acel domeniu care confera discursului maximul de coerenta). Cu alte cuvinte, pe $la(u_2)$ putem aplica definitiile CT de calculare a tranzitiilor pentru orice u între u_1+1 si u_2 .

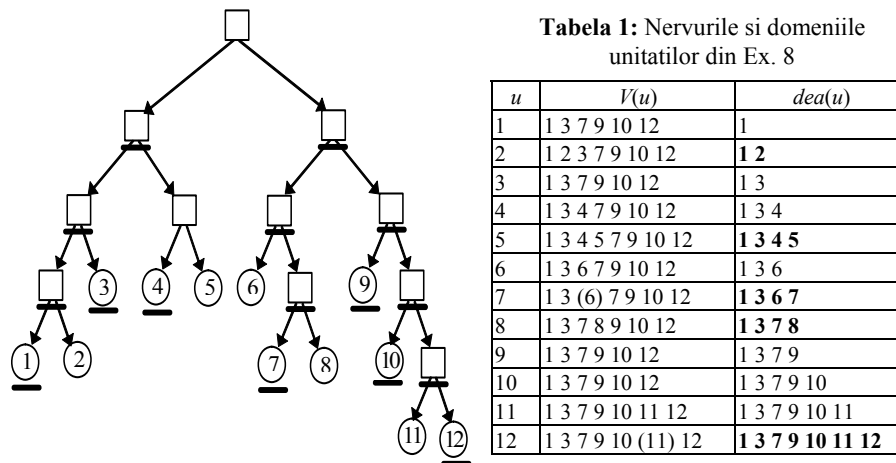


Figura 8: Structura de discurs a Ex. 8

5.2. O generalizare a CT

Urmând recomandările teoriei centrelor, sa presupunem ca marcam tranzitiile ce apar între unitati de discurs cu scoruri care sa dea un grad al usurintei de prelucrare:

CONTINUARE	(CON)	4
RETINERE	(RET)	3
SCHIMBARE LINA	(SSH)	2
SCHIMBARE ABRUPTA	(ASH)	1
LIPSA Cb	(-)	0

În felul acesta, tranzițiile line primesc scoruri mari, cele abrupte, scoruri mici. Însușind aceste scoruri pentru fiecare unitate a unui segment (segment, în spiritul AST) vom avea un scor al segmentului. Sa notam un scor în spiritul CT al unui segment s cu S_{CCT}^s (CCT de la *Classical Centering Theory*). El ne va da o masura a usurintei de interpretare a segmentului: cu cât un segment s , în totalitatea lui, e mai fluent, cu atât scorul lui va fi mai mare și cu cât el este mai abrupt, mai dificil de prelucrat, cu atât scorul lui va fi mai scăzut. În fine, sa adunam aceste scoruri pentru toate segmentele discursului, într-un scor al sumei segmentelor S_{CCT} :

$$S_{CCT} = \sum_s S_{CCT}^s$$

Sa ne imaginam acum ca fortam nota și calculam aceste scoruri și dincolo de granitele de segment, deci inclusiv în punctele de frontiera dintre segmente. Sa notam acest scor global cu S_{CCT}^G . În scorul global S_{CCT}^G contribuie cu scoruri de tranziții toate unitatile cuprinse între a doua unitate și ultima. În mod normal tranzițiile în punctele de trecere între segmente ar trebui sa fie foarte abrupte, cotate deci slab ori zero, și deci scorul global atasat textului n-ar trebui sa fie modificat semnificativ. Daca apare totuși o diferenta, ea trebuie sa fie datorata unor tranziții accidentale peste granta de segment. În orice caz trebuie sa avem $S_{CCT}^G \geq S_{CCT}$.

Sa procedam acum în mod analog, ca suport folosind de data aceasta liniile de argumentatie iar nu secventele liniare de unitati ale segmentelor în sensul clasic. Datorita comportamentului lor similar segmentelor, putem numi liniile de argumentatie **segmente în sens ierarhic**. Sa notam S_{HCT}^s (HCT de la *Hierarchical Centering Theory*) suma scorurilor unitatilor aparținând unei linii de argumentatie (segment ierarhic) s . Ca sa dam o masura a fluentei discursului în acceptiunea ierarhica, similara scorului global S_{CCT}^G , în calculul scorului global al discursului în sens ierarhic nu va trebui sa repetam contributiile unitatilor ce apar în mai mult decât o singura linie de argumentatie. Daca notam $S_{HCT}^{s'}$ scorul unui segment ierarhic s' în care am pastrat numai unitatile noi fata de segmentul anterior, atunci scorul global ierarhic al discursului este:

$$S_{HCTG} = \sum_{s'} S_{HCT}^{s'}$$

Cea de a doua conjectura a VT (a coerenței): Scorul global în sensul ierarhic al unui discurs este mai bun sau cel puțin egal decât scorul global în sensul clasic: $S_{HCT}^G \geq S_{CCT}^G$.

Pentru un anumit detaliu de granularitate în definirea segmentelor în sens clasic, unui segment în sens clasic îi corespunde o secventa de nervura, deci o portiune a unei linii de argumentatie. În spiritul acestei observatii, cea de a doua conjectura enunta prezumtia ca tranzițiile la distanta lunga calculate în lungul nervurilor sunt sistematic mai line decât tranzițiile accidentale la granitele dintre segmente. Sa notam ca aceasta presupozitie este

conforma unor observatii facute de autori precum Passonneau [1995] si Walker [1998], furnizând totodata o explicatie pentru rezultatele lor.

În cele ce urmeaza prezentam o analiza comparativa clasic-ierarhic care probeaza ipoteza coerenței, pe discursul din Ex. 8.

Tabela 2: Analiza Ex. 1 în maniera CCT

n	u_n	$C_r(u_n)$	$C_b(u_n)$	Traz.	Scor
1	<i>Piton primise-n taină poruncă de la Hera să-l pîndească pe Apolo,</i>	[Piton], [Hera], [Apolo],	[Piton]	-	-
2	<i>cînd Ø va trece prin munte,</i>	Ø = [Apolo], [munte]	[Apolo]	SSH	2
3	<i>și Ø să-i răpună viața.</i>	Ø = [Piton], i = [Apolo], [viața]	[Apolo]	RET	3
4	<i>Hera-l ura pe fiul cel nou născut al Letei,</i>	[Hera], [Leta], <i>fiul cel nou-născut al Letei</i> = [Apolo]	-	-	0
5	<i>pentru că soțul său, prea puternicul Zeus, ținea mai mult la dînsul decît la fiii ei: Hefaistos și Ares.</i>	[Zeus], <i>său</i> = [Hera], <i>dînsul</i> = [Apolo], [Hefaistos], [Ares]	[Hera]	ASH	1
6	<i>Cînd a ajuns Apolo în muntele Parnas,</i>	[Apolo], [munte]	[Apolo]	SSH	2
7	<i>dihania uriașă s-a avîntat spre dînsul,</i>	<i>dihania uriașă</i> = [Piton], <i>dînsul</i> = [Apolo]	[Apolo]	RET	3
8	<i>Ø (era) dornică să-l ucidă.</i>	Ø = [Piton], l = [Apolo]	[Piton]	SSH	2
9	<i>Dar zeul și-a întins arcul.</i>	<i>zeul</i> = [Apolo], [arcu]	[Apolo]	SSH	2
10	<i>Ø A tras prima săgeată.</i>	Ø = [Apolo], [săgeata]	[Apolo]	CON	4
11	<i>Erau doar patru zile de cînd Ø văzuse lumea,</i>	[4 zile], Ø = [Apolo], [lumea]	[Apolo]	RET	3
12	<i>și întia lui săgeată a și nimerit monstrul.</i>	<i>lui</i> = [Apolo], [săgeata], <i>monstrul</i> = [Piton]	[Apolo]	CON	4
Total					26

În constructia tabelului de mai sus am presupus ca toate referintele anaforice au fost corect rezolvate. Unitatile carora le corespund tranzitiile listate în tabela sunt cele ale caror numere apar în caractere aldine în prima coloana, adica 2-12, în numar total de 11. Scorul total de 23 corespunde unei scor mediu pe tranzitie de $26/11=2,36$, ceea ce înseamna ca textul, conform aprecierii CT, se comporta, în medie, intermediar între o schimbare lina (SSH) și o retinere (RET), mai apropiat de o schimbare lina.

Daca luam în calcul liniile de argumentatie indicate de nervuri, pot fi puse în evidenta 5 sub-discursuri, în lungul carora vom calcula, de asemenea, tranzitiile. În tabelele 3÷7 de mai jos unitatile pentru care consideram tranzitiile sunt, de asemenea, indicate în caractere aldine în prima coloana. Sa remarcam ca citirea textelor date de liniile de

argumentatie produce, în toate cazurile, discursuri perfect coerente. În ansamblu, doar câte o tranzitie este calculata pentru fiecare unitate, la fel ca si în interpretarea clasica.

Tabela 3: Analiza HCT a primei linii de argumentatie, secventa de unitati 1-2

n	u_n	$C(u_n)$	$C_b(u_n)$	Traz.	Scor
1	<i>Piton promise-n taină poruncă de la Hera să-l pîndească pe Apolo,</i>	[Piton], [Hera], [Apolo],	[Piton]	-	-
2	<i>cînd \emptyset va trece prin munte,</i>	\emptyset = [Apolo], [munte]	[Apolo]	SSH	2
Total					2

Tabela 4: Analiza HCT a celei de a doua linii de argumentatie, secventa de unitati 1-3-4-5

n	u_n	$C(u_n)$	$C_b(u_n)$	Traz.	Scor
1	<i>Piton promise-n taină poruncă de la Hera să-l pîndească pe Apolo,</i>	[Piton], [Hera], [Apolo],	[Piton]	-	-
3	<i>și \emptyset să-i răpună viața.</i>	\emptyset = [Piton], i =[Apolo], [viața]	[Piton]	CON	4
4	<i>Hera-l ura pe fiul cel nou născut al Letei,</i>	[Hera], [Leta], <i>fiul cel nou-născut al Letei</i> =[Apolo]	-	-	0
5	<i>pentru că soțul său, prea puternicul Zeus, ținea mai mult la dînsul decît la fiii ei: Hefaistos și Ares.</i>	[Zeus], <i>său</i> =[Hera], <i>dînsul</i> =[Apolo], [Hefaistos], [Ares]	[Hera]	ASH	1
Total					5

Se constata ca tranzitia RET a unitatii 3 catre 2 din analiza CCT s-a transformat într-o tranzitie CON, pe nervura, dinspre 3 catre 1.

Tabela 5: Analiza HCT a celei de a treia linii de argumentatie, secventa de unitati 1-3-6-7

n	u_n	$C_f(u_n)$	$C_b(u_n)$	Traz.	Scor
1	<i>Piton primise-n taină poruncă de la Hera să-l pîndească pe Apolo,</i>	[Piton], [Hera], [Apolo],	[Piton]	-	-
3	<i>și Ø să-i răpună viața.</i>	Ø= [Piton], i=[Apolo], [viața]	[Piton]	-	-
6	<i>Cînd a ajuns Apolo în muntele Parnas,</i>	[Apolo], [munte]	[Apolo]	SSH	2
7	<i>dihania uriașă s-a avîntat spre dînsul,</i>	dihania uriașă=[Piton], dînsul=[Apolo]	[Apolo]	RET	3
Total					5

Tabela 6: Analiza HCT a celei de a patra linii de argumentatie, secventa de unitati 1-3-7-8

n	u_n	$C_f(u_n)$	$C_b(u_n)$	Traz.	Scor
1	<i>Piton primise-n taină poruncă de la Hera să-l pîndească pe Apolo,</i>	[Piton], [Hera], [Apolo],	[Piton]	-	-
3	<i>și Ø să-i răpună viața.</i>	Ø= [Piton], i=[Apolo], [viața]	[Piton]	-	-
7	<i>dihania uriașă s-a avîntat spre dînsul,</i>	dihania uriașă=[Piton], dînsul=[Apolo]	[Piton]	-	-
8	<i>Ø (era) dornică să-l ucidă.</i>	Ø= [Piton], l=[Apolo]	[Piton]	CON	4
Total					4

Se constata ca tranzitia SSH a unitatii 8 catre 7 din analiza CCT s-a transformat într-o tranzitie CON, pe nervura, tot între 8 si 7 (C_b -ul unitatii 7 s-a schimbat din [Apolo] în [Piton], pentru ca, pe nervura lui 8, precedenta unitate a lui 7 este acum 3, iar nu 6 ca în secventa liniara).

Tabela 7: Analiza HCT a ultimei linii de argumentatie, secventa de unitati 1-3-7-9-10-11-12

n	u_n	$C(u_n)$	$C_b(u_n)$	Traz.	Scor
1	<i>Piton primise-n taină poruncă de la Hera să-l pîndească pe Apolo,</i>	[Piton], [Hera], [Apolo],	[Piton]	-	-
3	<i>și Ø să-i răpună viața.</i>	Ø = [Piton], i=[Apolo], [viața]	[Piton]	-	-
7	<i>dihania uriașă s-a avîntat spre dînsul,</i>	<i>dihania uriașă</i> =[Piton], <i>dînsul</i> =[Apolo]	[Piton]	-	-
9	<i>Dar zeul și-a întins arcul.</i>	<i>zeul</i> = [Apolo], [arcul]	[Apolo]	SSH	2
10	<i>Ø A tras prima săgeată.</i>	Ø = [Apolo], [săgeata]	[Apolo]	CON	4
11	<i>Erau doar patru zile de cînd Ø văzuse lumea,</i>	[4 zile], Ø = [Apolo], [lumea]	[Apolo]	RET	3
12	<i>și întia lui săgeată a și nimerit monstrul.</i>	<i>lui</i> =[Apolo], [săgeata], <i>monstrul</i> =[Piton]	[Apolo]	CON	4
Total					13

Însumând scorurile tranzitiilor pentru toate liniile de argumentatie se obtine scorul total: 29, ceea ce corespunde unei tranzitii medii a discursului, calculata conform HCT de $29/11=2,63$, asadar o tranzitie mai apropiata de retinere, mai buna decât scorul mediu calculat conform CCT.

6. Validarea conjecturilor VT

Validarea conjecturilor VT s-a realizat pe corpusuri adnotate la structura si la legaturi co-referentiale. Astfel în [Cristea *et al.*, 1998] se raporteaza o investigatie efectuata pe texte în limbile engleza, franceza si româna ce au însumat un total de 176 de unitati de discurs. Plecând de o adnotare în maniera RST a structurii de discurs, un program a calculat expresiile nervurilor unitatilor. Pentru verificarea conjecturii coeziunii, utilizând adnotarea legaturilor referentiale s-a calculat apoi procentajul referintelor directe, indirecte si pragmatice. În medie 99,1% dintre referinte se încadreaza acestor trei categorii (87,1% directe, 8,5% indirecte si 3,5% pragmatice). Pentru verificarea conjecturii coerentei, suplimentar marcajelor de structura si lanturi co-referentiale s-au marcat manual, pentru fiecare unitate, C_b -ul, în varianta clasica si în varianta ierarhica, si s-au calculat tranzitiile în cele doua variante. Scorul S_{HCT} a fost mai bun decât scorul S_{CCT} în toate cazurile (scorurile medii pe tranzitie au fost de 2,03 în varianta ierarhica fata de 1,89 în cea clasica).

În [Cristea *et al.*, 2000] se raporteaza experimente care au urmarit sa compare potentialul modelelor ierarhice, precum cele bazate pe VT, de a regasi un antecedent într-o plaja de cautare data fata de modelele lineare (modele ce presupun o parcurgere lineara a textului dînspre unitatea anafurului spre începutul textului). Pentru aceasta s-au utilizat 30 de texte englezesti (însumând aproximativ 1560 de unitati de discurs), adnotate la structura

RST și lanțuri co-referențiale. Presupunând o plajă de căutare de doar 2 unități, căutarea pe nervură a adus cu aproximativ 16% mai mulți antecedente decât căutarea liniară. După cum era de așteptat, pe măsură ce lungimea textului căutat crește cele două tipuri de modele se apropie în ceea ce privește potențialul de a regăsi legături co-referențiale. O căutare ierarhică înapoi într-o plajă de 5 unități rezolvă potențial doar 70% dintre anafore, pentru că o performanță potențială de 90% să poată fi atinsă doar dacă se organizează o căutare într-o lungime de 12 unități pe nervură. O altă investigație a urmărit compararea efortului necesar regăsirii unui anumit antecedent în cele două tipuri de abordări (liniară și ierarhică), unde prin efortul necesar găsirii unui antecedent se înțelege numărul de unități de discurs ce separă, în domeniu, unitatea anaforului de unitatea celei mai recente ancorări în text a unui antecedent. Din nou modelele ierarhice, de tipul celui dat de VT, s-au dovedit superioare celor liniare: în corpusul folosit în experiment, care a conținut 1200 de expresii referențiale, spațiul de căutare pentru legături co-referențiale s-a redus cu aproximativ 800 de unități.

Un alt tip de investigație empirică [Ide, Cristea, 2000] a urmărit frecvența referințelor evocative în comparație cu cele post-evocative și depistarea unor corelații între tipul de referințe și puterea de evocare a anafurilor. Studiul a comparat prezicerile avansate de VT relativ la domeniul de referențialitate evocativă cu cele ale modelului stivă al AST, corelând excepțiile (referințe ce nu se supun prevederilor celor două teorii) cu puterea de evocare a anafurilor (pentru VT excepțiile marchează, evident, referințe din categoria celor inferențiale). Într-o ordine descendentă a puterii de evocare (v. și [Gundel *et al.*, 1993]) tipurile de anafuri care dau naștere la excepții sunt: referințe pragmatice > nume proprii > substantive comune > pronume. Pronumele constituie mijloace de referire foarte fragile. Un emițent al unui mesaj utilizează un pronume când e sigur că structura permite recuperarea cu ușurință a entității referite de pronume. Practic, exceptând câteva cazuri în care un pronume putea fi înțeles fără un antecedent (*our* în *our streets*, de exemplu), este imposibilă utilizarea unui pronume pentru a referi o entitate aflată în afara *de*. La extrema cealaltă se plasează referințele pragmatice ce-și recuperează antecedentul din cunoscutele discursului și numele proprii. Interesant este că această sortare descrescătoare a tipurilor de anafuri data de puterea de evocare se aliniază numărului de excepții raportate în cazul VT (56,3% – pragmatice, 22,7% – nume proprii, 16,0% – substantive comune și 5,0% – pronume) și nu are nici o semnificație în cazul AST (0,0% – pragmatice, 26,1% – nume proprii, 39,1% – substantive comune și 34,8% – pronume). Ea probează corectitudinea conjecturii coeziunii.

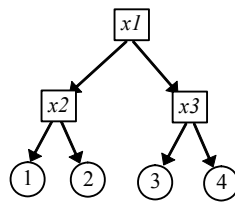
7. O proprietate de granularitate

Atunci când arborele de structură al discursului se modifică trecându-se de la o granularitate mai fină la una mai grosieră, constrângerea de accesibilitate, conjecturată de VT, se pastrează.

Demonstratie

Sa presupunem un arbore de discurs D , notat la structura si pe care s-au calculat expresiile *head* si nervura ale nodurilor. O operatie de marire a granularitatii poate fi efectuata daca o întindere de text, initial repartizata în mai multe unitati, si **careia îi corespunde un nod în structura initiala**, este "compactata" într-o singura unitate de discurs mai mare ce va lua locul nodului radacina din structura initiala. Pentru a vedea în ce masura o astfel de operatie poate afecta accesibilitatea vom investiga rezultatul aplicarii ei asupra expresiilor *head* si nervura.

Definitia expresiei *head*, punctul 1, obliga ca expresia *head* a ceea ce înainte de compactare era un nod interior, fie el n , sa fi fost data de concatenarea unui sir de etichete de noduri nucleare aflate în secventa de text subîntinsa de n . Sa notam acest nod, dupa compactare, cu o eticheta compusa din secventa nodurilor terminale pe care le acopera. De exemplu, pentru arborele:



daca subarborele cu radacina $x2$ ar fi compactat, atunci eticheta sa ar trebui sa fie notata 1-2, iar daca întregul arbore aflat sub $x1$ ar fi compactat, atunci eticheta sa ar trebui sa fie notata 1-2-3-4 (e imposibil sa avem un nod notat 2-3).

Acest lucru înseamna ca, aplicând o compactare asupra unui arbore, în expresiile *head* ale nodurilor sale secvente de noduri vor fi acum înlocuite cu etichete compuse care contin cel puțin aceleasi noduri, eventual mai multe, decât în expresiile originale. De exemplu, presupunând ca în arborele de mai sus, nodurile nucleare sunt $x2$, $x3$, 1 si 3, atunci, daca înainte de compactare am fi avut $head(x1)=1\ 3$, o expresie compusa din doua etichete, dupa compactarea întregului arbore vom avea $head(x1)=1-2-3-4$, adica o eticheta compusa, dar care include etichetele nodurilor ce apareau în expresia *head* originala. Vom numi astfel de expresii – expresii contrase si le vom nota cu $contr(e)$, unde e este expresia corespunzatoare înainte de compactare (avem deci $contr(1\ 3) = 1-2-3-4$). Sa remarcam ca secventele de etichete din expresiile contrase sunt formate întotdeauna din etichete de noduri adiacente, ceea ce permite comutarea functiilor seq si $contr$: $seq(contr(e_1), contr(e_2)) = contr(seq(e_1, e_2))$.

Vom demonstra mai întâi ca expresiile nervura ale nodurilor din arborele compactat sunt obtinute din expresiile nervura originale prin înlocuirea expresiilor *head* originale cu expresiile contrase. Investigând definițiile expresiilor nervura, se poate constata ca nici o alta modificare nu apare în expresiile nervura cu exceptia expresiilor contrase. Într-adevar, cazul 1 se transcrie: expresia nervura a radacinii arborelui compactat reprezinta expresia *head* contrasa a arborelui original, adica $contr(h)$, cu h – expresia *head* a radacinii arborelui original.

Sa presupunem acum ca ne aflam într-un nod n ale carui expresii *head* si nervura pe arborele original, necompactat sunt, respectiv h si v , iar $contr(h)$ este expresia *head* pe arborele compactat. Consideram mai întâi cazul când n este fiu al nodului radacina, a carui expresie *head* este $contr(h_0)$, unde h_0 reprezinta expresia *head* pe arborele necompactat. Daca n este nuclear, atunci conform cazului 2 (sectiunea 3), avem doua subcazuri:

- a) n nu are un frate nenuclear în stânga: atunci nervura sa este chiar nervura parintelui, adica $contr(h_0)$;
- b) n are un frate nenuclear în stânga de *head* $contr(h_1)$: nervura nodului n va fi $seq(mark(contr(h_1)), contr(h_0)) = seq(contr(mark(h_1)), contr(h_0)) = contr(seq(mark(h_1), h_0)) = contr(v)$;

Daca n este un nod nenuclear, atunci conform cazului 3, avem:

- a) n este în stânga: nervura sa este $seq(contr(h_0), contr(h)) = contr(seq(h_0, h)) = contr(v)$;
- b) n este în dreapta: nervura sa este $seq(simpl(contr(h_0)), contr(h)) = seq(contr(simpl(h_0)), contr(h)) = contr(seq(simpl(h_0), h)) = contr(v)$.

Folosind inductia, se probeaza în mod analog ca expresia nervura a nodului n este o expresie contrasa si pentru cazul în care n este un nod interior, nu neaparat imediat sub radacina, fiu al unui nod de nervura $contr(v_0)$.

Cum expresia accesibilitatii este definita ca un prefix al expresiei nervura din care au fost îndepartate marcasele, iar nervurile sunt expresii contrase, deci eventual continând mai multe etichete de noduri, înseamna ca orice referinta care pe arborele original satisface prima conjectura, cu alte cuvinte are loc între ultima unitate a unei expresii nervura si alta ce o precede, dupa compactare va satisface de asemenea conjectura, pentru ca nici o unitate nu a disparut din domeniu.

8. Discutii, aplicatii ale teoriei

Plecând de la o reprezentare a structurii de discurs similara celei din RST si în care esentiala este distinctia dintre nucleu si satelit, VT defineste nervura unui nod al arborelui ca secventa de unitati ale discursului ce sunt suficiente pentru a rezuma/interpreta întinderea de text acoperita de nod în contextul întregului discurs. Presupunerea principala pe care se bazeaza notiunea de nervura este ca *referintele inter-unitati sunt posibile cu*

precadere între unitati ce se afla într-o relatie structurala, chiar daca acestea sunt dispuse la distanta una de alta în text. Mai departe, referintele se realizeaza cu precadere spre unitati nucleare si doar în putine cazuri catre sateliti, reflectând intuitia ca nucleeele gazduiesc ideile principale ale discursului. Acest lucru se regaseste în calculul expresiei nervurii pe arbori (binari) polarizati-stânga (pe orice nivel exista un nucleu în stânga), în care orice referinta se realizeaza dinspre un nucleu sau un satelit catre un nucleu aflat în stânga (desi, nu orice nucleu). Facând uz de echivalarea modelului stiva al lui Grosz si Sidner [1986] cu structura de arbore utilizata de RST [Mann, Thompson, 1988], similaritate demonstrata de Moser si Moore [1996] si Marcu [1999], predictiile VT asupra accesibilitatii referentiale sunt consistente cu cele ale modelului stiva. În cazurile în care însa arborele de discurs nu e polarizat-stânga (exista cel puțin un satelit care precede nucleul sau, deci care apare ca frate stâng pe un nivel al structurii) VT ofera o interpretare mai naturala a accesibilitatii decât modelul stiva, corectând totodata slabiciunile acestuia. Într-adevar, într-o secventa *A*-satelit, *B*-nucleu, deci în care *B* domina *A* în termenii AST, *B* ar trebui sa apara în stiva positionat sub *A*, desi el este procesat în secventa dupa *A*. Totodata, VT formalizeaza intuitia ca într-o secventa de unitati *A*, *B*, *C*, unde *A* si *C* sunt sateliti ai lui *B*, *C* nu poate accesa *A* din cauza interpunerii unui nucleu, ce capteaza întreaga atentie.

Referentialitatea în lungul nervurilor este una naturala, usor de procesat si care, în general, nu necesita mijloace de evocare foarte puternice. Dimpotriva, iesirea din acest domeniu incumba utilizarea unor mijloace de evocare anaforica viguroase. Pe acest criteriu se face distinctia dintre referentialitate evocativa si ne-evocativa (sau inferentiala), referintele evocative fiind detaliate în directe si indirecte, iar între cele ne-evocative remarcându-se referintele pragmatice, ce nu necesita un antecedent pentru înțelegere.

În privinta coerenței discursului, VT utilizeaza domeniile de referentialitate pentru a introduce notiunea de linie de argumentatie si a deduce din ea pe cea de segment în sens ierarhic ce generalizeaza segmentul în sens clasic (asa cum este el utilizat în AST si CT). Totodata VT avanseaza conjectura ca segmentul în sens ierarhic da o mai corecta interpretare a portiunilor de discurs ce se comporta din punctul de vedere al coeziunii si coerenței ca un tot unitar. Aplicând concluziile CT relative la coerența discursului în lungul segmentelor în sens ierarhic CT poate fi generalizata pentru a o transforma într-o teorie globala a coerenței.

Au fost trecute în revista o seama de experimente care probeaza ca prezumtiile VT sunt corecte si independente de limba. Un aspect important îl constituie, de asemenea, faptul ca prezumtiile VT sunt stabile la schimbarea granularitatii în segmentarea discursului.

Aplicatiile VT se înscriu în trei directii importante: rezolutia anaforei, parsarea discursului si rezumarea automata. În [Cristea *et al.*, 2002a] si [Cristea *et al.*, 2002b] este descrisa o arhitectura care actioneaza ca un motor general si configurabil de rezolutie anaforica. Una dintre componentele oricarui model de rezolutie este o definitie a domeniului de referentialitate. Rezolutia anaforica se realizeaza, asadar, ghidata de structura de discurs.

În [Seretan, Cristea, 2002] se propune o abordare inversa, în care cunostinte asupra legarilor anaforice pot fi utilizate pentru corectarea structurii. Noi credem ca procesul de rezolutie anaforica si de construire a structurii de discurs sunt interdependente într-un asemenea grad încât în analiza de discurs ele trebuie sa aiba loc simultan. În interpretarea unui text exista o interconditionare reciproca între referinte si structura care trebuie sa conduca la obtinerea acelei reprezentari în care constrângerile, actionând ca forte, produc o stare de echilibru, ce trebuie sa fie un fel de stare de energie potentiala minima a sistemului. Oamenii dispun de un mecanism cognitiv care le permite sa ajunga în mod natural la cea mai plauzibila interpretare a unui text. Acest lucru este rasplatit de atingerea unei stari mentale „comfortabile” ce trebuie sa-si aiba suportul în satisfacerea la maxim a unui sistem de constrângeri. În [Tablan *et al.*, 1998] si [Cristea, 2000] se descrie un mecanism de parsare care modeleaza acest comportament uman. Prin combinarea unor scoruri contribuie de referinte cu scoruri contribuie de o analiza HCT se obtine cea mai fluida posibil structura de discurs (deci manifestând maxim de coerenta) si care prezinta maximul de referinte pe nervuri (fiind deci cea mai coeziva posibil).

Notiunea de *head* din VT este similara celei de multime de promovare (*promotion set*) pe care Marcu [Marcu, 2000] o utilizeaza pentru a obtine un rezumat ghidat de structura de discurs. Sa remarcam ca definitia nervurii presupune rezumarea ca o alternativa a înțelegerii unei unitati de discurs în context. Credem ca valentele teoriei nervurilor în realizarea unei strategii de rezumare focalizata [Mani, 2001] pe o anumita entitate sau segment de discurs au fost doar tangential studiate pâna acum [Sofronie, 1999], [Postolache, 2001] si merita atentie în abordarile viitoare. Credem, de asemenea, ca fiind interesanta o directie de studiu care sa aprecieze maniera în care nervura poate constitui un cadru de sub-specificare a structurii [Schilder, 2001], plecând de la observatia ca structuri diferite (dau nu fundamentale diferite) pot manifesta aceeasi expresie a nervurilor.

Bibliografie

- Brennan, S.E.; Walker Fredman, M. and Pollard, C.J. 1987. A centering approach to pronouns. *Proc. 25th Annual Meeting of ACL*, Stanford, p. 155-162.
- Cornea, P. 1998. Introducere în teoria lecturii, Editura Polirom, Iași.
- Cristea, D., and Webber, B.L. 1997. Expectations in Incremental Discourse Processing. *Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics*, Madrid.
- Cristea, D., Ide, N., and Romary, L. 1998. Veins Theory: A Model of Global Discourse Cohesion and Coherence, *Proceedings of the 17th Coling and the 36th Annual Meeting of the ACL (COLING-ACL'98)*, Montreal, Canada, p.281-185.
- Cristea, D., Ide, N., Marcu, D., and Tablan, M.V. 2000. An Empirical Investigation of the Relation Between Discourse Structure and Co-Reference. *Proceedings of the 18th*

-
- International Conference on Computational Linguistics COLING'2000*, Saarbrueken, p. 208-214.
- Cristea, D. 2000. An Incremental Discourse Parser Architecture. Christodoulakis, D. (Ed.) *Natural Language Processing - NLP 2000*, Second International Conference, Patras, Greece, Lecture Notes in Artificial Intelligence 1835, Springer, p. 162-175.
- Cristea, D. and Dima, G.E. 2001. An Integrating Framework for Anaphora Resolution. *Information Science and Technology*, Romanian Academy Publishing House, Bucharest, vol. 4, no. 3-4, p. 259-372.
- Cristea, D., Postolache, O.D., Dima, D.E., Barbu C. 2002a. AR-Engine – a framework for unrestricted co-reference resolution. *Proceedings of the Third International Conference on Language Resources and Evaluation, LREC'2002*, Las Palmas, Spain, p. 2000-2006.
- Cristea, D., Dima, D.E., Postolache, O.D., Mitkov, R. 2002b. Handling complex anaphora resolution cases. *Proceedings of the Discourse Anaphora and Anaphor Resolution Colloquium*, Lisbon, Portugal.
- deEugenio, B. 1990. Centering theory and the Italian pronominal system. *Proceeding of Coling*, p. 270-275.
- deEugenio, B. 1998. Centering in Italian. Prince, E., Joshi, A. and Walker, L. (eds.) *Centering in Discourse*, Oxford University Press.
- Fox, B. 1987. *Discourse Structure and Anaphora. Written and Conversational English*. Cambridge Studies in Linguistics, Cambridge University Press.
- Grosz, B.J. 1981. Focusing and description in natural language dialogues. Joshi, A., Webber, B. and Sag, I. (eds.) *Elements of Discourse Understanding*, Cambridge University Press, England, P. 85-105.
- Grosz, B.J., Joshi, A.K. and Weinstein, S. 1995 Centering: A framework for modeling the local coherence of discourse. *Computational Linguistics*, 12(2), p. 203-225.
- Grosz, B.J. and Sidner, C. 1986. Attention, intentions, and the structure of discourse. *Computational Linguistics*, 12(3), p. 175-204.
- Gundel, J., Hedberg, N. and Zacharski, R. 1993. Cognitive Status and the Form of Referring Expressions. *Language*, 69, P. 274-307.
- Halliday, M.A.K. and Hassan, Ruqaiya. 1976. *Cohesion in English*, Longman, London and New York.
- Hovy, E. 1988. Planning coherent multisentential text. *Proceedings of the 26th Annual Meeting of the Association for Computational Linguistics (ACL-88)*, State University of New York, Buffalo, p. 163-169.

- Ide, N. and Cristea, D. 2000. A Hierarchical Account of Referential Accessibility. *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics, ACL'2000*, Hong Kong, p. 416-424.
- Kameyama, M. 1998. Intrasentential Centering: A Case Study. Prince, E., Joshi, A. and Walker, L. (eds.) *Centering in Discourse*, Oxford University Press, p. 89-112.
- Kintsch, W. and Van Dijk, T.A. 1975. Comment on se rappelled et on résume les histoires, *Langages*, 40.
- Mani, I. 2001. Automatic Summarization. John Benkamin Publishing Company, Amsterdam/Philadelphia.
- Mann, W.C. and Thompson, S.A. 1988. Rhetorical Structure Theory: Toward a Functional Theory of Text Organization. *Text*, 8(3), p. 243–281.
- Marcu, D., 1999. A formal and computational synthesis of Grosz and Sidner's and Mann and Thompson's theories. *Proceedings of the Workshop on Levels of Representation in Discourse*. Edinburgh.
- Marcu, D. 2000. The theory and practice of discourse parsing and summarization, The MIT Press, Cambridge, Massachusetts.
- Miller, G. 1956. The magical number Seven, Plus or Minus Two: Some Limits on Our Capacity for Processing Information, *The Psychological Review*, vol. 63, p. 81-97.
- Moser, M. and Moore, J.D. 1996. Toward a synthesis of two accounts of discourse structure. *Computational Linguistics*, 22(3), p. 409–419.
- Passonneau, R., J. 1995. Integrating gricean and attentional constraints. *Proceedings of IJCAI*.
- Postolache, O. 2001. Sumarizarea textelor. Lucrare de licență. Universitatea „Al.I.Cuza” Iași, Facultatea de Informatică.
- Richadeau, F. 1969. La lisibilité. Langage-Typographie-Signes-Lecture, Paris.
- Schank, R. and Abelson, R. 1977. Scripts, plans, goals and understanding, Hillsdale, N.J.
- Schilder, F. 2001. Robust Discourse Parsing Via Discourse Markers, Topicality and Position. *Natural Language Engineering* 1, (1), p.1-22.
- Scott, D.R., de Souza, C.S. 1990. Getting the message across in RST-based text generation. Dale, R., Mellish, C. and Zock, M. (eds.) *Current Research in Natural Language Generation*, Academic Press, New York.
- Serețan, V. and Cristea, D., 2002. The Use of Referential Constraints in Structuring Discourse, *Proceedings of the Third International Conference on Language Resources and Evaluation, LREC'2002*, Las Palmas, Spain, p.1231-1237.

-
- Sidner, C. 1983. Focusing in the comprehension of definite anaphora. Brady, M. and Berwick, R.C. (eds.) *Computational Models of Discourse*, MIT Press.
- Sofronie, V. 1999. Implementări existente în sumarizarea textelor. SumVT. Lucrare de licență. Universitatea „Al.I.Cuza” Iași, Facultatea de Informatică.
- Strube, M. and Hahn, U. 1996. Functional Centering. *Proceedings of the 34th Annual Meeting of the Association for Computational Linguistics*, Santa Cruz, California., p. 270-277.
- Tablan, M.V., Barbu, C., Popescu, H., Hamza, R.O., Nita, C.I., Bocaniala, C.D., Ciobanu C. and Cristea, D. 1988. Co-operation and Detachment in Discourse Understanding. *Proceedings of the Workshop on Lexical Semantics and Discourse Structure, ESSLI'98*, Saarbruecken.
- Walker, M., Iida, M., Cote, S. 1994. Japanese Discourse and the Process of Centering. *Computational Linguistics*, 20(2), p. 193-232.
- Walker, M.A. 1996. Limited attention and discourse structure. *Computational Linguistics*, 22-2.
- Walker, M.A. 1998. Centering, anaphora resolution, and discourse structure. In Marilyn A. Walker, Aravind K. Joshi, and Ellen F. Prince, editors, *Centering in Discourse*. Oxford University Press.

DLIR - un sistem de cautare documentara multilingv

Amalia TODIRASCU
INRIA Lorraine, LORIA, Campus scientifique BP 239,
54506 Vandoeuvre-lcs-Nancy Cedex, France,
todirasc@loria.fr

Abstract

Aceasta lucrare prezinta un sistem de cautare documentara bilingv francez-roman pentru un domeniu limitat, cel al securitatii computerelor. Cautarea si indexarea se fac pe baza unei ontologii comune domeniului. Identificarea instantelor conceptelor în texte sau în întrebările utilizator se realizeaza cu ajutorul unor tehnici robuste de analiza limbajului natural, combinate cu o ontologie specifica domeniului.

Introducere

Sistemele de cautare de informatii indexeaza o baza de documente, în general pe baza unor liste de cuvinte cheie extrase din documentele respective. Ele primesc întrebările utilizatorului, încearca o mapare a întrebării cu indexul construit care regaseste documentele. Raspunsul sistemului contine un numar de documente care sunt relevante în raport cu întrebarea utilizatorului. Fiecare sistem defineste un criteriu de relevanta specific. Aceste sisteme sunt evaluate pe baza a doi parametri: rapel (numarul de documente regasite/numarul total de documente relevante) si precizie (numarul de documente relevante regasite/numarul de documente regasite). În cazul unui sistem de cautare multilingv, raspunsul la o întrebare poate contine mai multe documente relevante, chiar daca sunt scrise în alte limbi decât cea în care a fost formulata cererea.

Sistemele de cautare de informatii clasice ofera utilizatorului solutii imprecise sau vide. Aceste rezultate se datoreaza utilizarii ca indecsi doar a cuvintelor cheie, extrase pentru fiecare document în parte. Majoritatea sistemelor de cautare ignora fenomene caracteristice limbajului natural: fenomenul de ambiguitate (un cuvânt poate avea mai multe sensuri) sau de polimorfism (un concept poate fi exprimat în mai multe moduri). În plus, un sistem care își propune sa faca cautare într-o baza de date multilingva trebuie sa fie capabil sa gaseasca informatia ceruta în orice document disponibil, indiferent de limba în care a fost scris. Rezolvarea problemelor specifice limbajului natural (ambiguitate, traducere automata) necesita resurse lingvistice importante pentru fiecare limba care este tratata de

catre sistem, daca aplicam tehnicile clasice de analiza limbajului natural. Tehnicile clasice de analiza sintactica nu sunt adaptate sistemelor de cautare documentara, datorita dimensiunilor prea mari ale bazei documentare.

Tehnicile robuste de analiza sintactica, inspirate de domeniul extragerii de informatii (GATE [6], FASTUS) sunt dedicate rezolvarii unor probleme dedicate, precise (identificarea numelor proprii, ale grupurilor nominale simple). Printre acestea, automatele cu stari finite [5], colocatii [9] sau liste de pattern-uri sintactice reprezentind structura sintactica a grupului nominal simplu sunt resursele lingvistice necesare pentru aceste componente. Aceste tehnici au avantajul de a fi robuste, de a putea trata o cantitate importanta de informatii în timp real, precum si de a fi portabile de la un domeniu si/sau o limba la alta.

O alternativa la sistemele de indexare clasice sunt cele care folosesc structuri sintactice sau conceptuale pentru a indexa baza de documente. Acestea nu sunt foarte numeroase, pentru ca pe de o parte, ontologiile generice nu sunt disponibile decât în numar prea restrâns (WordNet \cite{vossen} si Corelex [3] sunt doar doua exemple de resurse libere). Pe de alta parte, textele nespecifice pun probleme analizoarelor existente, datorita faptului ca necesita resurse adaptate domeniului: dictionare, gramatici locale.

Într-o aplicatie de cautare de informatii pe un domeniu restrâns, asteptarile utilizatorului sunt altele decât cele pentru texte ne-specifice. Ne asteptam la o precizie mai buna a sistemelor. Aceasta impune folosirea de tehnici adaptate acestor sisteme, bazate pe existenta unui model redus al domeniului. În acest context, voi prezenta o metodologie de extragere a conceptelor candidat din corpus. Acestea sunt folosite de catre un expert uman pentru a îmbogati o ontologie existenta. De asemenea voi prezenta o metoda de indexare a documentelor pe baza unei ontologii, metoda care modifica metoda clasica de indexare semantica latentă.

2. Ontologii

Notiunea de ontologie este dificil de definit, mai multe puncte de vedere coexista. Pentru a simplifica, o ontologie este un model restrâns al unui domeniu specific, format din multimea claselor de obiecte ce populeaza acest domeniu si relatia lor cu celelalte clase.

Ontologiile reflecta un anumit grad de subiectivitate din partea expertului ce a definit-o. Fiecare expert poate avea o anumita viziune a claselor de obiecte ce trebuie incluse în descrierea ontologiei.

O problema a acestor ontologii este legata de portabilitate. O aplicatie definita pentru un anumit domeniu dat va trebui adaptata unui alt domeniu prin construirea unei ontologii corespunzatoare. Construirea lor manuala este dificila si trebuie tinut cont de posibilele redundante, incoerente, ce pot fi introduse în baza de cunostinte de catre expertul uman care o construiește. În ultimii ani, s-au facut eforturi deosebite pentru a putea reutiliza ontologiile existente: dezvoltarea unor formate de interschimb dedicate

(Knowledge Interchange Format - KIF), a standardelor (Ontology Interface Layer - OIL) [8], dezvoltate în cadrul proiectului Semantic Web (<http://www.semweb.org>).

Pentru a evita problemele legate de formatul în care a fost reprezentată ontologia, au fost propuse mai multe metode de extragere a ontologiilor din corpusuri. Acestea disting mai multe etape:

- identificarea termenilor (posibilele instanțe ale conceptelor exprimate în limbaj natural);
- identificarea relațiilor între termi;
- identificarea relațiilor între termi și concepte.

Majoritatea acestor etape necesită validarea rezultatelor de către un expert uman, sau asignarea unor interpretări (identificarea unei relații între două mulțimi de termi). Metodele statistice interpretează contextele existente și regroupează termii cu același context în clase diferite [1], [7]. Relațiile între termi sunt interpretate pe baza informațiilor de subcategorizare asociate verbelor. Dezavantajul metodelor statistice este acela că necesită corpusuri adnotate de talie importantă pentru a putea învăța.

Metodele bazate pe inferențe logice propun proceduri semi-automate pentru a verifica validitatea cunoașterii existente. Conceptele noi, deduse de către regulile de inferență, sunt adăugate ierarhiei domeniului dacă sunt coerente cu cunoașterea existentă. Relațiile pot fi identificate printre cunoașterea sintactică (subcategorizare [4]). Problemele acestei metodologii sunt supragenerarea de concepte și costul verificării incoerențelor și inconsistențelor cunoașterii sunt principalele neajunsuri ale metodei. Mai multe formalisme au fost dezvoltate în acest scop, și printre acestea logicile terminologice joacă un rol important.

2.1. Logici Terminologice

Logicile terminologice (LT) sunt formalisme de reprezentare a cunoștințelor care sunt derivate din rețelele semantice, dar sintaxa și semantica lor sunt bine definite. Ele combină proprietăți ale sistemelor orientate-obiect, ale sistemelor bazate pe frame-uri și ale logicilor modale.

LT propun o organizare ierarhică a cunoașterii, pe două nivele: unul conceptual (T-Box), care descrie clasele abstracte conținând obiectele relevante pentru modelarea domeniului și un nivel asertional (A-Box), conținând instanțele claselor. Clasele de obiecte (concepte) sunt descrise de relații (numite roluri) cu alte concepte, și cu atributele lor (rolurile cu valori atomice).

2.1.1. Sintaxa și semantica logicilor terminologice

Operatorii LT sunt inspirați de logica de prim ordin:

Operator	Operator Logic	Semantica
----------	----------------	-----------

D = SOME R C	$\exists x (xRC)$	Există cel puțin o instanță a lui C în relație cu R
D = ALL R C	$\forall x (xRC)$	restricționează co-domeniul relației R
D = AND C1 C2	$C1 \wedge C2$	Conjunția de descrieri conceptuale
D = OR C1 C2	$C1 \vee C2$	Disjuncția de descrieri conceptuale
$C1 \subseteq C2$	$C1 \subseteq C2$	Axiom: C1 conține condiții necesare pentru C2
D = NOT C	$\neg C$	complementul conceptului C
D = (n.R.C	$(y1...yn (1 (i(n, R(x, yi)(C(yi))$	Există cel puțin n obiecte de tip C în relația R cu D

Figura 1. Operatori în LT

Folosind toti acesti operatori, sau doar o parte a acestora, mai multe expresivitati sunt posibile: definirea conceptelor si a rolurilor ALC (folosind SOME, ALL, AND, OR, NOT ca operatori, axiomele conceptuale), posibilitatea utilizarii rolurilor tranzitive (R+), a rolurilor inversabile (I), a ierahiilor de roluri (H), a atributelor (f) sau a restrictiilor numerice.

Unele comenzi LT sunt explicate mai jos. CN este un nume de concept, C este o descriere conceptuala (orice combinatie de operatori AND, SOME, NOT, ALL). Comenzile LT sunt inspirate de formalismul KRSS ([2]):

1. (define-concept CN C) - definește un nou concept ca o descriere conceptuala;
2. (instance IN C) - definește o instanța a unui concept dat;
3. (implies C1 C2) - introduce o noua axioma conceptuala, definind condițiile C1 necesare pentru descrierea conceptuala C2;

LT sunt fragmente decidabile ale logicii de prim ordin. Acestea propun algoritmi decidabili pentru verificarea coerenței și consistenței cunoștințelor. LT propune mecanisme logice pentru a identifica subsumarea, regasirea instanțelor, drumurile care unesc mai multe concepte. Clasificarea este o ordonare parțială a ierarhiei de concepte, în raport cu relația de subsumare.

Câteva exemple de comenzi:

(concept-subsumes? C1 C2) testează dacă C1 subsumează C2

(concept-parents C) regăsește strămoșii direcți ai conceptului C

(concept-children C) regăsește fiii direcți ai lui C

(classify-tbox) calculează toate relațiile de subsumare între toate conceptele definite în T-Box

(concept-instances C) regăsește toate instanțele conceptului C

2.2. Logici terminologice pentru sisteme de extragere si de regasire a informatiilor

Rolul cunostintelor din domeniu într-un sistem de extragere a informatiilor este acela de a valida reprezentarea semantica a entitatilor care sunt potential relevante, identificate în text prin tehnici de procesare a limbajului natural. Aceste entitati pot fi folosite pentru a adauga noi concepte la ontologia existenta. Cea mai mare parte a sistemelor de extragere a informatiilor foloseste tehnici NLP robuste pentru identificarea candidatilor si entitatile candidat nu sunt validate de catre o interpretare semantica. Sistemele de extragere a informatiilor pot folosi cunoastere implicita, cum ar fi relatiile de hiponimie/hiperonimie.

Logicile terminologice prezinta avantajul de a lucra cu date semi-structurate sau incomplete. Nu este necesara definirea explicita unor valori ca instante ale unor concepte. Valorile implicite nu sunt utilizate de catre logicile terminologice. Unele valori ale rolurilor sunt lasate nespecificate ca în urmatorul exemplu:

```
(define-concept computer (and physicalobject (some hasOperatingSystem
OSystem) (some hasType Type)))
(define-primitive-concept Type)
(define-primitive-concept OSystem)
(instance sun1 (and computer (some hasType SparcStation)))
```

În acest exemplu, vom ilustra faptul ca definitiile implicite sunt acceptate de catre logicile terminologice (SparcStation nu sunt definite explicit de catre o instanta sau un subconcept al conceptului Type). Nu este definita explicit nici o instanta a rolului **hasOperatingSystem**.

Aceste proprietati nu sunt interesante pentru aplicatia noastra, dar erorile sunt posibile, iar cunoasterea domeniului este incompleta.

Relatiile de hiperonimie sau hiponimie sunt tratate cu ajutorul relatiilor de subsumare între conceptele domeniului. De exemplu, daca un concept candidat este identificat în text ca:

```
(instance x (and PC (and hasOperatingSystem Linux)))
(define-concept PCcomputer (and computer (some hasType PC)))
x este de asemenea o instanta a conceptului computer.
(instance y (and Password (some hasUser Root)))
(define-concept Password (and String (some hasAtr secret) \\ (some hasBelongs
User)))
(define-concept System (some hasUser User))
(define-concept Root User)
```

Pentru aplicatia noastra avem nevoie de o logica terminologica care sa propuna rationament la nivel de instanta, sa permita lucrul în contextul unei lumi deschise, precum si proceduri optimizate de calcul a relatiilor de subsumare sau de clasificare. Printre putinele sisteme care implementeaza rationament la nivel de instanta am ales RACER ([10]), fiind unul dintre cele mai performante.

În sectiunea urmatoare voi prezenta metoda de extragere a termenilor din texte folosind sistemul DLIR [16]. Textele vor fi traduse într-o reprezentare conceptuala unica, permitând regasirea informatiilor în mai multe limbi.

3. Arhitectura

Sistemul DLIR contine mai multe module: un modul de analiza sintactica robusta, un modul de întretinere a ontologiei domeniului, un modul de indexare a documentelor bazat pe celelalte doua module. In cele ce urmeaza voi prezenta aceste module în detaliu.

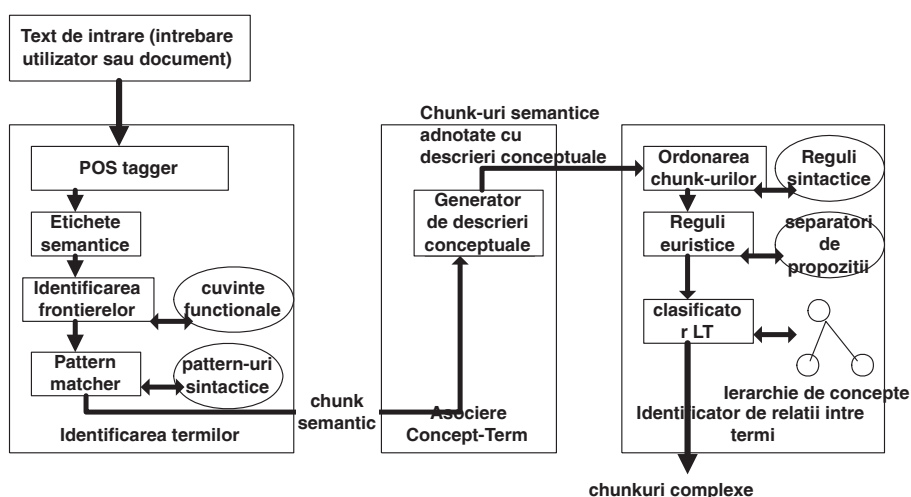


Figura 1: Instanțe ale conceptelor care apar în întrebare

3.1. Analiza sintactica robusta

Acest modul este dedicat identificării termenilor posibili, utilizând tehnici de analiză robuste, o serie de resurse specifice domeniului (o lista de corespondente cuvinte-concepte). Termii sunt combinați, conform unor reguli euristice pentru a crea concepte complexe. Acestea sunt validate ulterior, apelând modulul de acces la ontologia domeniului. Eventualele concepte valide sunt adăugate la ontologia existentă. Acest modul conține mai multe submodule implementate în Java, în Perl și CLIPS (modulul care

aplica regulile de combinare a termilor). Notiunea de chunk semantic a fost propusa pentru a identifica termii candidat [16]. Metoda a fost testata pentru limba franceza, dar cum resursele folosite pentru identificarea termilor sunt relativ independente de limba pentru care a fost construita aplicatia, este posibila extinderea ei si pentru limba româna, dupa cum voi arata mai jos.

3.2. Identificarea chunk-urilor semantice

Scopul principal al acestui modul este acela de a identifica secventele de cuvinte care corespund celor mai semnificative concepte ale domeniului (*chunk-uri semantice*).

Un *chunk semantic* contine un pattern sintactic simplu (grup substantival simplu, grup verbal) si este delimitat de doi separatori de clauze.

Separatorii sunt cuvinte functionale, verbe auxiliare, sau anumite sintagme prepozitionale.

Exemplu. "la victime d'une intrusion inattendue"

[victima unei intruziuni neasteptate]

În acest exemplu, "victima" si "unei intruziuni neasteptate" sunt chunk-uri semantice, care contin informatia relevanta.

Modulul contine mai multe submodule: un POS tagger, un tagger semantic, un identificator de frontiere si un pattern matcher. Identificarea chunk-urilor semantice este bazata pe informatia lexicala, propusa de POS tagger.

3.2.1. Part-Of-Speech tagging

Modulul care este dedicat identificarii partilor de vorbire asociate cuvintelor (folosind WinBrill, antrenat pentru franceza pe baza unui set de date propuse de Institut National pour la Langue Française [11]) identifica cuvintele care au un continut (substantive, adjective, verbe) si cuvintele functionale (prepozitii, conjunctii etc.).

Taggerul Brill foloseste un set de reguli contextuale si lexicale (bazate pe identificarea prefixelor si a sufixelor), învatate pe baza textelor adnotate, pentru a identifica partea de vorbire pentru cuvintele necunoscute.

Pentru limba româna, se foloseste QTAG adaptat pentru limba româna [17], datorita performantelor foarte bune (98% rezultate corecte).

3.2.2. Tagger-ul semantic

Tagger-ul semantic contine un pattern matcher, care consulta un dictionar de talie redusa. Acesta contine o lista cu cele mai frecvente cuvinte si un set de sintagme asociate descrierilor conceptuale corespunzatoare.

Setul de descrieri conceptuale a fost stabilit de catre un expert pe baza unei liste de cuvinte si segmente repetate obtinute dintr-un corpus reprezentativ (200,000 cuvinte). Un

segment repetat este o succesiune de cuvinte care intervin într-un text cel puțin de două ori [14].

Acest modul asociază fiecărui cuvânt conceptul sau descrierea conceptuală din dicționar. Un astfel de dicționar este creat pentru fiecare limbă care este tratată de către sistem.

3.2.3. Modulul pentru identificarea frontierelor

Acest modul identifică separatorii (cuvinte funcționale sau construcții sintactice mai complexe) care delimitează chunk-urile semantice. Acest modul folosește rezultatul POS tagger-ului (care identifică cuvintele funcționale), precum și un set de sintagme (constituenți sintactici care conțin auxiliare, prepoziții compuse). Setul de fraze este construit ca rezultat al studiilor corpusurilor de test pentru franceza și româna (200,000 cuvinte pentru fiecare limbă). Separatorii grupurilor nominale și prepoziționale (determinanți, prepoziții) sunt cei mai buni candidați pentru identificarea separatorilor de chunk-uri semantice; aceștia reprezintă anumite relații potențiale între concepte.

3.2.4. Pattern matcher

Scopul acestui modul este de a identifica nucleul chunk-urilor semantice, nucleu care este reprezentat de un grup nominal simplu sau un grup verbal.

Exemple. Un grup nominal simplu (în franceză) este identificat aplicând următoarele reguli:

N -> NP
 N ADJ -> NP
 Det N -> NP
 Det N ADJ -> NP

3.2.5. DLgen

Acest modul interpretează informația propusă de POS tagger și generează în mod automat o definiție de concept. Un expert trebuie să verifice rezultatele acestui modul. Câteva exemple de reguli propuse pentru generarea descrierilor DL simple (valabile pentru ambele limbi):

- S1/N S2/ADJ este asociat definiției (define-concept S1_S2 (AND S1 (SOME hasAtr "S2")))
- S1/N S2/NNP este asociat definiției (define-concept S1 (SOME hasName "S2"))
- S1/ADJ S2/N este asociat definiției (define-concept S2_S1 (AND S2 (SOME hasAtr "S1")))
- Verbele sunt traduse ca nume de roluri: S1/VB este asociat rolului **hasS1**.

Unele pattern-uri identifica negatiile, chiar daca este imposibil sa enumeram toate posibilitatile si sa detectam corect domeniul negatiei:

- sans/ADV S1/N este asociat definitiei (define-concept not_S1 (NOT S1
- nici_unul/ADV S1/N este asociat definitiei (define-concept not_S1 (NOT S1))

Rezultatele propuse de **DLgen** sunt 61% corecte. Iesirea este validata de un expert folosind clasificatorul LT pentru a verifica definitiile conceptuale obtinute în mod automat.

3.3. Relatii între termi

Acest modul foloseste inferentele LT, ca si regulile de sintaxa, pentru a combina descrierile conceptuale asociate fiecarui chunk semantic. Folosim un criteriu de ordonare al chunk-urilor, precum si reguli de combinare a conceptelor pentru a crea concepte complexe. Descrierile rezultante sunt validate de clasificatorul LT.

3.3.1. Ordonarea chunk-urilor

Modulul interpreteaza ordinea chunk-urilor si pozitia chunk-urilor în propozitie.

Clasificam chunk-urile în doua categorii: **chunk-uri principale** si **chunk-uri secundare**. Chunk-urile principale corespund notiunii de nucleu propuse de catre teoriile lingvistice clasice.

Chunk-urile secundare joaca rolul unui modificador, care adauga informatii suplimentare sensului nucleului. Chunk-urile secundare pot lipsi, dar restul propozitiilor este corect. Aceste exemple de reguli definesc chunk-uri diverse:

- chunk-urile care urmeaza dupa un verb la gerunziu sau un auxiliar plus un verb la participiu sunt *chunk-uri secundare*;
- verbele sunt întotdeauna chunk-uri principale.

Exemplu:

'[Main Les atacs Main] [Main ont commenc\{e} Main] [Second r utilisier
les faux comptes Second]'

'atacurile au început prin utilizarea unor conturi false'

Cele doua chunk-uri principale detectate în exemplul de mai sus sunt primul chunk al propozitiei si verbul principal. Chunk-ul secundar este adnotat astfel pentru ca urmeaza dupa prepozitia **r**.

3.3.2. Reguli euristice

Regulile sunt stabilite de catre expert pe baza unui studiu asupra corpusului reprezentativ pentru fiecare limba. Corpusul a fost adnotat cu categoria lexicala propusa de

POS tagger si adnotat manual cu descrierile conceptuale. Setul de reguli heuristice este stabilit pe baza unei liste de pattern-uri de forma $\langle Chunk1 \rangle ?x / FW \langle Chunk1 \rangle$.

Exemplu de reguli euristice sintactice: daca o prepozitie este un separator între doua chunk-uri semantice si prepozitia asociaza un substantiv cu un modificador, atunci putem combina descrierile conceptuale ale celor doua chunk-uri într-o descriere semantica mai complexa, rolul care leaga conceptele fiind cel de modificador:

```
if (<MainChunk1> <Border> <SecChunk2>)
  and (Noun in MainChunk1)
  and (Modifier in SecChunk2)
  then (and sem(MainChunk1) (some hasModifier sem(SecChunk2)))
```

Fiecare pattern este asociat unui cuvânt tinta care identifica conditiile pentru aplicarea regulilor. Prepozitiile, verbele la modul participiu, sunt câteva exemple de cuvinte asociate regulilor euristice. Un numar de 43 reguli (pentru franceza) si un numar de 21 de reguli pentru româna au fost descrise în CLIPS. Iesirea acestor reguli va fi o serie de *chunk-uri complexe*, ce trebuiesc validate de catre expert, cu ajutorul ontologiei domeniului, care este independenta de limba.

3.4. Indexare semantica}

O posibilitate de indexare a documentelor este aceea de folosi direct concepte drept index si nu cuvinte cheie. O metoda eficienta de indexare o reprezinta indexarea semantica latentă. Aceasta metoda construiesc o matrice document-cuvinte cheie si foloseste tehnici de descompunere a matricilor folosind metoda valorilor proprii. În acest fel se elimina coloanele si liniile care sunt vide. Propunem utilizarea conceptelor care fac parte din ontologie în locul cuvintelor cheie. Este posibil ca într-un sistem de cautare a informatiilor multilingv sa avem diferente între ontologiile dintr-o limba într-alta. Avantajul este ca putem folosi drept index concepte care sunt comune ambelor ontologii. Pentru aplicatia noastra am folosit o ontologie construita manual care contine 54 de concepte si 34 de relatii.

Numarul de concepte este mai redus decât numarul de termi, exploatând în special relatiile între termi.

Elementele matricii contin o pondere $weight(C,i)$ calculata astfel:

$$weight(C,i) = \frac{f(C,i)}{\sum_{j=1}^n f(C,j)}$$

pentru fiecare concept, codificând frecventa instantelor conceptului în document si frecventa instantelor în toate documentele indexate de sistem.

$f(C,i)$ - frecventa conceptului în documentul i ;

Conceptele sunt legate prin rolurile dintre acestea. Frecventa unui concept care este situat în ierarhie foarte sus poate fi compus din suma frecventelor instantelor sale. Instancele conceptelor în LT sunt instancele tuturor subconceptelor și ale instantelor sale directe.

Indexarea documentelor se face aplicând metodele de extragere a termenilor prezentate în secțiunea precedentă, înainte de a exploata sistemul. Se folosesc conceptele ontologiei care a fost construită manual. O serie de concepte mai generale ar putea fi obținute scufundând ontologia specifică domeniului cu WordNet ([16]).

Evaluarea acestui sistem a fost realizată pentru un set restrâns de întrebări (50) numai pentru limba franceză. Rezultatele au fost comparate cu cele furnizate de un sistem care folosește cuvinte-cheie pentru indexare. Pentru 74% din întrebări răspunsurile sistemului (rapel și precizie) au fost comparabile cu cele obținute prin metoda de indexare bazată pe cuvinte-cheie. În celelalte cazuri, răspunsurile au fost mai slabe decât indexarea pe baza de cuvinte-cheie. Ontologia folosită este departe de a fi completă, ceea ce a dus la neidentificarea unor termi.

4. Concluzii și perspective

Articolul prezintă o modalitate de a folosi ontologia unui domeniu pentru cautare de informații bilingvă: franceză și română.

Sistemul integrează tehnici de analiză sintactică robustă pentru extragerea celor mai relevante chunk-uri semantice. Metoda folosește o ontologie a domeniului construită manual. Pentru evaluarea pertinentă a metodelor de indexare pe baza de concepte, ontologia va fi actualizată și extinsă cu ajutorul raționamentelor propuse de logicile terminologice, ca și folosirea cunoștințelor sintactice, folosite pentru extragerea unei reprezentări semantice pentru texte și întrebări. Expertul uman trebuie să intervină pentru a decide dacă conceptele identificate în texte pot fi adăugate ontologiei domeniului.

Referințe bibliografice

- [1] Assadi, H., Bourigault, D., 2000, Analyse syntaxique et statistique pour la construction d'ontologies à partir des textes. In J.Charlet, M.Zacklad, G.Kassel, D.Bourigault (eds.) - Ingénierie des connaissances Evolutions récentes et nouveaux défis, Eyrolles Publishing House, pp. 243-256.
- [2] Baader, F., Hollunder, B., 1991. A Terminological Knowledge Representation Systems with Complete Inference Algorithms, Proceedings of the Workshop on Processing Declarative Knowledge.

-
- [3] Buitelaar, P., 1998. CORELEX: Systematic Polysemy and Under-specification, Ph.D. thesis, Brandeis University, Department of Computer Science
 - [4] Capponi, N., Toussaint, Y., 2000, Interprétation de classes de termes par généralisation de structures prédicat-argument. In J.Charlet, M.Zacklad, G.Kassel, D.Bourigault (eds.), Ingénierie des connaissances - Evolutions récentes et nouveaux défis, Eyrolles Publishing House, pp. 337-356.
 - [5] Chanod J.P., 1999. Natural Language Processing and Digital Libraries. In M.T.Pazienza (ed.), Information Extraction, Springer-Verlag, LNAI 1714, pp.17-31.
 - [6] Cunningham, H., Wilks, Y., Gaizauskas, R.J., 1996. New Methods, Current Trends and Software Infrastructure for NLP. In Proceedings of the conference on New Methods in Natural Language Processing (NeMLaP-2), Bilkent University, Turkey, 1996, pp.1-12.
 - [7] Daille, B., 1996, Study and Implementation of Combined Techniques for Automatic Extraction of Terminology. In J.Klavans, P.Resnik (eds.) - The Balancing Act - Combining Symbolic and Statistical Approaches to Language, MIT Press, pp. 49-66.
 - [8] Fensel D. et al., 2000, OIL in a nutshell. In R. Dieng et al. (eds.), Knowledge Acquisition, Modeling, and Management, Proceedings of the European Knowledge Acquisition Conference (EKAW-2000), Lecture Notes in Artificial Intelligence, LNAI, Springer-Verlag.
 - [9] Heid, U., 2000, A linguistic bootstrapping approach to the extraction of term candidates from German text, Terminology, pp 161-180.
 - [10] Haarslev V., Muller R, 2001, Description of the RACER System and its Applications, Proceedings of the International Workshop on Description Logics (DL-2001), Stanford, USA, 1.-3. August 2001, pp. 132-141
 - [11] Lecomte, J., Le Catégoriseur BRILL14-JL5/WINBRILL-0.3, InaLF, InaLF/ CNRS report, December 1998.
 - [12] Riloff,E., Lorenzen, J., 1999, Extraction-based Text Categorization Generating Domain-Specific Role Relationships Automatically. In ed. T.Strzalkowski, Natural Language Information Retrieval, Kluwer Academic Publishers, pp. 167-196.
 - [13] Riloff, E., Shepherd, J., 1997, A Corpus-Based Approach for Building Semantic Lexicons. In Proceedings of the Second Conference on Empirical Methods in Natural Language Processing.
 - [14] Rousselot, F., Frath, P., Oueslati, R., Extracting concepts and relations from Corpora. In *Proceedings of the Workshop on Corpus-oriented Semantic Analysis, European Conference on Artificial Intelligence, ECAI 96*, Budapest, 12 August 1996.
 - [15] Schimid, H., 1994, Probabilistic Part-of-Speech Tagging Using Decision Trees, Proceedings of the International Conference on New Methods in Language Processing, Manchester, United Kingdom

-
- [16] Todirascu, A., 2001, Semantic Indexing for Information Retrieval Systems, Ph.D. Thesis, University Louis Pasteur of Strasbourg, France, March 2001.
 - [17] Tufiş, D., Mason O., Tagging Romanian Texts: a Case Study for QTAG, a Language Independent Probabilistic Tagger. In *Proceedings of the First International Conference on Language Resources and Evaluation (LREC)*, Granada, Spain, 1998, pp. 589-596.
 - [18] Vilain, M., 1999, Inferential Information Extraction. In M.Pazienza (ed.), *Information Extraction*, LNAI 1714, Springer-Verlag, pp.95-119.
 - [19] P. Vossen, P., *Introduction to EuroWordNet*, Kluwer Academic Publisher, 1998.
 - [20] Zweigenbaum, P., Consortium MENELAS, 1995, MENELAS: Coding and Information Retrieval from Natural Language Patient Discharge Summaries. In M.-F.Laires, M.J.Ladeira, J.-P. Christensen (eds.) - *Advances in Health Telematics*, IOS Press, Amsterdam, pp.82-89.

Mediu hermenofor pentru asistarea învățării unor concepte dintr-o limba străină

Stefan TRAUSAN-MATU
Universitatea "Politehnica" Bucuresti, Facultatea de Automatica si Calculatoare,
Centrul de Cercetari Avansate în Învățare Automata,
Prelucrarea Limbajului Natural si Modelare Conceptuala al Academiei Române
email: trausan@cs.pub.ro, trausan@racai.ro
URL: www.racai.ro/~trausan

1. WWW, o prezenta din ce în ce mai comuna

În mai puțin de zece ani, rețeaua globală de documente World Wide Web (WWW sau, pe scurt web), a devenit omniprezentă și este posibil ca într-un timp nu prea lung să înlocuiască o mare parte din cărți, televizorul, cinematograful, ziarele și revistele (toate acestea fiind deja disponibile pe web) și, în plus să furnizeze chiar posibilitatea imersiunii în realități virtuale. Un singur exemplu cred că este suficient: anul trecut rezultatele bacalaureatului au fost publicate pe web.

WWW a atins deja dimensiuni comparabile cu imensa Bibliotecă a Congresului SUA. Extinderea sa este datorată ușurii cu care poate fi parcurs de către oricine are un calculator, pe de o parte, și de simplitatea cu care oricine poate publica ceva pe el. Pe de altă parte, costul accesului la resursele web este de cele mai multe ori infim.

WWW este un hipertext extins la scară întregului glob prin rețeaua mondială de calculatoare Internet. Pe fiecare calculator pot fi plasate unul sau mai multe documente care constituie noduri (pagini) în hipertext. Oriunde într-o astfel de pagină poate exista o legătură la o altă pagină, de pe același sau de pe alt calculator, în acest mod putând fi unite informații aflate în locuri diferite. O nouă pagină pentru web poate fi creată ușor chiar de utilizatori nu neapărat profesioniști în informatică, în acest scop existând mai multe editoare de texte specializate.

Termenul de hipertext se pare că provine de la termenul de spațiu hiperbolic sau hiperspațiu, apărut în 1704 și folosit de matematicianul F. Klein pentru geometria cu mai multe dimensiuni [Rad91]. Din această perspectivă, un hipertext este un text cu mai multe dimensiuni explicite (față de doar o dimensiune, în cazul textului liniar). De fapt, orice text are implicit mai multe dimensiuni, deoarece, chiar dacă forma de prezentare a unui text este liniară, pe hârtie, în el există o structură implicită, data de discurs. De asemenea, există conexiuni implicite, subiective între părți ale textului, concepte legate între ele, Hipertextul

este o organizare a unui text în care toate aceste legaturi sunt explicitate și pot fi exploatate în parcurgerea făcută pe un calculator.

În jurul anului 1962, Douglas Engelbart a dezvoltat primul sistem hipertext, prezentat atunci drept o arhitectura conceptuală destinată creșterii potențialului intelectului uman ("Conceptual Framework for Augmenting Human Intellect") [Eng95]. Sistemul era destinat manipulării de concepte structurate într-o rețea în care arcele sunt relațiile între concepte.

Primul sistem declarat ca fiind hipertext a fost creat de Theodor Nelson în 1967 sub numele de "Xanadu". Nelson își propunea atunci să dezvolte un sistem, masiv paralel, destinat muncii creative și studiului. El a plecat în îndeplinirea acestei idei de la dorința de a găsi cea mai bună abstracție care să unifice literatura și arta cinematografică.

Sistemele hipertext (hipermedia) permit accesul personalizat la volume imense de informații. În același timp, însă, ele suferă de problema aflului de informație cu care este bombardat un utilizator. O soluție este dezvoltarea de instrumente, aplicații, medii informatice pentru facilitarea accesului la cunoștințele dorite pe web. Aceste instrumente trebuie să faciliteze înțelegerea, abstractizarea textelor, extragerea informațiilor utile. Acesta este unul din motivele pentru care le-am denumit instrumente hermenofore. Trebuie remarcat faptul că ideea de a considera hipertextele ca instrumente de sprijinire a activităților cognitive a stat chiar la baza concepției acestora, după cum am precizat în paragrafele anterioare.

În continuare, după o trecere în revista a problematicii ontologiilor, în secțiunea următoare se va introduce conceptul de mediu hermenofor, se va justifica necesitatea acestuia și se vor prezenta caracteristicile acestora. Lucrarea va fi încheiată cu o exemplificare de sistem care are câteva trăsături ale unui mediu hermenofor și de o secțiune de concluzii.

2. Ontologii

Termenul de "ontologie" a fost, până nu de mult, folosit exclusiv în filosofie, pentru a denumi teoria asupra existenței, mai corect spus, asupra ceea ce consideră că există cel care întocmește teoria. Construirea multor sisteme filosofice pleacă de la o ontologie, adică de la definirea categoriilor fundamentale de entități din realitate și a relațiilor dintre ele. Chiar dacă ontologia nu este întotdeauna explicită, orice demers conceptual construiește o ontologie, chiar implicit, inconstent.

În ultimii ani, termenul de ontologie este folosit și în știința calculatoarelor. Cea mai frecventă extindere a folosirii acestui concept este în cadrul sistemelor de inteligență artificială bazate pe cunoștințe. Majoritatea programelor de calculator cu inteligență artificială prelucrează structuri de simboluri, care sunt menite să reprezinte conceptele, cunoștințele referitoare la domeniul considerat. Aceste structuri simbolice sunt grupate într-o așa numită bază de cunoștințe care constituie, de fapt, un model al domeniului respectiv.

În ultimii ani se considera ca aceasta baza de cunostinte trebuie vazuta ca o ontologie, o conceptualizare, o teorie asupra ceea ce exista în domeniul avut în vedere. O ontologie este, din aceasta perspectiva, o “specificare a unei conceptualizari ... Termenul este împrumutat din filosofie, unde însemna o considerare sistematica a existentei. În inteligenta artificiala se refera la precizarea a ceea ce se considera ca <<exista>>” [Gru96].

Între concepte pot exista diverse relatii. Cea mai importanta relatie este probabil cea hiperonimica [WN], taxonomică, între un concept si unul sau mai multe concepte mai generale, din care deriveaza, care îl subsumeaza, din a caror combinatie a fost generat. Prin aceasta relatie se pot “mosteni” proprietati de la conceptul (conceptele) mai general(e) la cel mai particular, daca aceste proprietati nu sunt redefinite la conceptul din urma. Alte relatii sunt cea meronimica [WN] (“parte-întreg”), între un concept si partile sale sau cea antonimica, între doua concepte (adjective) opuse.

O ontologie include, asadar:

- categoriile, conceptele fundamentale,
- proprietatile conceptelor,
- relatiile si distinctiile între concepte.

O ontologie este rezultatul unei experientieri, a unor experiente traite, în care sunt evidentiaste niste constante, niste regularitati, care ne îndreptatesc sa afirmam ca vor fi regasite în viitor. În urma investigatiei facuta pentru a gasi esenta regularitatilor se delimiteaza entitati mentale denumite concepte sau categorii, care pot fi diferentiate de alte categorii. Aceste entitati pot intra în combinatie cu altele formând noi concepte.

Un aspect deosebit de important în ceea ce priveste rolul ontologiile este faptul ca ele exprima o comuniune, (co)existenta unei diversitati de concepte, cu diferente si relatii între ele. O presupozitie este ca exista doar un numar limitat de concepte sau categorii, ceea ce înseamna ca se poate face un fel de cuantificare, de discretizare a realitatii. Acestea constituie un punct de sprijin pentru achizitia de noi concepte sau pentru rationamentele facute de om sau de calculator.

Partajarea unei ontologii este esentiala în sistemele bazate pe agenti (programe) inteligenti pentru, de exemplu, comertul electronic, pentru a le asigura autonomia, flexibilitatea si agilitatea. Ontologiile sunt liantul care integreaza sisteme de baze de date, sisteme de obiecte, sisteme bazate pe cunostinte, în diverse aplicatii integratoare si bazate pe colaborare. Ele reduc ambiguitatile semantice în partajarea si reutilizarea cunostintelor. “Scopul suprem este dezvoltarea de ontologii reutilizabile care pot fi aplicate pentru mai multe discipline”. [OORG]

“O ontologie are drept prim scop facilitarea comunicarii între calculatoare, independent de tehnologiile unui anumit sistem individual, arhitectura de prelucrare a informatiilor si domeniul aplicatiei. Ingredientii cheie care constituie o ontologie sunt un vocabular de termeni de baza si o specificare precisa a ceea ce înseamna acesti termeni.” [OORG] O ontologie este însa mai mult decât un vocabular. Ea este punctul de plecare

pentru dezvoltarea de structuri de cunostinte, nu numai taxonomii sau clasificari de concepte ci si relatii complexe. [OORG]

Din punct de vedere al programelor de calculator care folosesc ontologiile, exista doua tipuri de ontologii. Primul tip este cel al ontologiilor destinate sistemelor bazate pe cunostinte, de exemplu, al unui sistem de diagnostic medical. Aceste ontologii sunt caracterizate de un numar relativ redus de concepte, dar legate între ele printr-un numar mare si variat de relatii. Conceptele sunt grupate în scheme conceptuale complexe sau scenarii. Pentru fiecare concept pot exista una sau mai multe particularizari.

Spre deosebire de primul tip de ontologii, ontologiile lexicalizate includ un numar foarte mare de concepte, legate printr-un numar redus de tipuri de relatii (de exemplu, hiperonimica, meronimica etc.). Conceptele sunt reprezentate, de exemplu în WordNet [WN], prin multimi de cuvinte sinonime. Astfel de ontologii sunt folosite în sistemele de prelucrare a limbajului uman.

Corespondenta ontologiei WordNet (care este conceputa pentru limba engleza-americană) pentru limbile europene este EuroWordNet. Aceasta din urma aduce avantajul ca, fiind dezvoltata pentru mai multe limbi (engleza, franceza, germana, italiana, olandeza etc.), permite si dezvoltarea de aplicatii multilingve. În prezent, în cadrul Centrului de Cercetari Avansate în Învatare Automata, Prelucrarea Limbajului Natural si Modelare Conceptuala al Academiei Române este în desfasurare, în colaborare cu mai multe tari din regiunea balcanica proiectul BalkanNet pentru integrarea în EuroWordNet a limbilor din zona, inclusiv a limbii române.

3. Medii hermenofore

Denumim mediu hermenofor o colectie integrata de instrumente (pe care le vom numi hermenofore) si aplicatii informatice directionate catre facilitarea unor activitati de tip hermeneutic ale unui utilizator care exploreaza resurse aflate pe web. Termenul "hermenofor" [Tra01] poate fi parafrizat prin "generator de hermeneutica", pentru a sugera faptul ca un mediu hermenofor faciliteaza activitati hermeneutice, care acorda un rol important experientierii si sunt orientate spre descoperirea unor înțelesuri, a unor structuri profunde, greu detectabile.

Elaborarea de medii hermenofore este absolut necesara în contextul actual al exploziei numarului si volumului de resurse si a interconexiunilor între acestea pe web. Sistemele hipertext (hipermedia) aduc noi dimensiuni cum ar fi interactivitatea, posibilitatile cu totul remarcabile de vizualizare, accesul personalizat la volume imense de informatii. În acelasi timp, însa, ele introduc si unele probleme datorate afluxului de informatie, care poate duce la depasirea capacitatilor cognitive ale utilizatorului, la dezorientare si chiar la alienare. Este un fapt ca utilizatorul, chiar profesionist în informatica, poate fi dezorientat în "labirintul" de pagini de web si resurse de tot felul (baze de date, documente, imagini, ontologii, lexicoane etc.) interconectate.

O soluție la problemele enumerate mai sus este dezvoltarea de instrumente, aplicații, medii informatice pentru facilitarea accesului la cunoștințele dorite pe web. Se poate spune, din această perspectivă, că browserele de web, “motoarele de căutare”, agenții (asistenții) software sunt rudimente de medii hermenofore. Justificarea necesității considerării perspectivei hermenofore este lipsa abilităților hermeneutice ale acestor aplicații. Un exemplu tipic este faptul că “motoarele de căutare pe web” (de exemplu Google [Goo]) furnizează mii sau chiar zeci de mii de documente ca răspuns la o cerere. Alt exemplu este limita actuală a programelor de calculator în înțelegerea textelor cu scopul traducerii, sumarizării sau extragerii cunoștințelor. Aceste probleme sunt datorate, în primul rând, problemelor generate de ambiguitatea limbajului natural, a aspectelor legate de semantică, de pragmatică, de interpretare, de considerarea contextului, a metaforelor, a cunoștințelor de “bun simț”. Toate aceste probleme sunt recunoscute ca fiind “nodul gordian” al aplicațiilor de inteligență artificială. După cum remarcă Terry Winograd, programele de inteligență artificială nu pot depăși condiția unui birocrat, care nu poate să acționeze când nu are “reguli”, care nu se implică [Win87]. Putem spune că, de fapt, problema este că acestor aplicații le lipsesc abilitățile hermeneutice. Ideea noastră este de a oferi un cadru în care puterea oferită de tehnologia informației să fie integrată cu capacitățile specifice umane.

Hermeneutica este, după opinia lui P. Ricoeur, o abordare complementară celei structuraliste în analiza limbajului, a înțeleșului și simbolismului cultural. “Hermeneutica bazează înțelegerea textelor pe intențiile și istoria autorilor și relevanța acestor fapte pentru cititori. În contrast, filosofia analitică identifică de obicei înțeleșul cu referenți externi pentru texte iar structuralismul găsește înțeleșul în aranjarea cuvintelor. Hermeneutica privește textele ca mijloace pentru a transmite experiența, crezurile și judecățile de la un subiect sau comunitate către alții. Astfel, determinarea înțeleșurilor este o problemă de judecată practică și raționament de <<bun simț>> și nu privitor la o teorie a priori sau o demonstrație științifică.” [MHD].

Hermeneutica este studiul interpretării, inițial ea referindu-se doar la interpretarea textelor [MHD]. În prezent s-a extins accepțiunea termenului hermeneutica, vorbindu-se de o poziție hermeneutică în filosofie, care include pe Heidegger, Gadamer, Habermas și Ricoeur, deosebită de formalisti (filosofia analitică, neo-pozitivism sau pozitivismul logic), reprezentată prin Descartes, Leibniz și Russell [Wes97]. Distincția între cele două abordări pleacă de la problema capturării înțeleșului. Pe când formalistii pretind că pot reprezenta înțeleșul, semantică, doar prin identificarea unui denotat în lumea reală corespunzător unei expresii formale, adepții hermeneuticii neagă această posibilitate, pentru ei înțeleșul implicând și considerarea experienței, a credințelor subiectului. Se poate spune că, dintr-un punct de vedere se ajunge la aceeași dispută dintre Husserl și Heidegger sau dintre Dennett și Chalmers.

Mediile hermenofore furnizează informațiile dorite dintr-o perspectivă particulară, pentru un anumit utilizator, considerând un anumit domeniu și într-un anumit moment dat. Un mediu hermenofor trebuie conceput deci în scopul personalizării interfațării la resursele

web-ului, pentru a facilita înțelegerea. Dacă prezentările făcute într-un mediu hermenofor sunt structurate ca hipermedia, una din preocupările principale ce trebuie avute în vedere este faptul ca utilizatorul trebuie sa experimenteze parcurgerea unei secvente de pagini de web, secventa care trebuie sa respecte niste reguli de pragmatica.

În plus fata de furnizarea unei interfete adaptabile, o alta caracteristica a unui mediu hermenofor trebuie sa fie facilitarea initiativei utilizatorului. El trebuie sa poata experimenta, sa poata investiga resursele web-ului. Instrumentele hermenofore sunt destinate sprijinirii activitatii hermeneutice umane adica a unei atitudini directionate catre înțelegerea unor cunostinte sau structuri ascunse în texte (hipertextelor, hipermedia). Un rol important în procesul înțelegerii îl au modalitatile de a genera experientieri, adica experiente de traire, fapte de viata (conform teoriei ca înțelegerea necesita un proces empatic [Wri95], [Mar97]). Unul dintre cele mai uzitate mijloace de acest gen este folosirea metaforelor [LaJ80], [Tra00]. În acest sens se înscrie preocuparea de a dezvolta instrumente (hermenofore) pentru detectarea, adnotarea si prelucrarea metaforelor.

O caracteristica pe care o consideram esentiala la un mediu hermenofor, în contextul precizat mai sus, este si posibilitatea de vizualizare multipla, din perspective diferite, a aceluasi document. Enumeram aici, drept exemplu, în afara perspectivei continutului "brut" al unui document, alte perspective, date de concordante, adnotari (cu parti de vorbire, de exemplu), extrase, rezumate, arbori de analiza semantica, structuri care reprezinta continutul semantic. Remarcam, în acest context, rolul extraordinar de important al adnotarilor documentelor în limbajul extrem de versatil care este XML [XML].

Vom considera ca instrumentele hermenofore au ca scop revelarea si valorizarea unor cunostinte sau a unor structuri încorporate în volumele imense de hipertexte si hipermedia de pe web. Datorita faptului ca abordarea hermeneutica pune pe prim plan rolul experientierii umane, un instrument hermenofor trebuie neaparat considerat în relatie cu utilizatorul care îl foloseste. De aceea, el trebuie sa aiba asociat modelul utilizatorului, care sa contina cel putin urmatoarele informatii despre utilizator:

- ontologia sa,
- scopurile urmarite,
- profilul psihologic,
- istoricul actiunilor efectuate,
- preferintele sale (explicite sau implicite, derivate din observarea comportamentului sau).

Pe de alta parte, instrumentele hermenofore trebuie sa considere si aspectele legate de particularitatile autorilor documentelor :

- ontologiile considerate (de exemplu, ontologiile impuse de paradigmele sau de practicile domeniilor considerate),
- scopurile presupuse,

- elemente de istoric,
- aspecte psihologice general umane.

Instrumentele hermenofore pot fi împartite în mai multe clase, în funcție de acțiunile efectuate:

- cautare a documentelor relevante,
- categorizare a documentelor conform unei taxonomii predefinite,
- relevare de regularități (de exemplu, colocatii) sau structuri în documente,
- segmentarea textelor,
- extragere de informații sau cunoștințe din documente,
- sumarizare,
- relevare de structuri pe web [WSD97],
- instrumente de adnotare (la nivel sintactic, semantic sau pragmatic) a documentelor.

Spre deosebire de instrumentele de minerit al textelor (“text mining”), instrumentele hermenofore pun, în plus, accentul pe aspectele legate de istoricul interacțiunii, de experiența utilizatorului.

În secțiunea următoare se va prezenta sistemul GenWeb de instruire asistată a învățării terminologiei financiare într-o limbă străină [TMC02], [ABK02] care a fost dezvoltat ca un modul într-un proiect mai mare, denumit „Larflast” și finanțat de Comunitatea Europeană. GenWeb a implementat instrumente hermenofore care identifică și utilizează metafore pentru a facilita înțelegerea unui anumit concept [Tra00]. În acest scop, el caută metafore în texte considerate relevante. Metaforele sunt identificate printr-o pereche de cuvinte care corespund la concepte din ontologia domeniului considerat (finanțe) și din ontologia metaforelor, aceasta din urmă reflectând aspecte psihologice general umane [LaJ80]. Trecerea de la un concept la o mulțime de cuvinte (sinonime sau înrudite) se face pe baza ontologiei WordNet, derivată din investigații psiholingvistice [WN]. Metaforele sunt adnotate în XML [XML], unul din atributele folosite în adnotare fiind scopul urmărit de autor [Tra00].

Tot în GenWeb, textele adnotate cu metafore sunt folosite ulterior pentru a genera structuri (bazate pe principii retorice) de pagini de web personalizate conform modelului utilizatorului. Aceste structuri se constituie într-un sit în care cel care învață poate experimenta. Tot pe post de instrumente hermenofore, în GenWeb este disponibilă vizualizarea de concordante în context.

5. Sistem de instruire asistata cu calculatorul în înțelegerea unor termeni financiari

Exista mai multe puncte de vedere asupra modului cum are loc un proces de învățare. Suntem de partea abordării constructiviste [BIM96, Tra97, TNA98, Wil96] în conceperea proceselor educationale. Aceasta abordare considera ca fiecare dintre noi ne construim propria realitate, propriul bagaj de cunostinte, plecând de la experientele pe care le-am avut [ErK97]. Dupa cum remarca [BIM96], “Nucleul studiului este activitatea hermeneutica a constructiei de interpretari.” Învatarea poate fi si ea vazuta constructivist ca un proces hermeneutic, de înțelegere, de transpunere în domeniul studiat, de experimentare, de traire.

Plecând de la ideile învățării constructiviste se ajunge la urmatoarele principii [ErK97]:

- învățarea este un proces activ în care studentii experimenteaza, cauta sa înțeleaga singuri ceea ce învata, profesorul fiind mai mult un îndrumator;
- învățarea trebuie sa fie un proces auto-reglat de catre studenti;
- învățarea constructiva este un proces situational în sensul ca studentul trebuie introdus într-un mediu de învățare care îi permite sa experimenteze, în care se pot face simulari;
- învățarea trebuie sa fie sociala, trebuie sa existe o permanenta colaborare a studentului cu colegii lui.

Dintr-o alta perspectiva, învățarea poate fi considerata ca un proces de inducere de modele mentale adecvate [JoL83]. Înțelegerea poate fi vazuta astfel ca momentul în care realitatea supusa comprehensiunii este pusa în corespondenta cu un model mental complet si valid. Empatia, identificarea eu-lui cu starea de lucruri considerata poate fi, în acest caz, tocmai sentimentul de “traire” în lumea modelului mental.

O practica deja raspândita este de a dezvolta sisteme inteligente de asistare cu calculatorul a instruirii (“Intelligent Tutoring Systems”) care încearca sa monitorizeze procesul de învățare prin verificarea asimilarii conceptelor din ontologia domeniului considerat [Tra95]. Se considera ca un model adecvat al cunostintelor elevului poate fi construit prin raportare la aceasta ontologie. De fapt, aceasta metoda este folosita si în învățământul traditional: noii termeni sunt introdusi prin genul proxim si diferenta specifica. În termenii ontologiilor, noii termeni sunt definiti prin superconceptele care-i subsumeaza si prin particularitatile care-i diferentiaza.

Orice profesor stie însa ca astfel de definitii sunt necesare dar nu sunt suficiente. Pentru a aprofunda termenii definiti sunt necesare exemple, imagini cu un grad mai mare sau mai mic de iconicitate, plecând de la poze si schite, diagrame si grafice, pâna la imagini sugerate, pâna la metafore. Acest fapt este prezent nu numai în învățământ, el apare în orice proces de comunicare (învățământul fiind, bineînțeles, si el inclus).

În cele ce urmează nu ne vom referi la utilizarea imaginilor propriu-zise, care facilitează evident învățarea sau comunicarea. Vom considera un caz particular de imagini, mentale, sugerate, semne iconice lipsite de caracterul vizual dar care comunică o experiență (de multe ori chiar mai puternic, printr-un efect care ar putea face să ne gândim la percepția subliminală). Este cazul metaforelor, care sunt folosite într-o proporție de cele mai multe ori nebanuit de mare în comunicarea inter-umană.

Pentru a ilustra puterea de expresie a metaforelor și, bineînțeles, rolul lor în înțelegerea unor termeni, am să exemplific prin metafora “acțiunile la bursă sunt niște creaturi foarte sensibile” (gasită într-un text pe site-ul de web al Bursei din New York - <http://www.nyse.com>). Nu este nevoie să ne imaginăm o anumită creatură concretă pentru a înțelege ce sugerează metafora exemplificată. Succesul unei metafore, puterea ei expresivă, capacitatea de comunicare sunt date de măsura în care “rezonăm” la mesajul transmis. Ori ce este mai percutant pentru un om decât faptul că suntem creaturi extrem de sensibile? Prin urmare, succesul metaforei folosite într-un context foarte pragmatic, al discursului unui specialist în finanțe este determinat de inspirația vorbitorului de a se referi la un fapt general uman. Nici o definiție de tip gen proxim-diferență specifică nu poate comunica experiența referitoare la aspectul foarte fragil al acțiunilor la bursă precum o face metafora de mai sus.

Rolul covârșitor al metaforelor în viața noastră a fost remarcat și de Lucian Blaga (“omul este un animal metaforic” [Bla85]) și a fost foarte bine evidențiat de Lakoff și Johnson într-o lucrare cu un puternic impact (“Metaforele cu care trăim” - “Metaphors we live by” [LaJ80]). Cei doi autori americani consideră că “subcategorizarea și metaforele sunt două extremități ale unei continuu”, că metaforele “formează sisteme coerente în care ne conceptualizăm experiența” [LaJ80]. Putem spune deci că metaforele oferă alte mijloace expresive decât cele de categorizare oferite de ontologii. Ele nu țin de logica lui Ares, care categorizează, ci de logica lui Hermes, propusă de Noica [Noi86].

Dintr-o altă perspectivă, metaforele pot fi considerate instrumente empatice, care determină imersiunea cititorului (receptorului) în lumea experiențelor autorului. Acest fapt era evidențiat și de Lakoff și Johnson: “Esența metaforei este înțelegerea și experiențierea unui lucru prin altul” [LaJ80]. De exemplu, metafora amintită mai sus despre acțiunile la bursă ne comunică o informație pe care orice ființă vie o înțelege (sensibilitatea, perisabilitatea) dar care nu poate fi exprimată în categorizări.

Importanța metaforelor a fost revelată și de studiul preliminar făcut în cadrul proiectului Larflast (care a avut drept scop elaborarea unui sistem de asistare cu calculatorul a învățării terminologiei financiare într-o limbă străină [Lar], [TMC02], [ABK02]) de o profesoară de limbă engleză la o facultate economică din Sofia. Dânsa remarcă că o importantă dificultate “înțelegerea metaforelor. Limbajul economic și financiar este extrem de metaforic și, uneori, grupuri de metafore apar în imagini complexe. Deseori cuvinte uzuale sunt folosite în metafore elaborate, ... cum ar fi <<a sustine o pierdere>>” [Vit99].

Proiectul Larflast a inclus mai multe module tipice pentru sisteme inteligente de instruire, cum ar fi o ontologie, un mecanism de inferență, teste (grila) pentru diagnosticarea cunoștințelor elevului și actualizarea modelului acestuia. Sistemul dezvoltat include cinci servere de web, unul la București și altele la Leeds, Manchester, Montpellier și Sofia. Serverul de la București, după ce este lansat, accesează serverul de la Sofia pentru a prelua modelul elevului (ce concepte știe și ce concepte nu) și apoi generează pagini de web personalizate.

Metaforele sunt identificate în texte considerate relevante care au fost obținute în urma căutării cu o mașină de căutare uzuală (de exemplu, Google [Goo]). Textele găsite sunt grupate într-un corpus care este adnotat cu metaforele identificate. Acest corpus, împreună cu ontologia domeniului și cu modelul studentului (construit pe baza răspunsurilor date de student la teste) sunt folosite pentru generarea personalizată de pagini de web. În figura următoare este ilustrată arhitectura sistemului GenWeb.

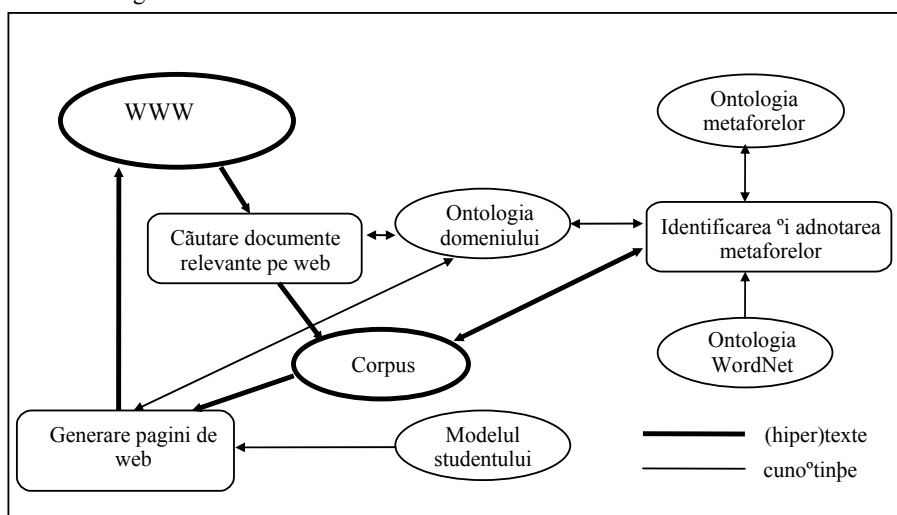


Figura 1

Pentru identificarea și adnotarea metaforelor a fost implementat un editor semantic specializat (fig.2) și un editor de concepte (fig.3).

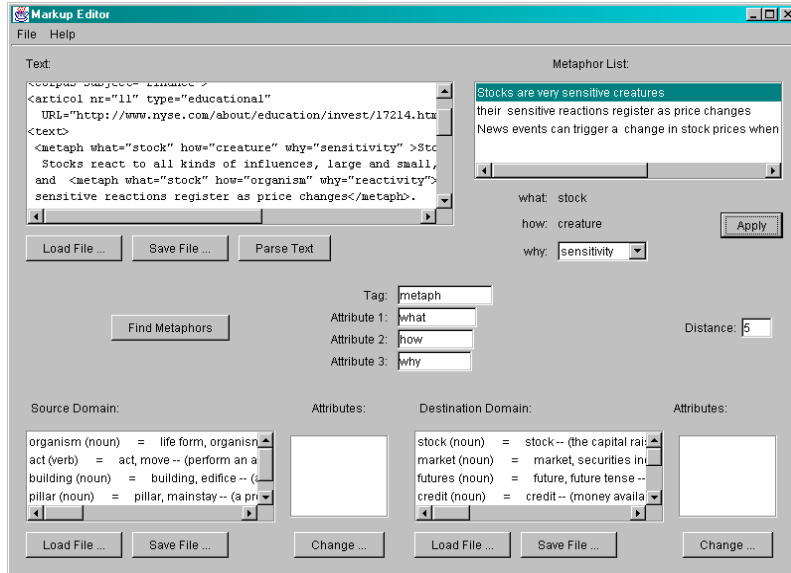


Figura 2

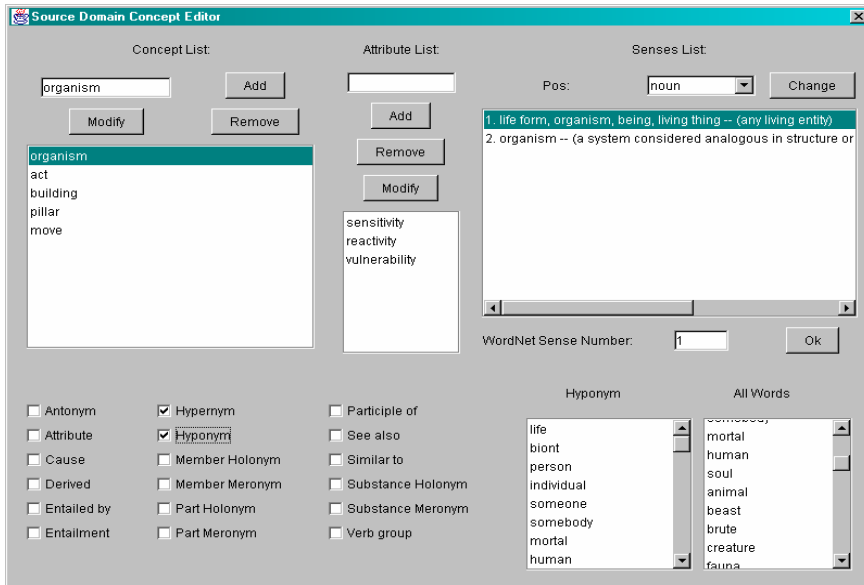


Figura 3

Modelul studentului este creat pe baza raspunsurilor la teste:

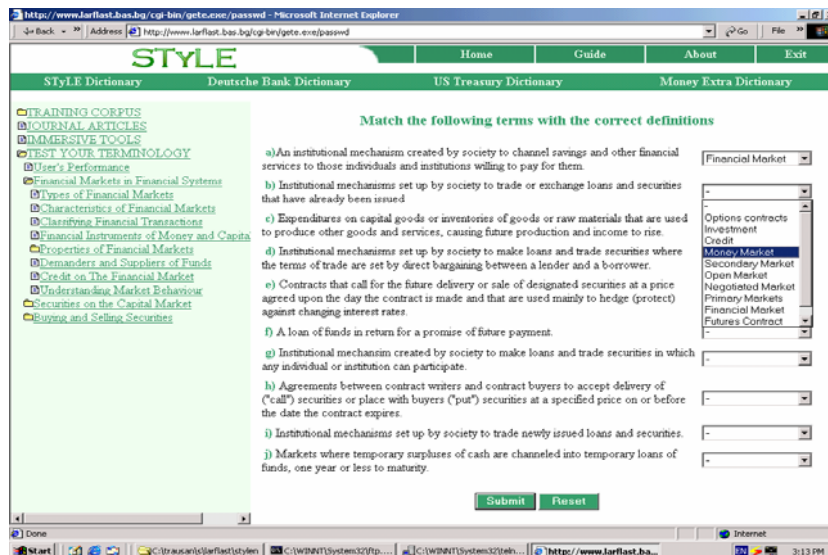


Figura 4

Paginile de web generate dinamic de modulul GenWeb, contributia româneasca la proiectul Larflast, se constituie în structuri care reflecta structura conceptuala (ontologia) a domeniului considerat. Parcurgerea acestora poate fi vazuta si în corespondenta cu facilitatile oferite de o Arta a memoriei [Cu194]. Din alta perspectiva, structurile trebuie concepute în ideea unei retorici specifice paginilor de web [Cl195], [THH95].

Sunt mai multe tipuri de pagini de web generate:

- pagini de diagnostic,
- pagini care definesc concepte, dau exemple de metafore si care includ structuri de paggini web care reflecta ontologia domeniului,
- pagini din structurile de mai sus,
- pagini cu concordante în context.

Aceste pagini sunt ilustrate în figurile urmatoare.

http://www.larflast.bas.bg/cgi-bin/gete.exe/passwd - Microsoft Internet Explorer

Back Address http://www.larflast.bas.bg/cgi-bin/gete.exe/passwd Go File

STYLE Home Guide About Exit

STyLE Dictionary Deutsche Bank Dictionary US Treasury Dictionary Money Extra Dictionary

[TRAINING CORPUS](#)
[JOURNAL ARTICLES](#)
[IMMERSIVE TOOLS](#)
[TEST YOUR TERMINOLOGY](#)

Diagnostics




--- Generated for Trausan ---

Trausan, you have correctly answered to some questions about **financial_market, secondary_market, futures_contract, option_contract, primary_market, investment, credit**, but it seems that you still do not correctly know the following concept(s):

- [Credit](#)
- [Futures contract](#)
- [Investment](#)
- [Primary market](#)
- [Option contract](#)
- [Money market](#)
- [Open market](#)
- [Negotiated market](#)

Please browse the web pages describing these concept(s).

Only the wrongly known and unknown concepts are detailed presented!

 [Back to main page!](#)
 [About LarFlast](#) Please send questions and remarks at trausan@yahhalla.racai.ro


Done Internet

Start C:\trausan\style\style C:\WINNT\System32\teh... Microsoft PowerPoint - [a... http://www.larflast.ba... 3:19 PM

http://www.larflast.bas.bg/cgi-bin/gete.exe/passwd - Microsoft Internet Explorer

Back Address http://www.larflast.bas.bg/cgi-bin/gete.exe/passwd Go File

STYLE Home Guide About Exit

STyLE Dictionary Deutsche Bank Dictionary US Treasury Dictionary Money Extra Dictionary

[TRAINING CORPUS](#)
[JOURNAL ARTICLES](#)
[IMMERSIVE TOOLS](#)
[TEST YOUR TERMINOLOGY](#)

t of New Taiwan dollar credit extended to each customer. In September 1991 amount of New Taiwan dollar credit extended to each customer were established. e foreign currency letters of credit, provide foreign currency guarantees, under full license domestic banks, credit cooperatives, and bill finance companies. T s vulnerability lies not in credit growth or reliance on external funding, but he central bank abandoned the credit allocation ceilings imposed on these banks, nally, China does not have a credit system. It remains a cash-based society, eve te banks continue to allocate credit base on the central plan and extend loans t ers from an over-extension of credit, a bank-dominated financial system and weak accounting standards, greater credit management processes, and greater avenues f insolvency Law Improve credit management Harden budget constraint tralization, but monetary and credit policy still implemented through a implemented through a credit plan. Establishment of special economic zone ents. Established first credit agency. SEC granted license to nine foreign mercial banks Abandoned credit allocation ceilings ;replaced with standard

Done Internet

Start C:\trausan\style\style C:\WINNT\System32\teh... Microsoft PowerPoint - [a... http://www.larflast.ba... 3:20 PM

The image shows three overlapping Microsoft Internet Explorer windows. The top-left window, titled 'FINANCIAL_MARKET', displays a definition of a financial market and a list of metaphorical phrases. The top-right window, titled 'DIAGNOSTICS', lists concepts that the user has not correctly understood, including 'Open market', 'Investment', 'Financial market', 'Secondary market', 'Options contracts', 'Cash', 'Futures contracts', and 'Negotiated market'. The bottom window, titled 'SECONDARY_MARKET', provides facts about the secondary market and lists similar concepts like 'Money market'.

Concluzii

În contextul dezvoltării explozive a numărului de documente pe web este absolut necesară existența unor medii care să permită utilizatorilor explorarea în scopul extragerii cunoștințelor din texte și structuri de documente web. Această activitate trebuie să fie integrată în scopul integrării ontologiilor de mari dimensiuni existente astăzi pe web. În concluzie, un mediu hermenofor integrează instrumente hermenofore cu ontologii într-o arhitectură în care utilizatorul trebuie să poată experimenta, să investigheze diverse transformări ale textelor. Se poate spune că un mediu hermenofor înglobează sinergic instrumente de prelucrare a cunoștințelor cu instrumente de prelucrare a textelor și cu tehnici specifice web.

Bibliografie

- [ABK02] G. Angelova, S. Boytcheva, O. Kalaydjiev, St. Trausan-Matu, P. Nakov, A. Strupchanska, Adaptivity in a web-based CALL system, in F. van Harmelen (ed.): ECAI 2002. Proceedings of the 15th European Conference on Artificial Intelligence, IOS Press, Amsterdam, 2002, to appear

-
- [Bla85] L. Blaga, Trilogia culturii, Ed. Minerva, 1985
- [BIM96] Black, J.B., McClintock, An Interpretation Construction Approach to Constructivist Design, in B.G. Wilson (ed.), Constructivist Learning Environments: Case Studies in Instructional Design, Education Technology Publications, 1996.
- [Cli95] Clibbon, K., Conceptually Adapted Hypertext For Learning, Proceedings of CHI'95, http://www.acm.org/sigchi/chi95/Electronic/documnts/kc_bdy.html
- [CTr01] Constandache, G.G., St. Trausan-Matu, Ontologia si hermeneutica calculatoarelor, Editura Tehnica, 2001.
- [Cul94] Culianu, I.P., Eros si magie în Renastere; 1484, Nemira, Bucuresti 1994.
- [Eng95] Engelbart, D.G., Toward Augmenting the Human Intellect and Boosting our Collective IQ, CACM No.8, Vol.38, Aug. 95, pp. 30-33.
- [ErK97] Ertl, B., Kraan, A.G., Internet-Based Learning Environments from a Constructivist point of view, Proceedings of RILW, Ilieni, 1997, p. 17-21.
- [Goo] <http://www.google.com>
- [Gru96] Gruber, T., What is an Ontology, <http://www.kr.org/top/definitions.html>
- [JoL83] Johnson-Laird, P.N., Mental Models - Towards a Cognitive Science of Language, Inference, and Consciousness, Cambridge Univ. Press, 1983.
- [LaJ80] Lakoff, G., Johnson, M., Metaphors We Live by, The University of Chicago Press, 1980.
- [Lar] LarFLaST, <http://www-it.fmi.uni-sofia.bg/larflast/>
- [Mar97] Marcus, S., Empatie si personalitate, Ed. Atos, 1997.
- [MHD] J.C. Mallery, R. Hurwitz, G. Duffy, Hermeneutics, Encyclopedia of Artificial Intelligence, pp. 596-611.
- [Noi86] C. Noica, Scrisori despre logica lui Hermes, Ed. Cartea Româneasca, 1986.
- [OORG] <http://www.ontology.org/main/papers/faq.html>
- [Sow99] J. Sowa, Knowledge Representation: Logical, Philosophical and Computational Foundations, Brooke Cole Publishing Co., Pacific Grove, CA, 1999, vezi si [CTr01].
- [THH95] Thiring, M., Hannemann, J., Haake, J.M., Hypermedia and Cognition: Designing for Comprehension, Communications of the ACM, vol.38, no. 8, pp. 57-66, aug. 1995.
- [TMC02] St. Trausan-Matu, D. Maraschi, S. Cerri, Ontology-Centered Personalized Presentation of Knowledge Extracted From the Web, in S.Cerri, G.Gouarderes (eds.), Intelligent Tutoring Systems 2002, Springer, Lecture Notes in Computer Science number 2363, to appear.
- [Tra95] St. Trausan-Matu, Programe inteligente pentru asistarea invatarii, in Revista

Romana de Informatica si Automatica, vol.5, nr.4, 1995, pag. 7-16.

- [Tra00] St. Trausan-Matu, *Metaphor Processing for Learning Terminology on the Web*, in S.A.Cerri (ed.), *Artificial Intelligence, Methodology, Systems, Applications 2000*, Springer-Verlag, ISBN 3-540-41044-9, 2000, pp.232-241
- [Tra01] St. Trausan-Matu, *Interfatarea evoluata om-calculator*, Ed. MatrixRom, 2001.
- [Wes97] D.West, *Hermeneutic Computer Science*, CACM, Vol.40, No.4, pp. 115-116, 1997, si în [CTr01].
- [Wil96] B.G. Wilson (ed.), *Constructivist Learning Environments: Case Studies in Instructional Design*, Education Technology Publications, 1996
- [Win87] T. Winograd, *Thinking machines: Can there be? Are we?*, Report No. STAN-CS-87-1161, Stanford, 1987.
- [WN] WordNet, <http://www.cogsci.princeton.edu/~wn/>
- [Wri95] von Wright, G.H., *Explicatie si intelegere*, Humanitas, 1995.
- [WSD97] <http://www.research.att.com/~suciu/workshop-papers.html>
- [XML] www.w3.org/xml

SECȚIUNEA III

**TEHNOLOGII ALE LIMBAJULUI
VORBIT**

Experimente în vederea recunoașterii vorbitorului

Corneliu BURILEANU,
Universitatea „Politehnica” din București, Spl.Independenței 303
cburileanu@messnet.pub.ro

Luigi BOJAN,
Graphco Technologies Inc., Newton, PA, USA

1. Introducere

Având în vedere funcția realizată și concomitent, sarcina de îndeplinit, tehnologia vorbirii se poate clasifica în mai multe domenii [1, 2]:

- Recunoașterea automată a vorbirii. Se bazează pe analiza automată a semnalului vocal și are în vedere informația transmisă de om mașinilor care “îl ascultă”. Din această informație, mașina este programată să extragă acele caracteristici ce îi vor permite să deceleze cine vorbește, ce vorbește, în ce fel și în ce condiții.
- Sinteza automată a vorbirii. Se realizează răspunsul “prin voce” al mașinilor către operatorul uman.
- Codificare/decodificare (analiză și sinteză) a vorbirii. Se referă la tehnici de compresie a informației conținută în semnalul vocal în vederea unor prelucrări ulterioare specifice sarcinii de îndeplinit.

Un domeniu interdisciplinar important, legat în mod esențial de aplicațiile de recunoaștere și sinteză automată ale vorbirii este cel al dialogului om-mașină.

Termenul “comunicare om - mașină” pare forțat: mașina nu este o entitate socială, nu are nici scop nici cultură. Ea nu poate acționa în lumea reală în sensul de a putea să răspundă corect la întrebări de genul: “ai putea să închizi ușa, te rog?”. Ea nu este “conștientă” decât de propria sa “lume”. Avem într-adevăr nevoie de a comunica cu mașinile? Au importanță intențiile lor, chiar dacă le-ar avea? Ce poate să-mi comunice sau să mă facă să știu o mașină?

Mașina îmi procură “unelte” pentru a realiza o sarcină, ea mă face să proiectez noi obiecte (eventual, virtuale), ea mă aduce într-un univers artificial, îmi permite să utilizez un mediu de programare împreună cu alți utilizatori umani, pentru a lucra într-o manieră cooperantă în același mediu informatic. Mașina se prezintă deci ca un *factor de interacțiune*. Ea trebuie să-mi furnizeze un spațiu de muncă, unelte și metode. Dar pentru aceasta, mașina trebuie adaptată sarcinii curente sau unor sarcini noi, să adopte un comportament “comprehensibil”, să se arate “prietenoasă” etc. Paradoxul este deci evident:

mașina trebuie să fie, dintr-un anumit punct de vedere, *socială* pentru a colabora eficient cu un utilizator în scopul îndeplinirii sarcinilor, din ce în ce mai complexe, care îi sunt încredințate.

Preocupările noastre în domeniul tehnologiei vorbirii au, între altele, scopul de a oferi mijloacele pentru o comunicare între om și mașină prin mesaje vorbite [3]. Această comunicare este doar un aspect al dialogului. Rămâne în continuare deschisă problema definirii conceptelor și cea a stabilirii unor strategii de dialog adecvate sarcinii de rezolvat.

Semnalul vocal conține o varietate de informații utile: ce se vorbește, cine vorbește, în ce fel și în ce condiții. În cadrul recunoașterii se pune problema identificării unui anumit tip de informații; de pildă, recunoașterea cuvintelor rostite înseamnă determinarea mesajului (ce se vorbește) indiferent (sau ajutându-se) de variabilitățile introduse de vorbitor (cine), maniera de a vorbi (în ce fel) și zgomotul ambiental (în ce condiții). Putem particulariza afirmând că **recunoașterea vorbirii** este procesul de transformare a semnalului acustic continuu produs de organul fonator uman într-o reprezentare discretă căreia i se poate atașa o semnificație și care, când e înțeleasă, poate fi folosită pentru a determina un răspuns.

Problemele majore pe care le ridică recunoașterea automată sunt legate de

- discretizarea semnalului vocal care, din punctul nostru de vedere înseamnă *segmentare*;
- caracterul adecvat al răspunsului ce depinde de natura sarcinii de îndeplinit; modalitatea de prelucrare este irelevantă.

Proiectarea unui sistem de recunoaștere presupune câteva opțiuni fundamentale de abordare. Punctul de vedere adoptat poate viza prelucrarea unui semnal acustic ca oricare altul, poate ține seama de mecanismul producerii vorbirii, poate simula recepția senzorială, sau poate folosi modelul uman al percepției vorbirii.

Termenul de **recunoaștere a vorbitorului** desemnează orice aplicație de discriminare a persoanelor pe baza vocii acestora. Procedurile de recunoaștere se desfășoară în două etape [4]:

- etapa de antrenare: colectarea de material vocal de la persoana care se dorește a fi recunoscută;
- etapa de testare: compararea unui fragment de vorbire neidentificat cu datele provenite din antrenare și luarea deciziei de recunoaștere.

Există două subclase de aplicații:

• **verificarea vorbitorului** își propune să determine dacă un fragment de semnal vocal aparține sau nu unui anumit vorbitor [5, 6, 7, 8]. Există doi parametri care caracterizează performanțele sistemului: respingerea adevăratului vorbitor și acceptarea unui impostor. Considerând un set de N vorbitori, informația (în biți) obținută este

$$I_{ver} = I \quad (1)$$

presupunând probabilitatea de verificare *a priori* egală cu 0.5;

• **identificarea vorbitorului** are ca scop punerea în corespondență a unei voci necunoscute cu un vorbitor dintr-un set dat [9, 10, 11, 12]. Pentru N vorbitori, informația (în biți) obținută este

$$I_{ident} = \log_2(N) \quad (2)$$

considerând probabilitatea de identificare *a priori* egală pentru toți vorbitorii.

Rezultă că, potențial, un sistem automat de verificarea vorbitorului are performanțe mai bune.

O clasificare suplimentară a automatelor de recunoaștere are în vedere natura sarcinii de îndeplinit și se reflectă în complexitatea sistemului [13]:

- sisteme de recunoașterea vorbitorului **dependente de text** - textul utilizat în faza de antrenare este același cu cel de testare;
- sisteme **independente de text** - indiferent de materialul vocal avut la dispoziție.

Setul de vorbitori vizat poate impune, de asemenea, o clasificare a automatelor:

- “*set închis*” – pentru procesul de identificare descris ca mai sus;
- “*set deschis*” - în cazul identificării există posibilitatea ca vocea necunoscută să nu aparțină niciunui dintre vorbitorii din setul dat, numărul de decizii posibile fiind în acest caz $N + 1$. Identificarea pe “set deschis” devine astfel o combinație a proceselor de verificare și identificare.

2. Reprezentarea parametrică

Variabilitățile pronunțării pentru diverși vorbitori, sau la un același vorbitor, la momente de timp diferite, constituie una dintre dificultățile majore ale sarcinii de recunoaștere a vorbitorului. Deosebiri de vorbire depind de dialect, context, stil de exprimare, stare emoțională etc. Mai mult, în opinia noastră, așa cum vom încerca să argumentăm mai departe, *limba în care se vorbește* impune deosebiri de abordare și diferențe ale performanțelor automatului [14].

Din acest motiv, alegerea judicioasă a **caracteristicilor acustice** care vor fi utilizate în procesul de recunoaștere este deosebit de importantă:

- să diferențieze vorbitori diferiți dar să fie tolerante pentru același vorbitor;
- să fie ușor măsurabile din semnalul vocal;
- să fie stabile în timp;
- să nu fie susceptibile de a fi contrafăcute de potențiali impostori.

Având în vedere cerințele formulate mai sus, am decis utilizarea *parametrilor cepstrali*.

Anumite abordări ale prelucrării semnalului vocal presupun adoptarea unor decizii fundamentale de dezvoltare a analizei: considerarea unui model de producere a vorbirii având ca prototip aparatul fonator uman, separarea efectelor sursei vorbirii de comportarea tractului vocal propriu-zis, o serie de aproximări care să facă analiza eficientă în condiții normale de procesare [15]. Variația (lentă) în timp a formei tractului vocal este aproximată printr-o serie de secvențe de durată suficient de mică pentru a presupune forma invariantă: este ceea ce se numește “analiza în timp scurt”. Dacă, în plus, în aceste durate “scurte” de timp se presupune că tractul este caracterizat în mod esențial de frecvențele sale de rezonanță, se ajunge la un model al cărui parametri se pot deduce prin rezolvarea unui sistem de ecuații liniare. Deși aproximările avute în vedere par destul de restrictive, *analiza prin predicție liniară* (LPC) dă rezultate deosebite pentru că semnalul vocal are o redundanță deosebită; este motivul pentru care metoda ne permite să aproximăm un eșantion de semnal printr-o combinație liniară (deci este liniar predictibil) dintr-un număr de eșantioane precedente. Desigur, principiile în sine ale metodei nu sunt noi; ele au permis însă, în decursul ultimilor ani, evoluția spre metode mai sofisticate [16, 17].

Nici principiile *analizei cepstrale* (analiză care, așa cum vom arăta, se poate baza pe rezultatele analizei LPC) nu sunt noi: se dezvoltă un mecanism care să permită decelarea mai amănunțită a influențelor diverselor elemente ale organului fonator. O serie de presupuneri fundamentale de abordare se păstrează (modelarea producerii vorbirii în maniera aparatului fonator uman, analiza “în timp scurt”); dar separarea efectelor excitației glotale, tractului vocal și radiației buzelor poate fi făcută într-o modalitate care ține seama mai detaliat de fiecare efect în parte [18, 19].

În concluzie, presupunerile fundamentale care stau la baza parametrizării propuse sunt:

- efectele excitației tractului vocal și ale tractului propriu-zis pot fi separate;
- tractul vocal este invariant pe durate scurte de timp, ceea ce are drept rezultat obținerea unui model descris de un sistem liniar ai cărui parametri variază lent în timp (constanți “în timp scurt”).

Fundamental pentru modul în care concepem abordarea analizei semnalului este asimilarea *analizei* cu *parametrizarea* semnalului și, în consecință, cu *compresia* sa. Alegerea parametrilor a avut în vedere și considerente pragmatice:

- complexitatea prelucrării;
- gradul de compresie,
- tipul de aplicație,
- în ce măsură parametrii sunt semnificativi și robuști.

O primă variantă a schemei bloc care descrie funcționarea sistemului de recunoașterea vorbitorului este prezentată în fig. 1. Blocul de *preprocesare* presupune filtrarea și achiziția semnalului în condiții normale pentru orice sistem de recunoaștere. În această secțiune vom descrie obținerea *cepstrului* pornind de la *analiza LPC*, iar în secțiunea următoare vom descrie principiile *cuantizării vectoriale* și deci procedura de recunoaștere propriu-zisă.

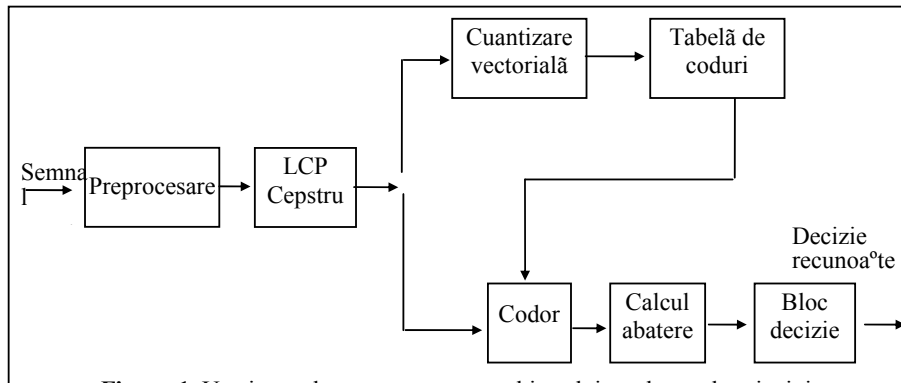


Figura 1. Un sistem de recunoaștere a vorbitorului – schema de principiu

Fie semnalul vocal presupus a fi convoluția unei excitații și a funcției de transfer a tractului vocal:

$$s(t) = e(t) * v(t) \quad (3)$$

Analiza homomorfică care duce la obținerea *cepstrului* presupune aplicarea unui operator neliniar “ H ”

$$s(n) \xrightarrow{H} \hat{s}(n) \quad (4)$$

în care $\hat{s}(n)$ va fi numit *cepstrul complex* asociat semnalului $s(n)$.

Prin definiție

$$\hat{S}(z) = \sum_n \hat{S}(n) \cdot z^{-n} \equiv \ln S(z) \quad (5)$$

Astfel, cepstrul complex asociat semnalului devine

$$\hat{S}(n) = \hat{e}(n) + \hat{v}(n) \quad (6)$$

ceea ce permite separarea componentelor printr-o “filtrare temporală” aplicată cepstrilor

$$s(n) \xrightarrow{H} \hat{s}(n) \begin{cases} \xrightarrow{L} \hat{e}(n) \xrightarrow{H^{-1}} e(n) \\ \xrightarrow{L} \hat{v}(n) \xrightarrow{H^{-1}} v(n) \end{cases} \quad (7)$$

Obținerea parametrilor cepstrali se poate realiza ținând seama de câteva proprietăți ale cepstrului.

Fie $c(n)$ partea pară a cepstrului complex al semnalului

$$c(n) = [\hat{s}(n) + \hat{s}(-n)] / 2 \quad (8)$$

Secvența $c(n)$ se numește *cepstrul real* al semnalului $s(n)$

$\hat{s}(n)$ este o secvență cauzală – ca și $s(n)$; rezultă

$$\hat{s}(n) = c(n) \cdot \begin{cases} 0 & \text{pentru } n < 0 \\ 1 & \text{pentru } n = 0 \\ 2 & \text{pentru } n > 0 \end{cases} \quad (9)$$

Cum transformata “z” a unei secvențe cauzale e determinată complet prin partea reală a transformatei sale Fourier, rezultă

$$c(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \ln |S(e^{j\omega})| e^{jn\omega} d\omega \quad (10)$$

Vom prefera calculul coeficienților cepstrali din coeficienți α_i ai analizei prin predicție liniară (LPC) conform relațiilor recursive:

$$\begin{aligned} c(1) &= -\alpha_1 \\ c(i) &= -\alpha_i - \sum_{n=1}^{i-1} \left(1 - \frac{n}{i}\right) \cdot \alpha_n \cdot c(i-n) \quad i > 0 \end{aligned} \quad (11)$$

Figura 2 prezintă evoluția coeficienților cepstrali pentru o voce feminină și una masculină.

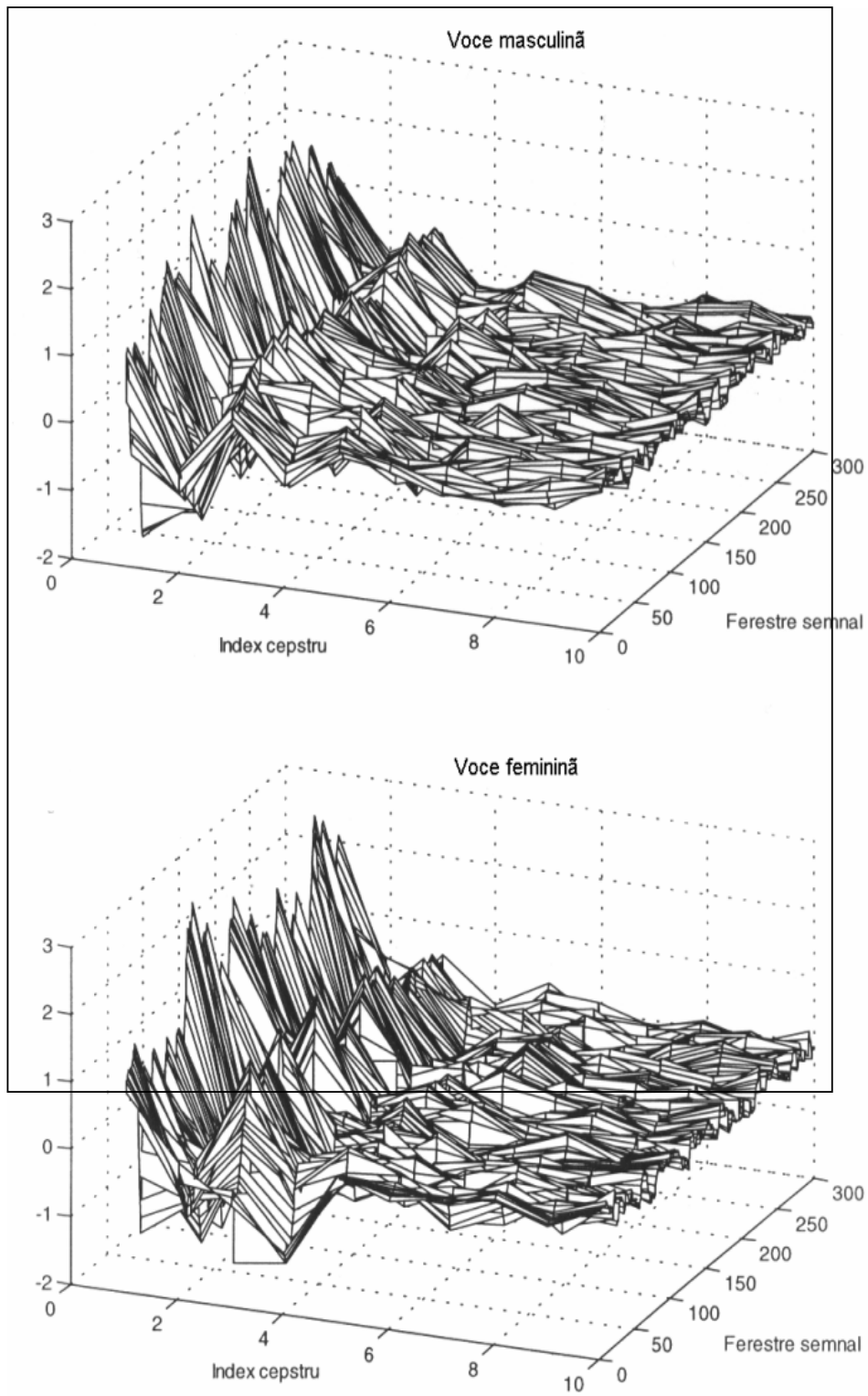
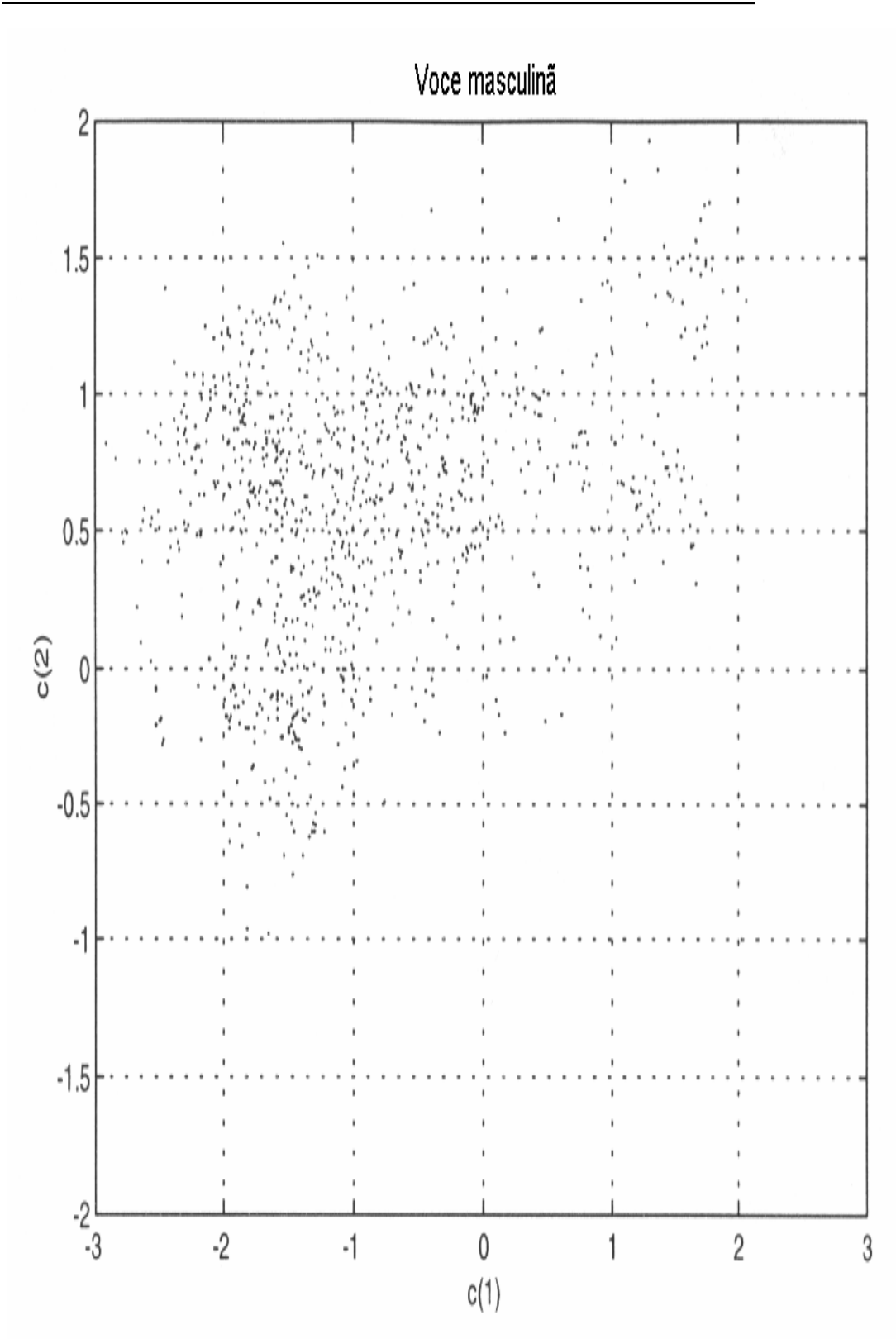


Figura 2. Evoluția în timp a coeficienților cepstrali ai semnalului vocal

Materialul vocal a fost achiziționat folosind un microfon de calitate (considerat fără zgomot) și a fost eșantionat cu frecvența de 8 kHz. Intervalele de analiză au lungimea de 240 ms, cu o suprapunere de 160 ms. Analiza prin predicție liniară s-a efectuat cu ordinul de predicție $p = 10$, iar pentru estimarea coeficienților de predicție liniară s-a folosit algoritmul Levinson-Durbin. O primă observație este aceea că modulul amplitudinii coeficienților este descrescător cu ordinul acestora. Pentru coeficienții de ordinul 5-10, evoluția coeficienților cepstrali tinde să devină uniformă. Amplitudinea redusă a acestora anunță existența unor dificultăți de estimare în condiții de zgomot.

În scopul unei aprecieri calitative, fig. 3 prezintă distribuția coeficienților cepstrali în planul $c(1)$ - $c(2)$, pentru aceiași doi vorbitori (masculin și feminin). Se poate observa distribuția diferită a principalilor coeficienți cepstrali pentru cei doi vorbitori. Se remarcă o concentrare a coeficienților în anumite zone ale planului $c(1)$ - $c(2)$.





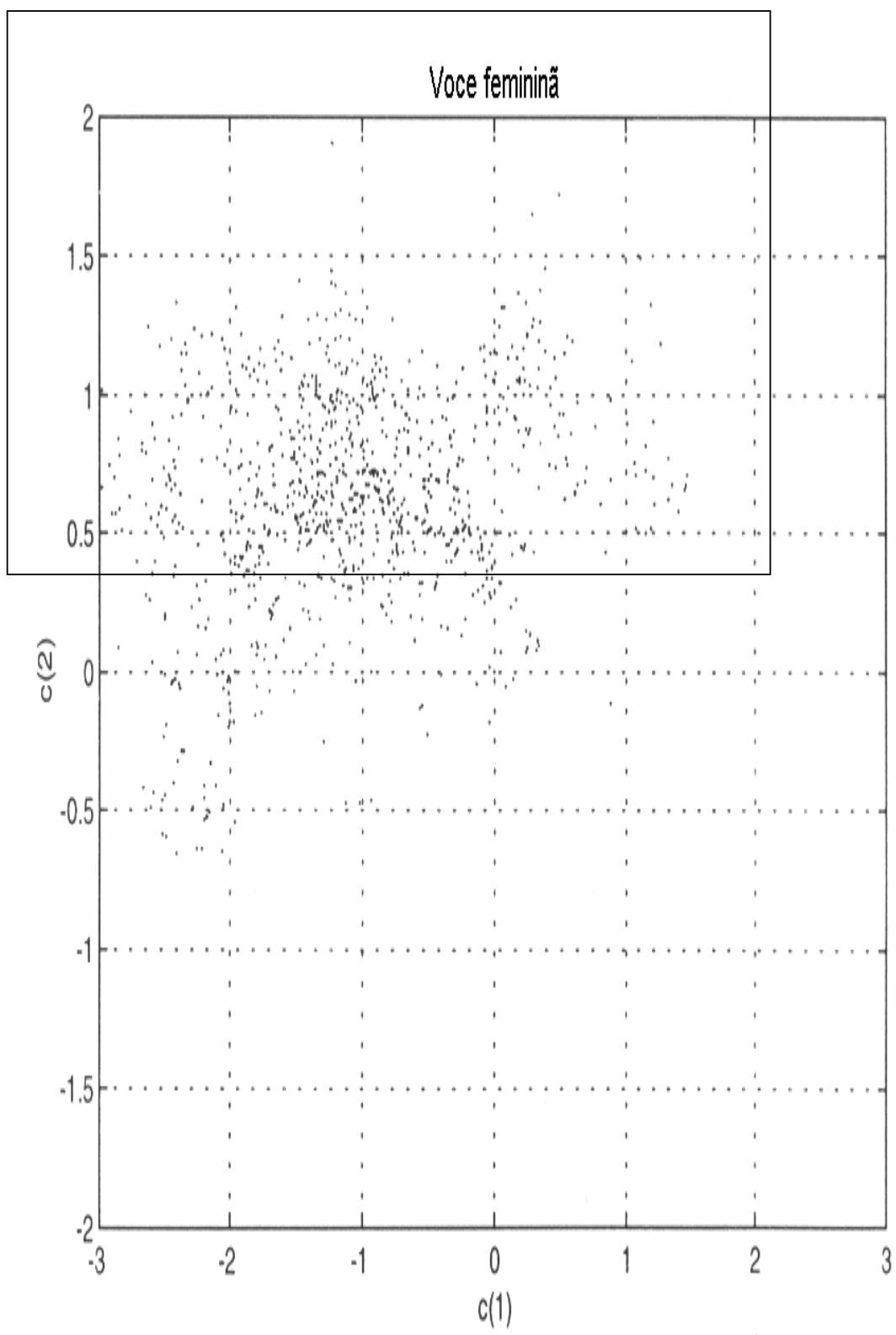


Figura 3. Reprezentarea coeficienților cepstrali în planul $c(1) - c(2)$

În fig. 4 este prezentată distribuția parametrilor cepstrali corespunzători unui semnal vocal compus numai din vocalele limbii române. Ordinul analizei cepstrale este $p = 12$. Reprezentarea grafică s-a făcut numai în planul $c(1) - c(2)$. Se observă faptul că vocalele sunt relativ ușor separabile în spațiul cepstral, într-o configurație asemănătoare celei din spațiul formantic. Această analiză oferă premise interesante și pentru recunoașterea vorbirii în limba română.

Figura 4. Semnal compus din vocale și parametrii cepstrali corespunzători

3. Cuantizarea vectorială

Din punctul de vedere al sistemelor de recunoaștere a vorbitorului, o persoană produce în timpul vorbirii o secvență de vectori de parametri. Aceștia caracterizează atât vorbitorul cât și cuvintele pronunțate. Pentru un interval de timp suficient de lung, ne așteptăm ca datele achiziționate să acopere spațiul vectorial într-un mod care depinde mai mult de caracteristicile vorbitorului și mai puțin de ceea ce a pronunțat. Se face presupunerea că, având la dispoziție un volum suficient de date, se poate genera un model al vorbitorului care să fie utilizat într-un proces de recunoaștere [20, 21].

Principiul cuantizării vectoriale este aplicat în sensul compresiei unui volum mare de vectori acustico-fonetici, reprezentând material vocal pronunțat de către un vorbitor, într-un set restrâns de vectori denumit *tabelă de coduri* (sau de *centroizi*). În etapa de antrenare, partiționarea spațiului acoperit de vectorii spectrali este făcută astfel încât media distanțelor minime dintre fiecare vector cepstral și cel mai apropiat centroid să fie minimizată. În etapa de testare, un set de vectori provenind de la un vorbitor necunoscut, este codat utilizând tabela de vectori corespunzătoare vorbitorului vizat. Distorsiunea totală medie este utilizată în decizia de recunoaștere [22].

Fie $\{X_n\}$ ansamblul de N versiuni cunoscute ale vectorului X .

Fie $\{G_k\}$ o partiție a acestui ansamblu în K clase; o clasă G_k cuprinde g_k elemente, astfel ca

$$\sum_{k=1}^K g_k = N \quad (12)$$

Notăm cu $X_p^{(k)}$ cuvântul prototip ("centroid", "vector-cod") al unei clase G_k

Distanța medie între centroizi este

$$\overline{\sigma_p} = \frac{1}{K \cdot (K-1)} \sum_{i,j=1}^K D(X_p^{(i)}, X_p^{(j)}) \quad (13)$$

Distanța medie între vectorii dintr-o aceeași clasă, parcurgând toate clasele este

$$\overline{\sigma_G} = \frac{1}{K} \sum_{k=1}^K \frac{1}{g_k \cdot (g_k - 1)} \sum_{i,j=1}^{g_k} D(X_i^{(k)}, X_j^{(k)}) \quad (14)$$

Raportul $\frac{\overline{\sigma_p}}{\overline{\sigma_G}}$ reprezintă calitatea partiției

Algoritmul utilizat pentru găsirea centrozilor este atunci următorul:

• dacă cei K centroizi sunt aleși la întâmplare, clasele sunt constituite asociind fiecare vector X centroidului cel mai apropiat:

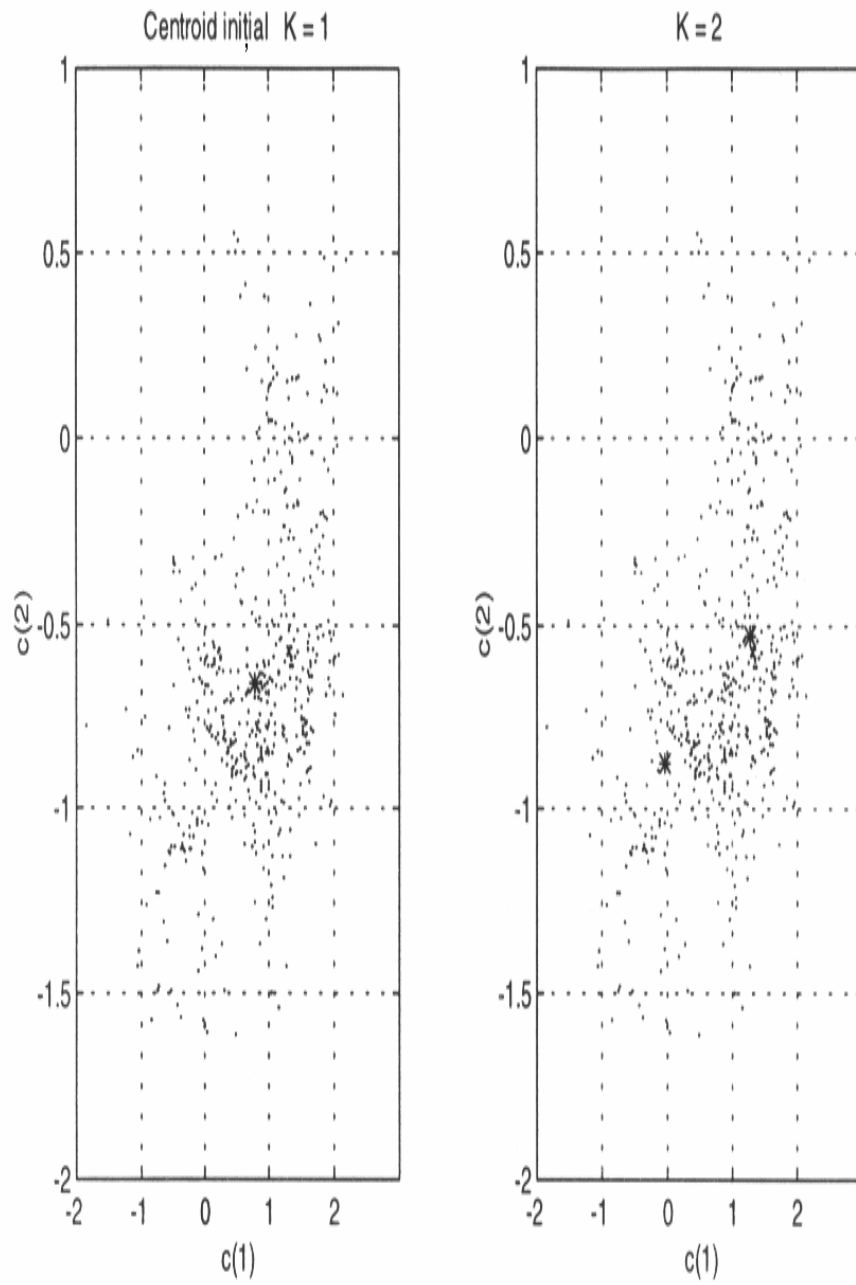
$$X_n \in G_k \quad \text{dacă} \quad D(X_n, X_p^{(k)}) < D(X_n, X_p^{(i)}), \quad \forall i \neq k \quad (15)$$

• se iterează găsirea centroizilor căutând în fiecare clasă k vectorul $X_n^{(k)}$ care are distanța față de vectorul cel mai depărtat al clasei minimă:

$$X_n^{(k)} \equiv X_p^{(k)} \quad \text{dacă} \quad \max_m D(X_n^{(k)}, X_m^{(k)}) \text{ e minimă} \quad (16)$$

• această procedură e iterată până când centroizii sunt stabiliți.

Prezentăm în fig. 5 un exemplu de cuantizare vectorială folosind algoritmul Linde-Buzo-Gray (LBG). Vectorii cuantizați sunt coeficienții cepstrali de predicție liniară. Pentru reprezentarea în plan s-a ales sistemul de coordonate $c(1)$ - $c(2)$. Dimensiunea tabelii de centroizi aleasă este $M = 8$. Se poate observa cum, în urma operației de optimizare, centroizii tind să “acopere” întregul spațiu ocupat de vectori. În mod evident, eroarea de cuantizare scade pe măsură ce dimensiunea tabelii de centroizi crește.



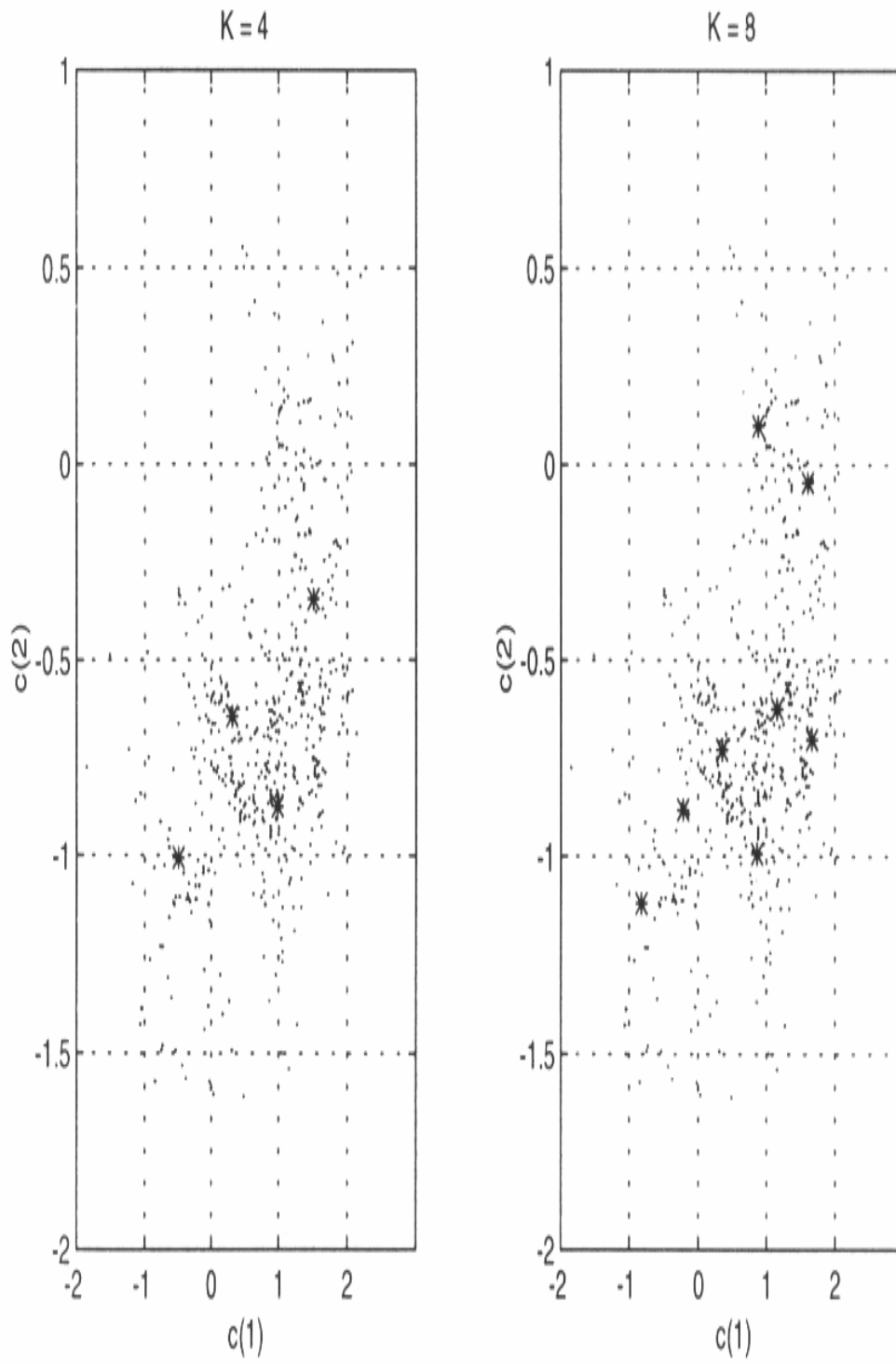


Figura 5. Evoluția algoritmului de cuantizare vectorială (8 centrozi)

“•” – vectori cepstrali; * – centrozi

Pe parcursul algoritmului se pot utiliza diverse strategii de divizare. De exemplu, dacă după o operație de divizare și reclasificare, una dintre clase devine subpopulată sau chiar vidă, o alta va fi divizată la pasul următor, pentru a menține constant numărul total de clase. Se pot folosi următoarele criterii de alegere a clasei care va fi divizată: clasa care posedă cel mai mare număr de elemente, clasa care prezintă distorsiunea totală cea mai mare, clasa care prezintă distorsiunea medie cea mai mare. Folosind această structură arborescentă, clasificarea unui vector se poate efectua prin asocieri succesive, printr-o parcurgere a claselor găsite pentru fiecare nivel de divizare. În aplicațiile care necesită o acuratețe de clasificare ridicată, se preferă o metodă de clasificare prin căutare exhaustivă.

4. Rezultate obținute

Un aspect important în proiectarea automatelor de recunoaștere a vorbitorului (eventual independent de text) îl reprezintă posibilitatea de evaluare a performanțelor acestora. Pentru a putea evalua un astfel de automat cu o precizie acceptabilă este nevoie de o bază de date corespunzătoare [23]. O astfel de bază de date trebuie să îndeplinească următoarele cerințe:

- să cuprindă material vocal achiziționat de la cât mai mulți vorbitori (de preferat, de ordinul zecilor sau sutelor);
- să conțină, eventual, dialecte diferite;
- să conțină fraze cât mai variate;
- frazele să fie rostite de mai multe ori, la intervale de timp
- pentru evaluare în condiții reale (de exemplu transmisie telefonică), materialul vocal trebuie să fie achiziționat prin intermediul mai multor aparate telefonice, în decursul mai multor legături, de preferat la distanțe diferite [24, 25].

Proiectarea și construirea unei astfel de baze de date este o sarcină dificilă.

Am folosit mai multe baze de date: internaționale, oarecum standard pentru procedurile de recunoaștere – “TIMIT” și “YOHO”, precum și o bază de date proprie, în română și engleză – “DiSPALL”.

Baza de date “TIMIT” conține eșantioane de voce provenind de la 630 de vorbitori, fiecare pronunțând 10 fraze. Experimentele descrise în lucrare au fost efectuate pe secțiunea TEST, care conține 168 vorbitori. Cele 10 fraze sunt: două fraze de calibrare (SA), cinci fraze compacte din punct de vedere fonetic (SX) și trei fraze variate contextual (SI). În experimente s-au folosit frazele SA și SX în faza de antrenare și frazele SI în cea de testare. Pentru evaluarea efectelor zgomotului telefonic în algoritmi de recunoaștere a

vorbitorului, s-a folosit o variantă a bazei de date numită “NTIMIT”. Aceasta conține același material vocal ca și baza “TIMIT” cu deosebirea că acesta a fost transmis prin intermediul rețelei telefonice. Transmisia s-a făcut folosind un echipament de simulare a tractului vocal uman, în legături telefonice reale, la diferite distanțe.

Baza de date “YOHO” cuprinde fraze rostite de 138 de vorbitori (106 bărbați și 32 femei), iar vocabularul folosit constă din numere de două cifre rostite în grupuri de câte trei. Pentru fiecare vorbitor am folosit 4 sesiuni de antrenare de câte 24 de enunțuri și 10 sesiuni de verificare de câte 4 enunțuri.

Baza de date proprie “DiSPPALL” [26] cuprinde materialul vocal de la 26 de vorbitori (23 de bărbați și 3 femei) cu vârsta ce variază de la 21 la 50 de ani. Fiecare vorbitor în parte a rostit 31 de fraze: 11 fraze echilibrate din punct de vedere fonetic au fost folosite pentru antrenare și 20 de fraze pentru verificare: 5 enunțuri de bază repetate de câte 4 ori. Frazele de verificare au fost înregistrate în două sesiuni diferite, 5 enunțuri de bază fiind repetate de două ori în fiecare sesiune. Prima sesiune de verificare a fost înregistrată în același timp cu sesiunea de antrenare, iar sesiunea a doua a fost înregistrată după două-trei săptămâni. Înregistrările s-au făcut cu un microfon de tip “head-set” într-o cameră cu zgomot ambiental normal: spre deosebire de baza “YOHO”, baza “DiSPPALL” conține material vocal alterat de zgomot pentru a face condițiile de test mai dificile și mai apropiate de o situație reală de recunoaștere a vorbitorilor

În experimentele de verificare a vorbitorului, o frază de test este comparată cu referința vorbitorului a cărui identitate se dorește verificată, calculându-se o distorsiune totală medie. Dacă aceasta este mai mică decât un prag dat, vorbitorul este considerat acceptat, altfel el este respins. Există două tipuri de erori asociate procesului de verificare: respingerea utilizatorului căruia îi aparține referința (denumită eroare de tip I) și acceptarea unui impostor (eroare de tip II) [27]. Fiecare frază de test este comparată cu referințele corespunzătoare tuturor vorbitorilor din baza de date aleasă pentru test. Pragurile de decizie nu sunt fixate *a priori* ci se determină distanța medie totală pentru care eroarea de tip I este egală cu cea de tip II (“rata-erorii-egale”). Valoarea corespunzătoare a erorii este considerată rezultatul final al procesului de evaluare. În fig. 6 sunt prezentate rezultatele procesului de verificare a vorbitorului, folosind cuantizarea vectorială, utilizând baza de date “TEST/TIMIT”. Ordinul de predicție (și implicit dimensiunea vectorilor cepstrali) este $P = 10$ iar dimensiunea tabelii de centroizi, $M = 64$. Ca distanță vectorială s-a folosit distanța euclidiană ponderată.

$$d(v_a, v_b) = \frac{1}{s_j^2} \sum_{j=1}^p (v_{aj} - v_{bj})^2 \quad (17)$$

unde s_j^2 este varianța componentei j calculată pe întreg setul vectorilor de antrenare. Ca metodă de cuantizare vectorială s-a folosit algoritmul LBG modificat.

Sunt evidente tendințele contrare ale erorilor de tip I, respectiv II. Rata-erorii-egale pentru evaluarea de mai sus este 6.8%, corezpunzând unui prag de decizie egal cu 2.8. În funcție de aplicația concretă, pragul de decizie se poate stabili *a posteriori* la o altă valoare, adecvată scopului propus. Spre exemplu, dacă se dorește limitarea acceptării impostorilor la 2%, respingerea adevăraților utilizatori va fi de 19.7%. Reciproc, pentru o eroare de respingere a utilizatorilor reali de 2%, acceptarea impostorilor va fi de 12.9%.

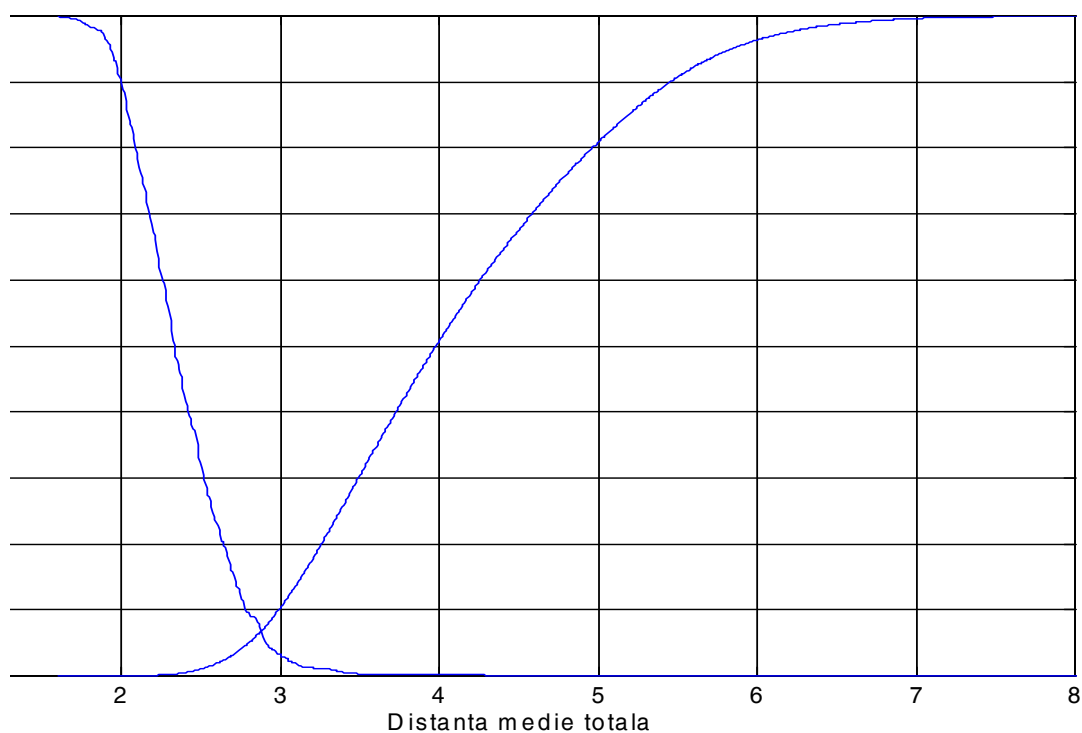


Figura 6. Eroarea de verificare a unui sistem de recunoaștere a vorbitorului utilizând cuantizarea vectorială

În experimentele de identificare a vorbitorului, fiecare frază de test provenind de la un vorbitor considerat necunoscut este comparată cu referințele fiecărui vorbitor din baza de date aleasă pentru test. Referința asociată cu cea mai mică distorsiune totală medie față de fraza de test este considerată ca aparținând vorbitorului identificat. În funcție de corespondența dintre apartenența frazei de test și a referinței aceluiași vorbitor sau unor

vorbitori diferiți, se decide dacă rezultatul procesului de identificare este adevărat sau fals. Eroarea de identificare este calculată ca raportul dintre numărul de identificări incorecte și numărul total de identificări [28, 29, 30].

5. Utilizarea frecvenței fundamentale în recunoașterea vorbitorului

Frecvența fundamentală poate fi utilizată ca parametru discriminator suplimentar în conjuncție cu algoritmi de cuantizare vectorială a vectorilor cepstrali.

Frecvența fundamentală F_0 sau *perioada fundamentală* T_0 (cunoscută și sub numele de "pitch"), constituie un parametru important al vocii umane, care își găsește utilizări practice în multe domenii ale procesării vorbirii. Încercări de utilizare a frecvenței fundamentale în procesul de recunoaștere a vorbitorului se cunosc încă de la începutul anilor '70, aceasta fiind pusă în corespondență directă cu prozodia. Majoritatea acestor experimente s-au desfășurat utilizând sisteme de recunoaștere dependente de text și metode de aliniere temporală. Sistemele de recunoaștere a vorbitorului independente de text bazate exclusiv pe frecvența fundamentală nu au dat rezultate satisfăcătoare.

Ideea prezentată în secțiunea de față este aceea de a folosi frecvența fundamentală ca un parametru discriminator suplimentar, în conjuncție cu algoritmi de cuantizare vectorială a vectorilor cepstrali [31]. Justificarea teoretică a acestei abordări rezidă în primul rând în modelul predicției liniare aplicat semnalului vocal, care presupune, așa cum am arătat o separare clară între sursa de semnal și tractul vocal. De asemenea, am arătat în secțiunea 2 că analiza cepstrală folosită pentru extragerea vectorilor cepstrali este un proces de deconvoluție, coeficienții cepstrali obținuți caracterizând în mod exclusiv tractul vocal. Ca atare, utilizarea ca date de intrare în același sistem atât a vectorilor cepstrali cât și a frecvenței fundamentale nu reprezintă o abordare redundantă.

Cerințele de bază ale unui algoritm de extragere a frecvenței fundamentale sunt: acuratețea de estimare (evitarea armonicilor), robustețea deciziei sonor/nesonor, insensibilitatea la zgomot, volumul de calcule minim. Se cunosc numeroși algoritmi de estimare a frecvenței fundamentale (AMDF, Dubnowski, Rabiner, SIFT, etc.), fiecare prezentând avantaje și dezavantaje. Trebuie arătat faptul că, din cauza, în principal, comportării nestaționare a semnalului vocal, niciunul din algoritmi cunoscuți nu este considerat perfect. Cu alte cuvinte, se acceptă ideea existenței erorilor atât în luarea deciziei sonor/nesonor cât și în obținerea valorilor propriu-zise ale frecvenței fundamentale. În experimentele prezentate mai jos s-a folosit algoritmul Rabiner, considerat ca fiind unul dintre cele mai robuste.

Ideea introdusă este aceea de a utiliza frecvența fundamentală în scopul unei clasificări grosiere a potențialilor candidați, atât pentru sarcina de verificare a vorbitorului, cât și pentru cea de identificare. În cazul verificării, scopul propus este acela de a reduce eroarea de tip II, prin eliminarea vorbitorilor a căror frecvență fundamentală nu "corespunde" cu cea a vorbitorului de referință. În cazul sarcinii de identificare, se dorește

reducerea numărului de candidați posibili, fără a afecta acuratețea de identificare. Aceasta poate conduce la o reducere majoră a volumului de calcule, dat fiind că estimarea frecvenței fundamentale se face o singură dată pentru fiecare vorbitor și este mai puțin consumatoare de timp decât clasificarea vectorială.

Având în vedere considerentele de mai sus, schema de principiu a sistemului de recunoaștere a vorbitorului modificat prin introducerea frecvenței fundamentale ca parametru discriminator este prezentată în fig. 7.

Un aspect important în utilizarea frecvenței fundamentale în aplicațiile de recunoaștere a vorbitorului îl reprezintă alegerea formei de prelucrare a datelor furnizate de estimator. “Conturul de pitch”, reprezentând evoluția în timp a parametrului F_0 , deși utilizat în sisteme de recunoaștere a vorbitorului dependente de text, conține un volum de date dificil de utilizat în operații de discriminare. În consecință, s-a încercat o reducere a datelor la câțiva parametri statistici. Au fost investigate patru valori statistice derivate din conturul de pitch: valoarea medie, valoarea maximă, valoarea minimă și dispersia (deviația standard). Pentru fiecare vorbitor, s-au calculat aceste valori pe ansamblul materialului vocal disponibil. Ca parametru de discriminare a fost utilizat raportul valorilor statistice de mai sus

$$R_{p,medie} = \frac{\text{media } F_0 \text{ pentru antrenare}}{\text{media } F_0 \text{ pentru test}} \quad (18)$$

tratându-se în mod similar toate celelalte valori statistice (maximă, minimă, dispersie). Pentru a evalua utilitatea acestor parametri în procesul de discriminare, s-a determinat distribuția fiecăruia atât pentru frazele pronunțate de aceiași vorbitori (intra-vorbitor) cât și pentru toate combinațiile de fraze pronunțate de vorbitori diferiți (inter-vorbitor).

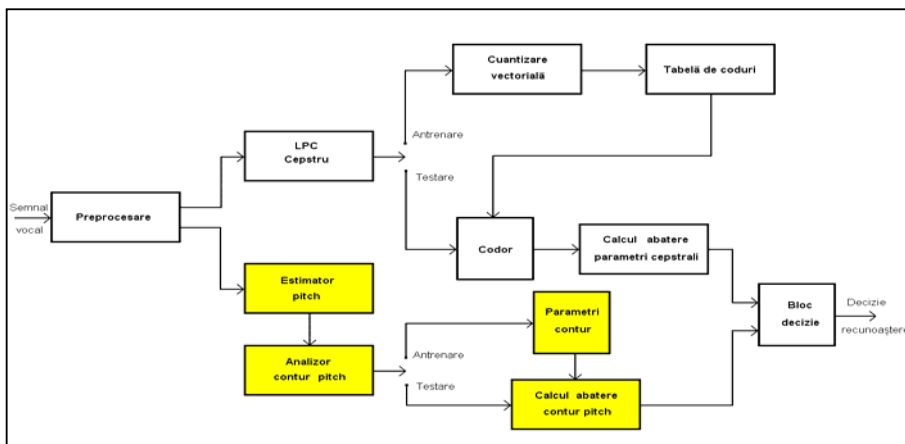


Figura 7. O variantă a sistemului de recunoaștere a vorbitorului – schema de principiu

Modul de discriminare a vorbitorilor este următorul: fixându-se un prag ε , dacă

$$R_{p,medie} \notin [l - \varepsilon, l + \varepsilon] \quad (19)$$

vorbitorul este rejectat și nu se execută clasificarea vectorială. În caz contrar, vorbitorul este considerat potențial candidat și urmează procesul de clasificare prin cuantizare vectorială.

Utilizând elementul de discriminare descris mai sus s-au obținut îmbunătățiri importante atât în procesul de verificare a vorbitorului cât și în cel de identificare. Rezultatele obținute pentru 14 coeficienți cepstrali și 128 centroizi sunt prezentate în tabelul 1.

Tabelul 1

ε	Neutilizat	0.30	0.25	0.20	0.15	0.10
EER la verificare (%)	6.3	6.1	5.3	3.9	2.7	6.5
Eroarea de identificare (%)	6.2	6.2	5.9	5.6	5.5	9.4
Candidați identificare (%)	100	57.2	49.1	43.4	32.3	26.5

Cele mai bune rezultate s-au obținut pentru $\varepsilon = 0.15$, caz în care eroarea de verificare obținută este de aproape 2.5 ori mai mică decât în cazul folosirii doar a clasificării vectoriale. În cazul identificării, deși îmbunătățirile de acuratețe nu sunt impresionante, cel mai important rezultat îl reprezintă reducerea numărului candidaților, cu peste 65%. Pentru valori ale lui ε mai mici decât 0.10, se observă o degradare abruptă a performanțelor de verificare și identificare, ceea ce indică faptul ca variația intra-vorbitor a frecvenței fundamentale medii este mai mare decât acest prag.

6. Concluzii

Lucrarea de față se ocupă de un aspect bine delimitat al tehnologiei vorbirii și anume recunoașterea vorbitorului ca parte integrantă a recunoașterii automate și mai departe a dialogului om-mașină. Tipurile de probleme care apar sunt similare pentru întreg domeniul recunoașterii automate.

Am precizat presupunerile fundamentale care au stat la baza analizei propuse (în special opțiunea de a aborda proiectarea ținând seama de mecanismul producerii vorbirii); insistăm asupra faptului că aceste abordări nu sunt obligatorii, ci constituie alternative care au avantaje și dezavantaje.

S-au trecut în revistă etapele esențiale ale procedurilor de recunoașterea vorbitorului: achiziția semnalului vocal, prelucrarea acustico-fonetică, recunoașterea propriu-zisă.

Am subliniat importanța parametrizării semnalului vocal. Analiza cepstrală care a fost aleasă pentru reprezentarea parametrică a semnalului vocal este legată de opțiunile fundamentale de analiză: separarea efectelor sursei de semnal și ale tractului, separarea efectelor diverselor porțiuni din tractul vocal, analiza “în timp scurt”

Am utilizat cuantizarea vectorială ca metodă de recunoaștere. Sunt prezentate o parte dintre rezultatele experimentelor realizate. Subliniem importanța utilizării unor baze de date specifice și, în consecință, am acordat spațiu prezentării acestora.

O contribuție pe care o considerăm interesantă la îmbunătățirea performanțelor recunoașterii vorbitorului o constituie utilizarea frecvenței fundamentale ca parametru discriminator grosier. Sunt prezentate o serie de rezultate care probează în ce mod anumite performanțe sunt superioare abordării “clasice”.

O parte dintre rezultatele obținute sunt susceptibile de a fi generalizate pentru recunoașterea vorbirii în limba română [32] (de pildă, coeficienții cepstrali pentru foneme ale limbii române). De asemenea, utilizarea frecvenței fundamentale apropie recunoașterea vorbitorului de o anumită dependență de limba în care sunt rostite frazele de antrenare și de test.

Referințe bibliografice

- [1] M.Dragănescu, C.Burileanu, coordonatori (1986). Analiza și sinteza semnalului vocal – Editura Academiei Române, București.
- [2] M.Dragănescu, G.Ștefan, C.Burileanu (1991). Electronica funcțională – vol. I, Editura tehnică, București, ISBN 973-31-0290-3.
- [3] G. Yu and H. Gish (1993). Identification of Speakers Engaged in Dialog, Proc. of IEEE Int. Conf. Acoust., Speech, Signal Processing, Vol.II, p. 383-386.
- [4] Sadaoki Furui (1994). An Overview of Speaker Recognition Technology, Proc. of ESCA Workshop on Automatic Speaker Recognition, Identification and Verification, p. 1-9.
- [5] Y. Bennani, P. Gallinari (1994). Connectionist Approaches for Automatic Speaker Verification – ESCA Workshop on Speaker Recognition, Identification and Verification, p. 95-103.
- [6] M. Hanah s.a. (1994). The Role of the Reference Template in Speaker Verification – ESCA Workshop on Speaker Recognition, Identification and Verification, p. 181-184.
- [7] Chi-Shi Liu; Hsiao-Chuan Wang; Lee, C. (1996) Speaker Verification Using Normalized Log-Likelihood Score, IEEE Tr. on Speech and Audio Processing, Vol. 4. Issue 1, p. 56

-
- [8] S. Nakagawa, K. P. Markov (1997). Speaker Verification Using Frame and Utterance Level Likelihood Normalization, Proc. of SPCHL97 ,Vol. 2, p. 1087.
- [9] K.T. Assaleh, R.J. Mammone (1994). New LP – Derived Features for Speaker Identification – IEEE Tr.on SAP, vol.2, no.4, p. 630-638.
- [10] H. Gish, M. Schmidt (1994). Text-Independent Speker Identification – IEEE Signal Proc. Mag., vol.11, nr.4, p. 18-32.
- [11] Q. Lin s.a. (1994). Microphon Array Speaker Identification – IEEE tr. on ASSP, vol.2. nr.4, p. 622-629.
- [12] D. Reynolds (1994). Experimental Evaluation of Features for Robust Speaker Identification – IEEE Tr. on ASP, vol.2, nr.4, p. 639-643.
- [13] F. Bimbot, G. Chollet, A. Paoloni (1994). Assesment Methodology for Speaker Identification and Verification – ESCA Workshop on Speaker Recognition, Identification and Verification, p. 75-83.
- [14] M. Abe, S. Sagayama (1990). Statistical Study on Voice Individuality Conservation Across Different Languages – Proc. of ICSLP, p. 157-160.
- [15] Y. Gong, J.P. Haton (1994). Non-Linear Interpolation Methods for Speaker Recognition – ESCA Workshop on Speaker Recognition, Identification and Verification, p .23-26.
- [16] J. He s.a. (1995). On the Use of Features from Prediction Rersidual Signal in Speaker Identification Proc. of EUROSPEECH95, p. 313-316.
- [17] D.Naik s.a. (1994). Robust Speaker Identification Using Pole Filtering – ESCA Workshop on Speaker Recognition, Identification and Verification, p. 225-228.
- [18] J. Openshaw, J. Masson (1994). Optimal Noise-Masking of Cepstral Features for Robust Speaker Identification – ESCA Workshop on Speaker Recognition, Identification and Verification, p. 231-234.
- [19] J. Thompson, J.S. Masson (1993). Within Class Optimization of Cepstra for Speaker Recognition, Proc. of EUROSPEECH, p. 165-168.
- [20] K. Sonmez, L. Heck, M. Weintraub (2000). Multiple Speaker Tracking and Detection: Handset Normalization and Duration Scoring, Digital Signal Processing, 10(1/ 2/3), p. 133-143.
- [21] T. Isobe, J. Takahashi (1999). A New Cohort Normalization Using Local Acoustic Information for Speaker Verification, Proc. IEEE Int. Conf. Acoust. Speech, Signal Processing, 26.8, vol. 2, p. 841-844.
- [22] X. Zhu s.a (1994). Text-Independent Speaker Recognition Using VQ, Mixture Gaussian VQ and Ergodic HMMs – ESCA Workshop on Speaker Recognition, Identification and Verification, p. 55-58.
- [23] L. Boves s.a. (1994). Design and Recording of Large Data-Bases for Use in Speaker

-
- Recognition and Identification – ESCA Workshop on Speaker Recognition, Identification and Verification, p. 43-46.
- [24] A. Federico (1995). Parametric Speaker Recognition Over Large Population of Telephonic Voices – Proc. of EUROSPEECH95, p. 329-332.
- [25] J.L. Gauvain s.a (1995). Experiments with Speker Verification over the Telephone – Proc. of EUROSPEECH95, p. 651-654.
- [26] C. Burileanu, D. Burileanu s.a.(2000). Cohort Normalisation Methods for Speaker Verification – Proc. of International Conference “Communications 2000”, Bucharest, Romania, p.118-121.
- [27] M. Wagner s.a. (1994). Analysis of Type-II Errors for VQ-Distortion Based Speaker Verification – ESCA Workshop on Speaker Recognition, Identification and Verification, p. 83-86.
- [28] J.F. Bonastre (1993). Automatic Spaker Recognition and Analytic Process – Proc. of EUROSPEECH93, p. 441-444.
- [29] M. Sugiyama s.a. (1993). Speech Segmentation, Clustering Based on Speaker Features – Proc. of ICASSP, p.395-398.
- [30] H. Beigi, S. Maes and J. Sorensen (1998.) A Distance Measure Between Collections of Distributions and Its Application to Speaker Recognition, Proc. of ICASSP, Vol. 2, p. 753-756.
- [31] L.E. Bojan, C. Burileanu s.a. (1996). Enhancements in Automatic Speaker Verification and Identification for Large Data-bases Using Pitch Contour Analysis – Proc. of ICSPAT96, Boston, SUA, p. 1796-1800
- [32] C. Burileanu, L.E. Bojan s.a. (1993). A Representation for Recognition of Isolated Words Spoken in the Romanian Language – Proc. of ICSPAT93, Santa Clara, USA, p. 1478-1484.

Prelucrarea inițială a textului de intrare în cadrul unui sistem de sinteză a vorbirii pornind de la text în limba română

Dragoș BURILEANU
Laboratorul de "Tehnologia vorbirii și prelucrarea digitală a semnalelor",
Facultatea de Electronică și Telecomunicații, Universitatea "POLITEHNICA"
București
Bdul Iuliu Maniu 1-3, Sector 6, 77202 București
bdragos@mESsnet.pub.ro

1. Introducere

Limbajul reprezintă modalitatea de exprimare a ideilor prin intermediul unui ansamblu de semne, fie grafic, fie prin gesturi, sau sunete, un astfel de sistem structurat fiind specific doar oamenilor. Fără îndoială, *vorbirea* este una din principalele sale componente; ea este cea mai veche modalitate de comunicare între oameni și este și astăzi cea mai răspândită. Este deci ușor de înțeles faptul că vorbirea a fost studiată intens și s-a încercat adesea să fie prelucrată într-un mod automat. Pentru mulți ingineri și specialiști din domeniu, posibilitatea de a conversa liber cu o mașină reprezintă de fapt o adevărată provocare pentru înțelegerea cât mai deplină a proceselor de producere și percepție implicate în comunicarea prin voce între oameni. Ceea ce este însa și mai important este faptul ca *interfețele de comunicare prin voce* devin tot mai mult o necesitate. În viitorul apropiat, sistemele și rețelele interactive vor oferi un acces simplu și ieftin la cantități mari de informație și servicii, ceea ce va afecta fundamental viața noastră zilnică.

Deși principiile de bază ale producerii și recepționării vorbirii au început să fie studiate încă de la sfârșitul secolului al XVIII-lea, când s-au înregistrat primele cercetări în domeniul dezvoltării sintetizoarelor mecanice de sunete asemănătoare vocii umane, *tehnologiile de prelucrare a vorbirii* au obținut rezultate semnificative doar în ultimele decenii (fiind denumite în sens larg *tehnici de analiză și sinteză a semnalului vocal*). Aceste rezultate au fost posibile datorită progreselor făcute în domeniile acusticii și lingvisticii, modelării matematice a producerii și percepției vorbirii, prelucrării semnalelor și tehnologiilor VLSI. Putem evidenția în acest sens dezvoltarea procesoarelor numerice de semnal pe un singur chip, realizarea de capsule de memorie mai mari și mai ieftine, apariția unor algoritmi îmbunătățiți pentru prelucrare de semnal, iar în domeniul comunicațiilor

crearea de standarde globale pentru transmisie, compresie de semnal și protocoale de comunicație.

Prin urmare, putem aprecia că cercetările actuale în domeniul prelucrării vorbirii au ca scop larg îmbunătățirea calității, securității și costului comunicațiilor și a accesului uman la informații. Pe de o parte, este de așteptat în viitorul apropiat o extindere importantă a serviciilor integrate de voce, poșta electronică, FAX, paging și transmisiuni de date pe canale fără fir. Pe de altă parte însă, comunicarea verbală între om și mașini, în ambele sensuri, tinde deja să devină o realitate, fiind vizibilă tendința actuală de a apropia caracteristicile mașinii de cele ale utilizatorului uman.

În acest ultim sens, trebuie observat faptul că tendința menționată anterior este absolut firească. Filozoful grec Aristotel (384 - 322 î.C., fondator al logicii formale), afirma: "*Rațiunea de a fi a oricărui lucru constă în funcția sa*". Ori este evident faptul că o interfață de dialog prin voce reprezintă o modalitate ideală de comunicare cu mașina, vorbirea fiind cea mai naturală, flexibilă, eficientă și economică modalitate de comunicare utilizată de oameni.

Aceste idei legate de posibilitatea comunicării prin voce între om și mașina nu sunt noi; totuși, doar în ultimii ani a început să prindă contur conceptul ce a căpătat denumirea de "dialog om-mașină", iar tehnologia necesară implementării acestui concept a părăsit deja laboratoarele și a pătruns în lumea reală, într-o gamă largă de aplicații.

Pentru a realiza un mod de comunicare cât mai natural și pentru a permite o utilizare cât mai largă, calculatorul trebuie să înțeleagă și să producă singur vorbirea; acesta este motivul principal pentru care *recunoașterea și sinteza vorbirii* au devenit în ultimii ani tehnologii de un interes special și constituie subiecte pentru cercetări intense și aprofundate. Ambele tehnologii prelucrează vorbirea în primul rând sub aspectul conținutului informațional: recunoașterea transformă vocea omului în text ce poate fi folosit literal (de exemplu pentru dictare), sau o interpretează sub forma unor comenzi de control pentru diverse aplicații, iar sinteza permite generarea limbajului vorbit pornind de la text sau de la anumite concepte.

Cu toate că s-au făcut pași importanți în aceste domenii, rezultatele sunt încă departe de așteptări. Sarcinile enunțate inițial s-au dovedit în timp a fi deosebit de dificile, în primul rând datorită complexității semnalului vocal ca și a dificultăților legate de prelucrarea acestuia, dificultăți legate fie de recunoașterea conținutului său informațional (semnalul vocal depinzând puternic de vorbitor și de condițiile în care acesta rostește un mesaj), fie de producerea sa, fie de transmiterea acestui semnal la distanță [1].

În acest context, producerea vorbirii artificiale și în special conversia text – voce, care constituie obiectul principal al lucrării de față, este astăzi un obiectiv de bază al domeniului prelucrării vorbirii și subiect al unor cercetări intense. Un *sistem de sinteză pornind de la text* (TTS – "*Text-to-Speech*") poate oferi o gamă variată de aplicații, de la accesul la poșta electronică și diferite tipuri de baze de date, la pronunțarea unui text pentru persoane cu handicap vizual.

Este important de observat faptul că tehnologia de răspuns prin voce prezintă o serie de avantaje fundamentale pentru transmiterea informației:

- oricine poate înțelege un mesaj, fără antrenare sau concentrare deosebită;
- mesajul poate fi recepționat chiar dacă cel ce ascultă este implicat în alte activități, cum ar fi mersul, manipularea unor obiecte, sau citirea altor informații;
- rețeaua telefonică convențională poate fi utilizată pentru accesul rapid la distanță la o bază de informații;
- această formă de comunicare este mai economică decât cea tradițională prin mesaje scrise.

Toți acești factori precum și numeroasele aplicații cerute de industrie au creat premisele unor cercetări aprofundate, obținându-se astfel în multe țări sisteme comerciale care pot produce vorbire sintetică pornind de la text, cu o inteligibilitate acceptabilă.

Într-adevăr, scopul principal al celor mai multe sisteme de sinteză existente este de a produce o vorbire *inteligibilă*. Din acest punct de vedere, sinteza pare a fi de mai multă vreme o tehnologie "stabilă", ieftină și ușor de implementat; se spune chiar, uneori, că acest domeniu este în prezent suficient de bine dezvoltat, iar problemele rămase sunt minore din punct de vedere științific. Dacă însă scopul este sinteza în timp real, pornind de la un vocabular nelimitat de cuvinte și fără restricții asupra textului, iar vorbirea să fie nu numai inteligibilă, ci și la fel de *naturală* ca cea umană, atunci se constată că performanțele actuale sunt departe de a fi satisfăcătoare. Rămân încă multe probleme importante de rezolvat: extinderea vocabularului oferit, înlăturarea restricțiilor impuse textului în privința unor caractere speciale, îmbunătățirea caracteristicilor de *prozodie*, posibilitatea de modificare a *ritmului* și *stilului* vorbirii sintetizate, sau elaborarea unor sisteme de sinteză în mai multe limbi. Aceste sarcini se dovedesc a fi deosebit de dificile și cer, evident, eforturi interdisciplinare susținute [2].

2. Sinteza automată a vorbirii

Etimologic, cuvântul "sinteză" provine din limba greacă și semnifică îmbinarea mai multor elemente diferite într-un tot.

În ceea ce privește *sinteza vorbirii*, nu există o definiție precisă și unanim acceptată de către specialiștii în tehnologia vorbirii. Acest termen a avut în decursul timpului mai multe accepțiuni, majoritatea depinzând de nivelul tehnologic al momentului și de elementele constitutive ale semnalului vocal care au fost folosite pentru sinteză. De exemplu, primele circuite integrate care permiteau simpla restituire a unui mesaj vocal înregistrat și stocat digital au purtat denumirea de "sintetizoare vocale", fie că se făcea sau nu o compresie a semnalului. Este evident că în acest caz nu se poate vorbi de sinteză, din moment ce **textul este fix** și astfel de sisteme nu pot rosti decât mesaje preînregistrate;

chiar dacă vocea umană este comprimată cu ajutorul unui algoritm, nu este cu adevărat "sintetică", ci poate fi numită mai curând o "înregistrare cu număr redus de biți".

Aceeași situație este în cazul sintezei la recepție a unor mesaje transmise pe canale de comunicație standard (caracteristică sistemelor de tip "vocoder"), care este de obicei considerată ca făcând parte din domeniul *codării vorbirii* și cuprinde tehnici de reducere a debitului semnalului vocal pentru transmisie; cu alte cuvinte, și acest tip de sinteză, care reface **același** mesaj analizat la emisie, deci nu generează fraze **noi**, nu este tratat ca o sinteză automată propriu-zisă.

O categorie distinctă de sinteză vocală este aceea care implică sisteme ce concatenează cuvinte sau fraze preînregistrate, dar generează fraze noi, acestea nefiind niciodată pronunțate ca atare; astfel de sisteme cer utilizarea unor reguli lingvistice mai mult sau mai puțin complicate pentru a funcționa corespunzător.

În sfârșit, o categorie specială o reprezintă sinteza vorbirii pornind de la text; aceasta reprezintă, în esență, transformarea unui text oarecare, scris într-un anumit limbaj, în semnal vocal. Trebuie remarcat faptul că în prezent, în multe lucrări științifice, acest tip de sinteză este sinonim chiar conceptului de sinteză automată a vorbirii.

Analizând exemplele de mai sus, putem defini trei noțiuni generale [3], pe care le vom utiliza pe parcursul lucrării de față:

Definiția 2.1 *Sinteza automată a vorbirii* este "tehnologia integrată care simulează procesul uman de generare a vorbirii, mergând de la sisteme simple ce pot genera automat fraze noi și cuprind un formalism lingvistic minimal și până la sisteme care transformă în vorbire reprezentări simbolice sau lingvistice ale limbajului".

Definiția 2.2 Un *sistem de sinteză pornind de la text* este "un sistem automat care poate produce vorbirea plecând de la un text scris, prin intermediul unei reprezentări fonetice a mesajului".

Definiția 2.3 *Sintetizorul vocal* este "etajul unui sistem de sinteză automată a vorbirii care realizează conversia finală în semnal vocal, pornind de obicei de la o reprezentare parametrică a unor segmente acustice fundamentale".

3. Sinteza vorbirii pornind de la text

Pentru a înțelege mai bine dificultatea sarcinii unui sistem de sinteză pornind de la text, considerăm că este util să punem în evidență mai întâi modul (fiziologic) în care o persoană citește cu voce tare un text. Imaginea textului este sesizată de neuronii sistemului vizual, transmisă creierului sub forma unor stimuli electrici, aici fiind prelucrată pentru a putea permite comanda neuronilor responsabili de corecta activare a plămânilor, coardelor vocale și organelor articulatorii. În acest fel se produce vorbirea, ea fiind permanent

monitorizată de creier (în special prin intermediul organelor auditive), în scopul ajustării configurației tractului vocal în timp real.

Desigur, cunoaștem încă prea puțin despre organizarea de ansamblu a sistemului nervos uman, care este capabil de această activitate complexă; putem propune totuși următorul *model funcțional* prin care este prelucrată informația optică și apoi este dată comanda de generare a vorbirii:

- Atunci când citim un text, efectuăm practic o sarcină de *recunoaștere de caractere*, ignorând, parțial inconștient, anumite erori de redactare a cuvintelor (caractere lipsă sau înlocuite cu altele) și decodificând mai degrabă cuvântul ca un întreg; are loc un proces de inferență a informației dintr-un context posibil incomplet. De asemenea, recunoaștem cu ușurință caractere speciale sau abrevieri.
- Considerând *fonemele* ca fiind cele mai mici elemente sonore care permit diferențierea între ele a cuvintelor, este evident că secvența fonemică corespunzătoare unui cuvânt diferă de șirul de caractere grafice din care este compus cuvântul; creierul trebuie să facă prin urmare o *transcriere fonetică pornind de la litere*, această operație practic instinctivă permițând pronunția unui număr nelimitat de cuvinte.
- În cele mai multe situații, suntem capabili să începem pronunția unei fraze mult înainte de terminarea ei; cu alte cuvinte, putem face o *structurare sintactică*, descompunând fiecare propoziție în grupuri de cuvinte și asociindu-le intonația corespunzătoare. Și acest proces este practic inconștient, fiind bazat pe educație și experiență.
- În sfârșit, putem discrimina cu ușurință cuvinte ce se scriu asemănător dar se pronunță diferit, după *înțelesul semantic*, fapt posibil datorită aceleiași capacități de deducție a creierului de care am vorbit mai sus.

Concluzia este simplă: pe baza experienței lingvistice căpătate în urma educației, o persoană familiară cu limbajul în care este scris un text depășește imediat pașii descriși anterior și poate cu ușurință să citească cu voce tare textul scris, în primul rând pentru că înțelege ceea ce citește.

Având în vedere considerațiile expuse anterior, devine evident faptul că o mașină care trebuie să pronunțe un text scris nu va putea adopta o schemă de prelucrare atât de complicată cum este cea care caracterizează acțiunea citirii cu voce tare a unui text de către o persoană. Sunetele vorbirii sunt inerent guvernate de ecuații diferențiale ale mecanicii fluidelor, aplicate într-un context nestaționar, deoarece presiunea aerului la nivelul plămânilor, tensiunea glotală, ca și configurațiile tractului vocal și nazal, evoluează în timp. Toate acestea sunt controlate de creierul uman, care beneficiază de avantajul puterii sale de prelucrare paralelă pentru extragerea esenței textului citit: **înțelesul**. Chiar și la nivelul la care a ajuns știința astăzi (cercetări intense în domeniile sintezei articulatorii, rețelelor neuronale artificiale și prelucrării limbajului natural), construirea unui sistem de sinteză

pornind de la text cu un model atât de complex rămâne practic nerealizabilă; chiar dacă, să spunem, s-ar ajunge foarte aproape de aceste cerințe, sistemul rezultat nu ar fi de loc compatibil cu criteriile economice normale.

Figura 1 introduce o diagramă funcțională foarte generală a unui sistem TTS, bazată pe observațiile anterioare.

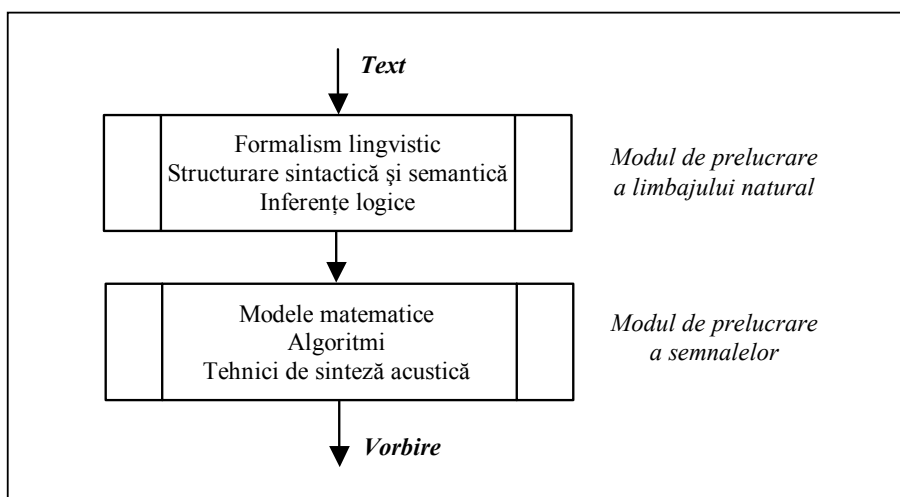


Figura 1. Diagramă funcțională pentru un sistem TTS

Ca și pentru un cititor uman, schema cuprinde un *modul de prelucrare a limbajului natural*, capabil să producă o transcriere fonetică a textului citit, împreună cu informații despre intonație, accente, durate și de asemenea un *modul de prelucrare a semnalelor*, care transformă informația simbolică primită în vorbire sintetică, pe baza unor tehnici de sinteză adecvate și a unor structuri stocate în urma unei analize preliminare. Etapele de bază ale sintezei pornind de la text pot fi astfel descrise printr-un număr de transformări succesive ce trebuie aplicate asupra șirului de caractere ce reprezintă textul de intrare; scopul este de a se obține o vorbire **de calitate**, într-o limbă oarecare, fără constrângeri asupra textului introdus.

Trebuie menționat faptul că formalismul descris poate "sări" uneori peste anumiți pași, dacă se utilizează în mod adecvat cunoașterea lingvistică și matematică; acest lucru se întâmplă atunci când punem anumite restricții asupra textului ce trebuie pronunțat, sau impunem vorbirii sintetizate o inteligibilitate și o naturalitate moderate. Cu alte cuvinte, proiectarea sistemului TTS se poate simplifica dacă se impun sistemului sarcini precise, corespunzătoare unor aplicații concrete.

Colectivul nostru de cercetare a început acum câțiva ani dezvoltarea unui sistem complet TTS în limba română, bazat pe *concatenare de difoneme*. Arhitectura acestui sistem este prezentată în Figura 2. Sistemul cuprinde o parte importantă de prelucrare lingvistică și un modul de generare a semnalului de vorbire având la bază un algoritm de tip PSOLA [4]. După realizarea unei prime variante a sistemului, se depun în continuare eforturi pentru creșterea naturaleții vorbirii sintetizate, prin îmbunătățirea performanțelor la diferite nivele de prelucrare.

Modulul de prelucrare a limbajului într-un sistem TTS are ca sarcină transformarea textului de intrare într-o reprezentare fonetică și prozodică, care trebuie să descrie cât mai fidel posibil pronunția sa. Acest lucru poate fi realizat parcurgând mai multe etape succesive, puse în evidență și în figura anterioară. Vom discuta în cele ce urmează **modalitățile de rezolvare a părții de prelucrare inițială (preprocesare) a textului în cadrul sistemului nostru de sinteză în limba română.**

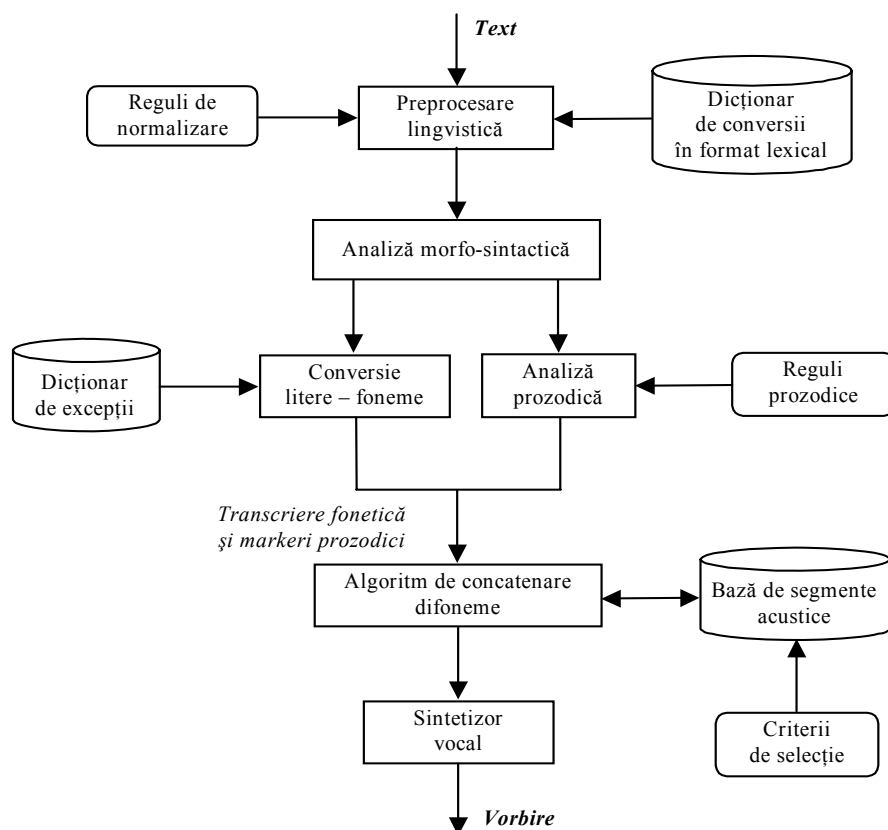


Figura 2. Arhitectura sistemului TTS în limba română

4. Preprocesarea textului de intrare în cadrul sistemului TTS în limba română

4.1 Probleme generale

Una dintre dificultățile majore ale sistemelor TTS constă în faptul că aceste sisteme trebuie să poată prelucra practic **orice text**, plecând de la propoziții simple izolate

și mergând până la paragrafe complexe, care pot cuprinde un număr mare de propoziții, cu posibile structuri negramaticale și simboluri speciale. Ca atare, partea de preprocesare lingvistică a textului are un rol extrem de important, deoarece detectarea corectă și interpretarea șirurilor de caractere de intrare influențează acuratețea întregului sistem de sinteză și contribuie la conversia unui text fără restricții în vorbire sintetică.

Uzual, un text scris se prezintă sub forma unei secvențe de caractere ASCII; el este alcătuit din cuvinte compuse cu ajutorul literelor alfabetului, dar și din alte tipuri de caractere: spații albe, semne de punctuație, șiruri de numere, sau alte simboluri speciale (de exemplu operatori matematici). Textul poate conține numerale (*12*, *12.450*, *1,245*), abrevieri (*prof.*, *dr.*, *ing.*), sau acronime (*IBM*, *S.R.L.*, *TTS*). Aceste secvențe sunt de obicei "anormale" din punct de vedere lingvistic față de majoritatea cuvintelor din text și trebuie mai întâi transformate într-un format ce poate fi recunoscut de partea de analiză lingvistică. Această sarcină revine modulului de preprocesare, care trebuie de asemenea să realizeze o segmentare a textului de intrare (detectarea cuvintelor și a sfârșitului frazelor) și o prelucrare a semnelor de punctuație și a simbolurilor speciale [5, 6, 7, 8].

La prima vedere, preprocesarea unui text pentru un sistem TTS poate părea banală; în realitate însă, lucrurile sunt destul de complicate. Spre exemplu, nu este totdeauna posibilă determinarea marginilor unei fraze pe baza semnelor de punctuație. Astfel, punctul (.) poate apare și la sfârșitul unei fraze, dar și în multe alte situații, ca de exemplu în abrevieri (*ing.*), acronime (*S.R.L.*), indicația că se omite un anumit fragment de text (...), sau numerale (*12.450 – douăsprezece mii patru sute cincizeci*), situații care trebuie diferențiate prin procedee adecvate [9]. De asemenea, cratima creează dificultăți în operația de segmentare; ea poate fi folosită pentru despărțirea în silabe, pentru scrierea cuvintelor compuse, pentru delimitarea unui nou paragraf, sau în enumerări.

O sarcină dificilă este și conversia anumitor secvențe de simboluri în cuvinte care să poată fi analizate lingvistic. Dacă unele abrevieri uzuale pot fi "expandate" imediat, cu ajutorul unui tabel de echivalențe, există multe situații în care secvențe de simboluri care nu se pot distinge pe baza ortografierii lor, cer tipuri diferite de conversii; de exemplu, numărul format din șapte cifre *6123456* poate reprezenta un număr întreg sau un număr de telefon și va trebui citit diferit în cele două situații. În general, prezența șirurilor de numere în text ridică numeroase dificultăți, deoarece ele pot apare în diferite contexte: ore, date, numere de telefon, expresii aritmetice etc.

Trebuie observat că aceste ambiguități create de natura multifuncțională a semnelor de punctuație sau de modul diferit de citire a acelorași secvențe de simboluri, pot avea implicații majore asupra acurateții întregului proces de prelucrare lingvistică și în final asupra pronunției corecte a textului de către sistemul de sinteză.

Evident, numărul secvențelor de caractere neuzuale dintr-un text ce se dorește a fi transformat în vorbire depinde mult de tipul și subiectul textului. Spre exemplu, textele literare dintr-un volum de proză sau comentariile politice dintr-un ziar au mult mai puține situații dificile decât comentariile economice, sportive, sau prezentările de spectacole. În ultimele situații menționate, construcțiile neuzuale, criptice sau chiar negramaticale,

abrevierile uneori ambigui, pot fi atât de numeroase, încât se poate spune chiar că astfel de texte nici nu sunt potrivite pentru o sinteză automată pornind de la text; singura soluție rezonabilă este, probabil, o reeditare a lor pentru a le face mai accesibile unui sistem de sinteză.

Problema enunțată anterior este de fapt mult mai generală. Părerea autorului acestei lucrări este că în orice aplicație TTS trebuie făcut un compromis între calitatea vorbirii sintetizate, dimensiunile vocabularului și complexitatea sistemului de sinteză. Cu alte cuvinte, nu trebuie încercat cu orice preț, prin orice mijloace, obținerea unei vorbiri "perfecte", cel puțin în acest moment.

4.2 Algoritm de preprocesare a textului

Pentru preprocesarea textului de intrare în cadrul sistemului TTS proiectat, am propus un set de definiții, reguli și proceduri, bazate pe o analiză detaliată a situațiilor cele mai întâlnite în limba română.

Definițiile propuse sunt prezentate în continuare.

- Definiția 4.1** Vom denumi *expresii* "secvențele de caractere care cuprind una sau mai multe din următoarele categorii: secvențe de litere dintre care cel puțin una este majusculă, secvențe de cifre, semne de punctuație, alte simboluri speciale".
- Definiția 4.2** Vom denumi *caractere extra-textuale* "acele semne de punctuație care îndeplinesc în text o funcție de punctuație propriu-zisă".
- Definiția 4.3** Vom denumi *caractere intra-textuale* "acele semne de punctuație care fac parte integrantă din expresii și ajută la pronunția lor".
- Definiția 4.4** Vom denumi *expandare* "procesul de conversie a unor expresii în format lexical (secvențe de litere alcătuind cuvinte uzuale, ce pot fi analizate lingvistic)".
- Definiția 4.5** Vom denumi o secvență de caractere *ambiguă* "dacă ea poate fi încadrată, având în vedere forma sa, în mai multe clase lingvistice".

Pornind de la aceste definiții, am proiectat un algoritm de preprocesare a textului, ce constă în principiu din trei etape de bază:

I. Segmentarea textului

Textul se segmentează de la stânga spre dreapta, în *grupuri de caractere*. Se obțin astfel secvențe de caractere ASCII delimitate de spații albe (blanc); semnele de punctuație se includ temporar în aceste grupe.

II. Conversia șirurilor de caractere de tip expresie în caractere ortografice

Se parcurg pe rând grupurile de caractere rezultate în urma segmentării și se realizează *expandarea* lor (acolo unde este cazul) sub forma unor cuvinte uzuale, pe baza unei analize contextuale simple la nivel de cuvânt sau segment de cuvânt și a unor dicționare de conversie în format lexical (pentru abrevieri și unele tipuri de acronime).

III. Interpretarea unor semne de punctuație

Se detectează și se memorează pozițiile unor *caractere extra-textuale* și a sfârșitului frazelor, pentru a fi folosite ulterior de modulele de analiză sintactică și prozodică.

Detaliind etapa I prezentată anterior și utilizând și definiția 4.1, putem observa că grupurile de caractere rezultate în urma segmentării textului de intrare pot fi de următoarele tipuri [10, 11, 12]:

a. Secvențe de litere alfabetice, scrise cu minuscule

- a1. Cuvinte uzuale;
- a2. Abrevieri scrise fără punct (de exemplu unități de măsură: *m*, *km*, *ms*).

b. Expresii

- b1. Cuvinte scrise cu o singură literă, majusculă: abrevieri (puncte cardinale: *E* – est, *V* – vest; simboluri chimice: *C* – carbon, *O* – oxigen; unități de măsură: *A* – amper, *V* – volt); cifre romane: *V* – cinci, *I* – unu etc.
- b2. Abrevieri scrise cu minuscule și puncte (*tel.* – telefon, *a.c.* – anul curent)
- b3. Secvențe de mai multe litere, scrise cu minuscule și inițială majusculă
 - b3.1. Cuvinte la început de frază;
 - b3.2. Nume proprii;
 - b3.3. Abrevieri scrise fără punct (de exemplu unități de măsură: *Hz*, *Mw*).
- b4. Secvențe de mai multe litere, scrise cu minuscule și o majusculă, pe altă poziție decât prima (unități de măsură: *mA*, *kV* etc.)
- b5. Secvențe de litere scrise cu mai mult de două majuscule, cu sau fără punct
 - b5.1. Acronime (NATO, S.R.L.);
 - b5.2. Abrevieri (P.S. – post scriptum);
 - b5.3. Unități de măsură (MHz, MByte);

- b5.4.** Cifre romane (VI, IX).
- b6.** Secvențe de cifre, scrise cu sau fără semne de punctuație
- b6.1.** Numere întregi;
- b6.2.** Numere zecimale;
- b6.3.** Numerale ordinale (al 2-lea);
- b6.4.** Ore și date;
- b6.5.** Numere de telefon.
- b7.** Semne de punctuație: . ? ! : ; ... , - / ' " () [] { }
- b8.** Simboluri speciale
- b8.1.** Simboluri matematice uzuale: + - * (sau ×) : (sau /) = < >
% ~
- b8.2.** Alte simboluri speciale: @ \$ &

Deoarece semnele de punctuație ridică cele mai serioase probleme, vom analiza în primul rând situațiile cele mai uzuale de apariție a lor (pe grupe de importanță), precum și soluțiile posibile de rezolvare a acestor situații. Vom discuta apoi câteva aspecte fundamentale legate de grupurile de cifre, abrevieri și acronime.

1. Punctul

Punctul (.) poate apare în abrevieri, acronime, numerale, sau poate semnifica sfârșitul unei fraze. Ambiguitățile create de punct sunt o problemă majoră pentru operația de preprocesare, datorită faptului că el poate reprezenta fie un caracter intra-textual, fie extra-textual, fie ambele în același timp; de exemplu, punctul după abreviere poate marca în același timp și sfârșitul frazei.

Este deosebit de utilă punerea în evidență a câtorva situații de utilizare corectă a punctului în limba română:

- Punctul se folosește în abrevierile provenite din cuvinte simple sau compuse în care **nu apare** litera finală a cuvântului; exemple: *id.* (**idem**), *etc.* (**etcetera**), *tel.* (**telefon**), *a.c.* (**anul curent**), *a.m.* (**ante meridian**), *d.a.* (**după-amiaza**), *P.S.* (**post scriptum**) – deci categoriile **b2**, **b5.2** puse în evidență anterior.
- Dacă în abreviere **apare** litera finală a cuvântului, nu se pune punct după abreviere; exemple: *cca* (**circa**), *dna* (**doamna**), *dl* (**domnul**), *dnei* (**doamnei**), *jr* (**junior**) – categoria **a2**.
- Nu se pune punct după simbolurile unor termeni de specialitate: *C* (carbon), *L* (lungime), *V* (volt sau volt), *mA* (miliamperi), *MHz* (mega hertzi) – categoriile **a2**, **b1**, **b3.3**, **b4**, **b5.3**.

- În acronime (abrevieri provenite din inițialele unor substantive compuse formate din mai mulți termeni), punctul este facultativ; sunt corecte atât formele *O.N.U.*, *S.U.A.*, cât și *ONU*, *SUA* (categoria **b5.1**).
- Nu se folosește punctul în abrevierile ce reprezintă indicative de state (*RO* – România), sau de județe (*CT* – Constanța) și în situațiile când abrevierea s-a transformat într-un cuvânt sudat, caracterizat prin lectură cursivă (*TAROM*) – categoria **b5.2**.
- Punctul se folosește de asemenea în scrierea unor numere și a datelor: numere întregi sau zecimale (*1.234*, *1.234,567*), date (*15.04.2002*) – categoriile **b6.1**, **b6.2**.

Considerațiile anterioare sugerează următoarea procedură: atunci când este detectat punctul într-un grup de caractere, se cercetează contextul în care apare și apoi se ia decizia corespunzătoare, astfel:

- Dacă există cifre la stânga și la dreapta, el este declarat caracter intra-textual și:
 - dacă mai există un punct în secvența de cifre, secvența reprezintă o dată și se expandează folosind un set de reguli (de exemplu: *15.04.2002* va deveni *cincisprezece aprilie două mii doi*);
 - dacă nu mai există un alt punct, secvența reprezintă un număr și se expandează folosind de asemenea reguli (de exemplu: *1.234* va deveni *o mie două sute treizeci și patru*).
- Dacă punctul este în poziție finală și este precedat de alte două puncte (...), această secvență se declară caracter extra-textual, fiind identificată cu semnul de punctuație corespunzător; acest caz îl vom discuta separat.
- Dacă punctul este precedat de o secvență de litere (minuscul sau majuscul) și eventual de alte puncte, se caută într-un dicționar de abrevieri și acronime și:
 - dacă grupul de caractere este găsit în dicționar, punctul este declarat caracter intra-textual și secvența se expandează conform echivalenței din dicționar;
 - dacă grupul de caractere nu este găsit în dicționar, dar conține majuscul, este un acronim – această situație o vom discuta separat;
 - dacă grupul de caractere nu este găsit în dicționar și nu conține majuscul și alte puncte, punctul (care este sigur în poziție finală) este declarat caracter extra-textual și va reprezenta sfârșitul unei fraze, poziția sa fiind memorată pentru modulele de analiză sintactică și prozodică.

Ultimele reguli prezentate nu pot însă elimina ambiguitatea situației în care punctul după o abreviere poate reprezenta în același timp și sfârșitul frazei (cazul lui *etc.* este tipic, dar există și numeroase alte exemple).

O soluție ar putea fi cercetarea grupului de caractere ce urmează după blank, ținând cont de faptul că la începutul unei noi fraze se află de regulă un cuvânt cu inițială majusculă. Această situație nu este însă complet edificatoare, deoarece în limba română majuscula apare ca inițială în multe cazuri: substantive nume proprii de persoană, nume de localități sau denumiri geografice, nume de planete și constelații, nume de instituții, nume de lucrări, nume de evenimente istorice sau de manifestări artistice și științifice, nume de sărbători, ca semn de respect etc.

Este clar că această ambiguitate nu va putea fi rezolvată numai de către preprocesor. Soluția pe care o propunem este următoarea:

- Dacă în urma cercetării contextului din dreapta rezultă că punctul din finalul unei abrevieri ar putea fi în același timp și sfârșitul frazei, punctul rămâne caracter intra-textual (și ajută la expandarea abrevierii), dar se adaugă un simbol special pentru marcarea provizorie a sfârșitului frazei, urmând ca acesta să fie validat sau nu de analiza sintactică ulterioară.

2. Semnele de punctuație ? ! : ; ...

Situațiile cele mai frecvente de apariție a lor sunt următoarele:

- Semnul întrebării (?) și semnul exclamării (!) se folosesc uzual în limba română la sfârșitul frazei. Ele apar foarte rar în interiorul frazelor, când pot reprezenta, de exemplu, considerații personale introduse în text, acestea fiind de obicei puse între paranteze; ca atare, cercetarea caracterului din dreapta lor (blank sau paranteză) poate diferenția simplu cele două situații.
- Semnele : și ; marchează și ele, de cele mai multe ori, finalul unui enunț. Deși nu constituie un sfârșit de frază propriu-zis, pot fi considerate în acest fel în contextul sintezei TTS, deoarece textele din partea stângă și din partea dreaptă se pot pronunța ca și cum ar fi izolate, fără să fie afectată naturalitatea pronunției.
- Prin urmare, cele patru semne menționate sunt importante în primul rând pentru modulul de analiză prozodică, deci locul lor trebuie detectat și memorat de către preprocesor, iar poziția în frază (finală sau intermediară) este utilă doar pentru a ușura analiza sintactică ulterioară a textului.
- Semnul ... semnifică faptul că se omite un anumit fragment de text (de exemplu finalul neprecizat al unei enumerări); el apare în mod obișnuit la sfârșitul unei fraze, dar poate apare și în poziție intermediară. Putem deci aplica aceeași regulă ca și pentru punctul final al unei abrevieri: cercetarea contextului din dreapta și, dacă este cazul, marcarea provizorie ca final de

frază, până la o analiză sintactică mai aprofundată; altfel, el nu modifică prozodia textului.

În toate situațiile menționate, semnele de punctuație vor fi interpretate drept caractere extra-textuale. Există însă și trei excepții, în care semnele ! și : au altă semnificație decât cea uzuală; aceste situații pot fi descrise de următoarele reguli:

- Dacă simbolul ! se găsește la finalul unei secvențe de numere, el semnifică cu mare probabilitate un "factorial" și va fi transcris ca atare.
- Dacă simbolul : se găsește în interiorul unei secvențe de numere, este considerat caracter intra-textual; secvența reprezintă o oră și se expandează folosind un set de reguli (de exemplu: *14:30* va deveni *ora paisprezece și treizeci de minute*).
- Dacă simbolul este înconjurat de blaturi, face parte dintr-o expresie matematică și va fi transcris conform dicționarului (*împărțit la*).

3. Virgula

Virgula (,) apare în mod uzual într-o frază în poziție intermediară, la finalul unui cuvânt, dar poate apare și în scrierea numerelor zecimale. Regula aplicată în cadrul algoritmului propus este următoarea:

- Se cercetează contextul în care apare virgula și:
 - dacă este înconjurată de cifre, se consideră caracter intra-textual; secvența reprezintă un număr zecimal și se expandează folosind un set de reguli (de exemplu: *1,234* va deveni *unu virgulă două sute treizeci și patru*).
 - dacă la stânga sa se găsește o literă sau un alt semn de punctuație (de exemplu punct după o abreviere), se consideră caracter extra-textual și poziția sa va fi memorată pentru modulul de analiză prozodică.

4. Cratima

Cratima (-) este un semn ortografic ce are în limba română două valori principale:

- *gramaticală*, atunci când servește la scrierea unor cuvinte compuse (*bună-cuviință, nord-vest, prim-plan, pare-mi-se, propriu-zis etc.*);
- *fonetică*, atunci când servește la marcarea pronunțării într-o singură silabă a două sunete din două cuvinte diferite, dar care se găsesc alăturate în vorbirea curentă (*de-a*).

În fapt, deoarece simbolurile uzuale folosite de calculator nu cuprind linii mediane de lungimi diferite, cratima devine practic un semn de punctuație și poate fi folosită atât

pentru scrierea cuvintelor compuse sau a unor numere ordinale, cât și pentru despărțirea în silabe, pentru delimitarea unui nou paragraf, sau în enumerări.

Determinarea caracterului intra sau extra-textual se poate face prin cercetarea contextului în care apare; ea este mărginită de obicei fie de litere, fie de blankuri, dar această informație este utilă doar pentru analiza sintactică, deoarece în mod uzual nu se citește (este suprimată de către preprocesor) și nu modifică prozodia textului. În numerele ordinale, expandarea se face simplu, pe bază de reguli (*al 2-lea – al doilea*).

5. Bara oblică

Bara oblică (/) are sensul prepoziției "pe" în abrevierile științifice (*km/h – kilometru pe oră, m/s – metru pe secundă*) și în exprimarea unei proporții (*2/3 – doi pe trei*), sau sensul conjuncției "sau" în textele uzuale (*c(e/i) – ce sau ci*); în ambele situații reprezintă un caracter intra-textual. De asemenea, poate semnifica o împărțire în expresiile matematice.

Regulile pe care le propunem pentru simbolul / sunt următoarele:

- Dacă este înconjurat de litere, grupul de caractere din care face parte se caută în dicționarul de abrevieri și:
 - dacă se găsește în dicționar, se transcrie *pe* și se folosește expresia completă găsită (*metru pe secundă*);
 - dacă nu este găsit în dicționar, se transcrie *sau*.
- Dacă este înconjurat de numere izolate, se transcrie *pe*.
- Dacă este înconjurat de secvențe de cifre și alte caractere matematice ($2 \times 3/4 \times 5$), sau de paranteze și secvențe de cifre ($(2+3)/(4+5)$), se transcrie *împărțit la*.

6. Apostroful

Apostroful (') este folosit în limba română în mai multe situații:

- pentru a reproduce în scris rostiri în care un sunet sau mai multe nu sunt pronunțate; aceste rostiri sunt însă rare, fiind practic neliterare, populare (*pân'deseară*);
- în nume proprii străine sau în neologisme neadaptate (*O'Neill, five o'clock*);
- în scrierea anilor, fără prima sau primele cifre (*'907, '99*).

Regulile pe care le propunem pentru simbolul ' sunt următoarele:

- Dacă se găsește într-o secvență de litere, el este eliminat (nu reprezintă propriu-zis un caracter intra-textual și nu ajută la pronunția cuvântului).

- Dacă în dreapta se găsește o secvență de cifre, în funcție de numărul acestor cifre, grupul de caractere se expandează folosind un set de reguli (de exemplu: '99 va deveni *o mie nouă sute nouăzeci și nouă*).

7. Alte semne de punctuație: " () [] { }

Alte semne de punctuație ce pot fi utilizate în textele obișnuite sunt ghilimelele (sau semnele citării) și parantezele rotunde; ele semnifică de obicei un citat, reprezintă porțiuni de text cărora li se dă un sens (stilistic) special sau asupra cărora autorul vrea să insiste, constituie traducerea ori sensul unui cuvânt, sau delimitează considerații personale introduse în text. Apar de obicei în perechi și vor fi declarate caractere extra-textuale, servind modulului de analiză prozodică pentru obținerea unei vorbiri sintetizate cât mai naturale.

Parantezele drepte și acoladele apar extrem de rar în textele românești uzuale; ele pot apare însă (ca și parantezele rotunde) în expresii matematice. Se identifică simplu, deoarece sunt alăturate unor secvențe de cifre și se expandează de obicei prin utilizarea cuvintelor corespunzătoare semnificației lor, cu ajutorul dicționarului de conversii în format lexical.

8. Secvențele de cifre

Secvențe de cifre pot apare și în texte obișnuite, dar mai ales în expresii matematice, împreună cu semne de punctuație sau simboluri matematice; evident, deoarece numărul lor posibil este practic infinit, ele trebuie expandate pe bază de regulii de conversie, în funcție de context.

Am propus anterior o serie de regulii pentru cazurile cele mai frecvente (numere întregi sau zecimale, numerale ordinale, ore, date). O situație specială (pe care de asemenea am menționat-o anterior), o reprezintă cazul în care o secvență de cifre, scrisă fără semne de punctuație, poate reprezenta fie un număr întreg, fie un număr de telefon. În acest caz, dacă din cercetarea contextului nu se poate elimina ambiguitatea (de exemplu prezența abrevierii *tel.*), această problemă rămâne în sarcina modulului de analiză sintactică, care poate realiza o cercetare contextuală mai amplă.

9. Simbolurile matematice uzuale: + - * (sau ×) : (sau /) = < > % ~

Simbolurile matematice au o situație oarecum privilegiată, deoarece ele sunt încadrate de obicei de blankuri în expresiile matematice uzuale și ca atare pot fi imediat identificate și expandate pe baza dicționarului de conversii în format lexical (de exemplu *plus, minus, înmulțit cu, împărțit la* etc.) Dacă totuși în scrierea expresiei nu apar blankuri, contextul secvențelor de cifre și al celorlalte simboluri duc practic la aceeași rezolvare.

10. Abrevierile

O serie de considerații privind abrevierile au fost expuse anterior la regulile ce privesc punctul. Situația lor este dificilă datorită faptului că în limba română abrevierile se pot scrie în multe feluri: cu majuscule și/sau minuscule, cu sau fără semne de punctuație (uzual punct).

Regula principală ce poate fi aplicată este evidentă:

- Dacă în grupul de caractere apare cel puțin un punct și/sau cel puțin o majusculă, se caută în dicționarul de abrevieri; dacă secvența este găsită, abrevierea se expandează punând-o în corespondență cu cuvântul corespunzător din dicționar.

Pot rămâne însă ambiguități, în special pentru abrevierile scurte (de exemplu *V* – unitatea de măsură "volt", dar și cifra romană "cinci" și punctul cardinal "vest"), sau pentru abrevierile scrise cu minuscule și fără punct (*km, cca, dl*), acestea din urmă nefiind căutate în dicționar (după regula expusă). Singurele soluții practice pentru rezolvarea unor astfel de cazuri ambigui este ca ele să fie preluate mai departe de analiza sintactică sau să fie recunoscute la etapa de conversie fonetică, prin căutarea într-un dicționar limitat de excepții.

11. Acronimele

Spre deosebire de abrevieri, cea mai mare parte a acronimelor nu trebuie stocate în dicționar, deoarece pronunția lor nu necesită informații textuale suplimentare. De obicei, pronunția lor se reduce la citirea secvențială a caracterelor ce compun acronimul, individual (ca pentru *S.R.L.*), sau la citirea normală a cuvintelor, atunci când pronunția lor s-a generalizat în limbaj într-o formă compactă (*NATO, TAROM*); pentru citirea secvențială a acronimelor, este necesar doar un set de reguli simple de transcriere a literelor rostite separat (de exemplu *S.R.L. – serele*).

Regula propusă pentru acronime este deci următoarea:

- Dacă secvența de caractere cuprinde cel puțin două majuscule și nu este găsită în dicționarul de abrevieri, se caută în dicționarul de acronime:
 - dacă este găsită aici, secvența se expandează conform echivalenței din dicționar;
 - dacă nu este găsită în dicționarul de acronime și nu cuprinde puncte, majusculele sunt (eventual) înlocuite cu minuscule și secvența nu va suferi altă prelucrare (se va citi ca atare);
 - dacă nu este găsită în dicționar, dar cuprinde puncte, secvența este expandată secvențial, utilizând un set minim de reguli de transcriere a literelor rostite separat.

Pentru toate situațiile menționate, preprocesorul va semnaliza acronimul modulelor ulterioare, pentru o corectă analiză sintactică și prozodică a textului.

5. Concluzii

Am discutat în această lucrare câteva aspecte fundamentale legate de sinteza automată a vorbirii, ca și un număr important de reguli și principiile generale pe baza cărora a fost proiectat preprocesorul de text pentru sistemul TTS în limba română. Nu am urmărit totuși să descriem complet și în detaliu funcționarea și implementarea acestuia; o serie de considerații suplimentare și totodată modalitatea concretă de implementare (pentru o variantă preliminară) au fost prezentate de autor în [13] și [14].

În varianta actuală, preprocesorul de text a fost îmbunătățit pentru a rezolva unele situații dificile legate de abrevieri, numerale urmate de unități de măsură etc. De asemenea, un mecanism de automat de corecție permite preprocesorului să fie "tolerant" cu anumite erori tipice de sintaxă, cum ar fi de exemplu fraze ce nu încep cu minuscule, sau un format "ușor" incorect pentru date sau numerale.

Putem spune, ca o concluzie a celor discutate anterior, că un preprocesor de complexitate medie, cum este și cel propus pentru sistemul TTS în limba română, poate rezolva cu succes (împreună cu analiza lingvistică ulterioară) o mare parte din problemele întâlnite într-un text obișnuit; el nu poate realiza însă normalizarea completă a **oricărui** text și nu poate soluționa **toate** ambiguitățile care se pot ivi, datorate în special numărului extrem de mare al abrevierilor, acronimelor – în general a secvențelor neuzuale care pot apare într-un text scris. De asemenea, nu poate face față unor construcții negramaticale (deși, de exemplu, unele simboluri speciale neașteptate sunt ignorate).

Desigur că un set mai mare de reguli și un dicționar de conversii în format lexical mai cuprinzător ar spori eficiența preprocesorului, dar este posibil ca el să devină atât de complicat, încât să fie practic neoperațional pentru un sistem TTS. Singura soluție practică pentru tratarea cazurilor ambigui este folosirea unui set minim de reguli, păstrarea în dicționar a celor mai uzuale situații (cu posibila adaptare a dicționarului la **tipul** textului ce se citește) și examinarea cazurilor rămase la un nivel superior, pe baza plauzibilității sintactice, semantice sau pragmatice a frazelor obținute după preprocesare.

Referințe bibliografice

- [1] R. Boite, H. Bourlard, T. Dutoit, J. Hancq, H. Leich (2000). *Traitement de la parole*. Presses Polytechniques et Universitaires Romandes, 2000.
- [2] G. Bailly (1996). Pistes de recherches en synthèse de la parole – în "Fondements et perspectives en traitement automatique de la parole" (H. Méloni – Coord.), Aupelf-Uref, pp. 109-120, 1996.
- [3] D. Burileanu (1999). Contribuții privind sinteza automată a vorbirii pornind de la text în limba română – Teză de doctorat. Universitatea "POLITEHNICA" București, 1999.
- [4] D. Burileanu (2002). Basic Research and Implementation Decisions for a Text-to-

- Speech Synthesis System in Romanian Language – Lucrare în curs de publicare în "International Journal of Speech Technology", Kluwer Academic Publishers, 2002.
- [5] G. Fries, A. Wirth (1997). FELIX - A TTS System with Improved Preprocessing and Source Signal Generation – Comunicare la "EUROSPEECH'97", Rodos, pp. 589-592, 1997.
- [6] E. Lewis, M. Tatham (1993). A Generic Front-End for Text-to-Speech Synthesis Systems – Comunicare la "EUROSPEECH'93", Berlin, vol. 2, pp. 913-916, 1993.
- [7] M.Y. Liberman, K.W. Church (1992). Text Analysis and Word Pronunciation in Text-to-Speech Synthesis – în "Advances in Speech Signal Processing" (S. Furui, M. Sondhi – Coord.), Dekker, pp. 791-832, 1992.
- [8] A. Lindstrom, M. Ljungqvist (1994). Text Processing within a Speech Synthesis System – Comunicare la "International Conference on Spoken language Processing", Yokohama, pp. 139-142, 1994.
- [9] M. McAllister (1989). The Problems of Punctuation Ambiguity in Full Automatic Text-to-Speech Conversion – Comunicare la "EUROSPEECH'89", Paris, pp. 538-541, 1989.
- [10] G. Beldescu (1984). Ortografia actuală a limbii române. Ed. Științifică și Enciclopedică, București, 1984.
- [11] T. Dutoit (1997). An Introduction to Text-to-Speech Synthesis. Kluwer Academic Publishers, 1997.
- [12] F. Șuteu, E. Șoșa (1993). Dicționar Ortografic al Limbii Române. Ed. ATOS, București, 1993.
- [13] D. Burileanu (1999). Natural Language Processing for Speech Synthesis in Romanian Language –Comunicare la "The 12th International Conference on Control System and Computer Science", București, vol. II, pp. 1-6, 1999.
- [14] D. Burileanu, C. Dan, M. Sima, C. Burileanu (1999). A Parser-Based Text Preprocessor for Romanian Language TTS Synthesis – Comunicare la "EUROSPEECH'99", Budapesta, vol. 5, pp. 2063-2066, 1999.

Utilizarea tehnicilor nuanțate (fuzzy) și de dinamică neliniară pentru sinteza adaptivă a vorbirii

Horia-Nicolai L. TEODORESCU
Academia Română, Secția Știința și Tehnologia Informației,
Calea Victoriei 125, București
E-mail: hteodor@etc.tuiasi.ro

1. Introducere

În timp ce mașina realizează tipic transmisie de date, omul comunică. Diferența constă în participarea intelectuală și afectivă a persoanei la actul comunicării, participare reflectată atât la nivelul limbajelor neverbale (gestică, mimică etc.), cât și la nivelul vocal. Această participare afectivă dă varietate, coloratură și sensuri suplimentare, nu neapărat pe plan semantic, semnalului vocal. Sinteza vocii, în prezent, este limitată de lipsa afectului, varietății și sensurilor suprapuse în planuri multiple. Vocea mașinii rămâne astfel cantonată într-o regiune “moartă” a comunicării, este monotună și obositoare pe termen lung.

În această lucrare, reluând unele idei din [1-12], precum și în contextul unor dezvoltări recente [13-27], în special legate de e-Voice și VXML, prezentăm și dezvoltăm unele concepte și tehnici care ar putea permite mașinii atingerea dezideratelor mai sus menționate. Realizarea unor mașini capabile să mimeze calitățile vocii umane și să *dialogheze* cu oamenii, sau măcar să comunice într-o manieră similară în care omul o face, este un deziderat în numeroase domenii, de la dialogul om-calculator, la sistemele auto și la sistemele de învățare asistată de calculator [13-15]. Rezolvarea acestei probleme are implicații semnificative pentru acceptarea sintezei vocii într-o varietate de aplicații, de la robotică la realitate virtuală, la industria de jocuri electronice și la protezare.

Prozodia, adică structura acustică ce se extinde pe mai multe segmente de semnal vocal, chiar peste mai multe cuvinte sau propoziții, implică ritm, accent, intonație, timbru, afect și alte caracteristici ale vocii încă insuficient înțelese sau vag definite în literatură. Informația paralingvistică ce este conținută de prozodie nu este nicăieri regăsită la nivelul “spus” prin cuvinte, dar – așa cum am subliniat în [2] – această informație poate fi chiar mai importantă pentru ascultător decât informația lingvistică propriu-zisă. Incapacitatea sistemelor actuale de sinteză vocală de a reda prozodia naturală este evidențiată chiar de marii producători de aplicații [25] și este bine cunoscută în mediul cercetătorilor în domeniul sintezei vorbirii: “*One of the most difficult problems in speech to date is prosodic modeling*” [25].

2. Soluții pentru sinteza adaptivă și varietală

Cele două calități ale vocii naturale, adaptivitatea – în sens larg – și variabilitatea se pot realiza, cu costuri nu neapărat mari, la nivelul sintetizoarelor actuale, cu adaptări minimale (sau deloc) la nivel hardware și cu îmbunătățiri ale programelor de control. Sinteza adaptivă se referă la adaptarea la:

- Condițiile sonore ambientale [1, 4];
- Contextul semantic-afectiv al cuvintelor și frazelor sintetizate [2, 3].
- Interlocutorul sistemului de sinteză automată, atunci când acesta este recunoscut [2].

Sinteza varietală se referă la modificările inter-pronunție, la repetarea unor fraze, chiar și în cazul în care condițiile ambientale și contextul (și interlocutorul) rămân neschimbate. Această variabilitate elimină monotonia și personalizează vocea (naturală sau sintetizată), în măsura în care variabilitatea se face după reguli adesea proprii individului (cum este cazul în realitate) – și nu doar aleatoare.

Variabilitatea intrinsecă a vorbirii derivă din mecanismele fizice de producere a semnalului vocal (curgere turbulentă a aerului prin organul fonator), precum și din mecanismele neurologice de control al producerii semnalului vocal (controlul neuronal este cunoscut ca având o dinamică cu o importantă componentă neliniară). Aceste caracteristici au fost documentate de mai multe grupuri de cercetare, inclusiv de noi și colaboratorii [5-9].

Adaptabilitatea și variabilitatea în sensurile de mai sus vor fi prezentate sumar în secțiunile următoare, sintetizând lucrările citate și unele cercetări mai noi, nepublicate încă.

3. Adaptabilitate la mediu

Una dintre cele mai elementare adaptări ale semnalului vocal generat de om este cea de adaptare la condițiile de mediu. Adaptarea la un mediu real, cu fond de zgomot, se realizează pe patru căi principale: prin modificarea amplitudinii semnalului (mai mare în mediul de zgomot ridicat), prin modificarea spectrului (crește contribuția frecvențelor înalte), prin modificarea ritmului (scăderea ritmului, creșterea duratei vocalelor), și prin creșterea duratei dintre cuvinte, care devin separate, segmentate în timp. Adaptările realizate – instinctiv de un vorbitor uman – se operează deci la un nivel relativ elementar, cu modificări de prozodie minimale.

Realizarea acestei adaptări este esențială în multe aplicații de sinteză a vocii, incluzând sinteza vocală pentru aplicații în medii industriale și în mijloace de transport, sau sinteza vocală pentru proteze laringiene. Este remarcabil că această adaptare se poate realiza, la pretenții reduse, cu foarte puțin hard suplimentar și/sau cu un soft minimal, aducând însă o îmbunătățire esențială în utilizare. În privința hardului, este necesar unul sau mai multe canale de culegere a semnalului de zgomot (semnal sonor ambiental).

Procesarea semnalului de zgomot, în vederea realizării controlului sistemului de sinteză automată, presupune determinarea puterii zgomotului ambiental într-o fereastră temporală și determinarea componenței spectrale a semnalului ambiental. Primul parametru de caracterizare a zgomotului se obține ca medie aritmetică a pătratului semnalului s , într-o fereastră dată, de lărgime de W eșantioane și caracterizată de momentul actual de timp, n :

$$P(t) = \sum_{k=0}^W s_{n-k}^2 \quad (1)$$

Caracterizarea spectrală se poate realiza sumar prin raportul HL dintre puterea la frecvențe “înalte” (frecvențele înalte corespunzând în mare benzii de frecvență ce include formanții nr. 2, 3, 4 și 5 din spectrul vocal) și puterea la frecvențele “joase” (până la aproximativ al doilea formant, deci până la frecvența de cca. 400–500 Hz, ținând cont și de vorbitorii feminini):

$$HL = \frac{\int_0^{500} S^2(\omega) \cdot d\omega}{\int_{500}^{10000} S^2(\omega) \cdot d\omega} \quad (2)$$

Deoarece parametrii respectivi sunt relaționați cu impactul pe care îl au asupra inteligibilității vorbirii, deci sunt dați de calități subiective, este natural să abordăm o definiție probabilistă sau fuzzy a lor. Dată fiind simplitatea controlului nuanțat⁹³ (fuzzy), vom prefera a doua variantă. Un exemplu de definiție⁹⁴ a funcțiilor de apartenență respective este prezentat în Figura 1. Este de presupus ca această definiție să constituie doar un punct de plecare, îmbunătățirea calității sintezei realizându-se și prin modificarea funcțiilor de apartenență.

⁹³ Deși nu este larg acceptat și are o traducere mai dificilă în alte limbi, vom utiliza aici termenul “nuanțat”, propus de Grigore C. Moisil, în locul englezescului “fuzzy”.

⁹⁴ Pentru a nu încălca prezentarea, ecuațiile funcțiilor respective sunt date în Anexa 1.

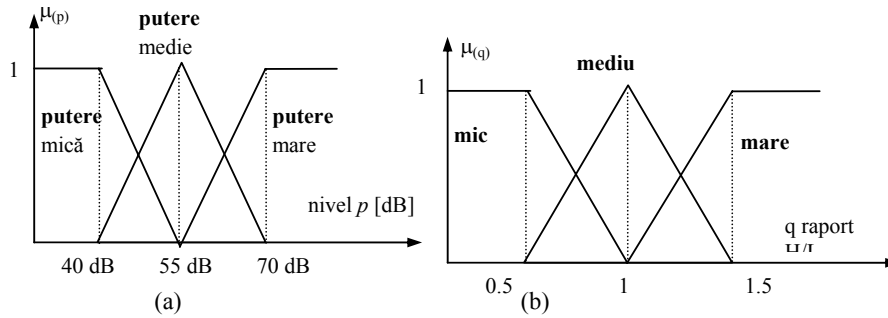


Figura 1. Funcțiile de apartenență ale premiselor regulilor folosite pentru determinarea modificărilor parametrilor de control ai sintetizorului

După cum s-a precizat deja, ca rezultat al aprecierii condițiilor de mediu, se controlează patru parametri ai semnalului sintetizat:

- creșterea amplitudinii (parametru notat AI)
- creșterea conținutului în frecvențe înalte (HFICI)
- creșterea duratei vocalelor (VLI)
- creșterea duratei dintre cuvinte (accentuarea segmentării pe cuvinte a frazei), notat IDBBW.

Controlul se realizează pe bază de reguli și poate fi rezumat în Tabelele 1-4 de mai jos.

Tabelul 1

Creșterea amplitudinii (AI – *Amplitude Increase*)

HL/P	mic	mediu	mare
mic	0,0	0,1	0,4
mediu	0,1	0,3	0,5
mare	0,4	0,5	0,6

Tabelul 2

Creșterea conținutului de frecvențe înalte (HFICI – *High Frequency Content Increase – F3 increase*)

HL/P	mic	mediu	mare
mic	0,0	0,1	0,4
mediu	0,1	0,3	0,5
mare	0,4	0,5	0,6

Tabelul 3

Creșterea duratei vocalelor (*Vowel Length Increase – VLH*)

HL/P	mic	mediu	mare
mic	0,0	0,1	0,4
mediu	0,1	0,3	0,5
mare	0,4	0,5	0,6

Tabelul 4

**Creșterea duratei dintre cuvinte
(*Increase of the Duration of the Break Between Words – DBBW*)**

HL/P	Mic	mediu	mare
mic	0,1	0,1	0,4
mediu	0,1	0,3	0,5
mare	0,4	0,5	0,6

Tabelele sunt interpretate în sensul uzual pentru sistemele nuanțate. Preferăm sistemele de tip Sugeno de ordin 0 (vezi Anexa 1), deoarece furnizează ca rezultat, direct, valori numerice, care vor fi interpretate ca și coeficienți de multiplicare ai valorilor nominale ale sintezei. De exemplu, prima linie și prima coloană din Tabelul 1 spun că:

DACĂ Puterea (zgomotului) este **medie** și parametrul LH este **mediu**

ATUNCI Amplitudinea crește cu **0.3** ori.

Toate regulile din Tabelul 1 și toate celelalte tabele se interpretează într-un mod similar.

Rezultatul final se obține prin agregarea rezultatelor parțiale, date de regulile respective. De exemplu, dacă valoarea intensității sonore este de 45 dB, iar raportul HL este de 0,7, prin aplicarea fuzificării⁹⁵ se obține gradul de adevăr al premisei (combinate) din regula respectivă, prin

$$\min(\mu_{putere=mica}(P_0), \mu_{LH=mic}(LH_0))$$

unde $P_0 = 45$ iar $LH_0 = 0,7$. Folosind expresiile funcțiilor (v. Anexa 1), se obțin valorile $\mu_{putere=mica}(P_0) \approx 0,67$, $\mu_{LH=mic}(LH_0) = 0,6$, deci valoarea minimă este 0,6 și reprezintă gradul de încredere în faptul că amplitudinea crește de 1,1 ori. Aceasta este valoarea de adevăr pentru singletonul (de la ieșirea sistemului) ce corespunde regulii respective, $\alpha_{1,1}^A$. În total, sunt 9 reguli per tabel, deci există 9 valori de singletoni. Într-adevăr, în același timp, valorile de intrare corespund funcțiilor de apartenență „mediu” pentru „putere” și LH, deci regulii:

DACĂ Puterea (zgomotului) este **mică** și parametrul LH este **mic**

⁹⁵ Termenul echivalent românesc ar fi “nuanțare”

ATUNCI *Amplitudinea crește cu 0,0 ori.*

cu gradul de încredere în rezultat:

$$\min(\mu_{putere=medie}(P_0), \mu_{LH=medie}(LH_0))$$

precum și regulilor:

DACĂ *Puterea (zgomotului) este mică și parametrul LH este mediu*

ATUNCI *Amplitudinea crește cu 0,1 ori.*

respectiv:

DACĂ *Puterea (zgomotului) este medie și parametrul LH este mic*

ATUNCI *Amplitudinea crește cu 0,1 ori.*

cu gradele de încredere

$$\min(\mu_{putere=mica}(P_0), \mu_{LH=medie}(LH_0))$$

și respectiv

$$\min(\mu_{putere=medie}(P_0), \mu_{LH=mic}(LH_0))$$

Celelalte cinci reguli din Tabelul 1 au gradele de încredere în rezultat nule, deoarece valorile funcțiilor de apartenență „mare” ale premiselor („puterea este mare” și „LH este mare”) sunt nule, pentru valorile date, $P_0 = 57$ și $LH_0 = 0,7$.

Prin agregare (defuzzificare), considerată aici conform formulei uzuale:

$$y = \frac{\sum_{k=1}^9 \alpha_k^A \mu_k^A(x_0)}{\sum_{k=1}^9 \mu_k^A(x_0)} \quad (3)$$

se obține valoarea de ieșire (amplitudinea, creșterea conținutului de frecvențe înalte, creșterea lungimii vocalelor, respectiv creșterea duratei pauzei dintre cuvinte). În relația de mai sus, α_k^X reprezintă abscisele singletonilor de ieșire din sistemele tip Sugeno respective, $\mu_k^X(\cdot)$ reprezintă gradele de încredere în concluzia regulilor respective, iar y reprezintă valoarea agregată (defuzzificată) de ieșire a sistemului Sugeno. Sumarea se face pentru toți singletonii de ieșire (notați de la 1 la 9). Indicele „A” arată că ne referim la parametrul controlat „amplitudine”, controlul fiind desigur diferențiat pentru cei patru parametri discutați.

Valorile astfel obținute sunt folosite, cum s-a precizat, ca factori de multiplicare ai parametrilor nominali⁹⁶. De exemplu, dacă amplitudinea nominală este A_0 , atunci, prin aplicarea controlului, amplitudinea efectivă a semnalului va fi:

$$A = A_0 \cdot \left(1 + \frac{\sum_{k=1}^9 \alpha_k^A \mu_k^A(x_0)}{\sum_{k=1}^9 \mu_k^A(x_0)} \right) \quad (4)$$

Sistemul de control este instantaneu, în sensul că nu ține cont decât de valorile recente (din fereastra prezentă, de lărgime W) ale zgomotului, nu și de valorile anterioare. Controlul de amplitudine și frecvență se poate exercita în afara sintetizorului propriu-zis, asupra unui amplificator și a unui filtru plasate la ieșirea sintetizorului. Aceste două controale se pot prevedea de altfel și în alte aplicații, precum sisteme de sonorizare mari (eventual distribuite, ca în cazul sonorizării unor spații mari, gen piețe sau stadioane), sau a unor sisteme de sonorizare locale (de exemplu, sisteme de interfonie). Controlul pauzelor dintre cuvinte și un control fin al spectrului vocalelor necesită comanda directă a sintetizorului.

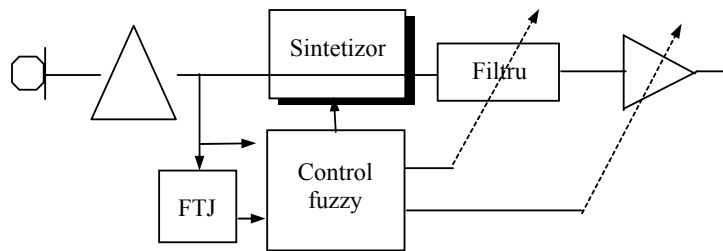


Figura 2. Schema bloc a unui sistem audio adaptiv la zgomotul ambiental

În cazul în care se utilizează doar primele două tipuri de adaptare, în amplitudine și spectral, adaptarea se poate realiza și cu mijloace hardware externe sintetizorului, putând, de altfel, fi utilizată în orice aplicație audio (de sonorizare etc.). Schema unui asemenea sistem de adaptare este cea prezentată în Figura 2, o variantă fiind inițial propusă în [4].

⁹⁶ *Nominali, în sensul că sunt valorile standard pentru sistemul de sinteză automată respectiv și pentru sunetul respectiv produs în condițiile contextuale date.*

4. Adaptare și variabilitate contextual-interpretativă

Interlocutorul uman răspunde cu afect, după cum consideră anormală, nepotrivită, sau oricum în alt fel “departe de așteptări” întrebarea sau afirmația făcută de partenerul la dialog. De asemenea, răspunsul este diferit atunci când vorbitorul uman este nesigur de răspuns, are un interes special în răspuns sau în topica discuției, sau, din contra, este dezinteresat. În plus, situarea interlocutorului față de partenerul sau partenerii de dialog, în context social sau afectiv, tonalizează discursul verbal și îi imprimă specificitate relativă. Toate aceste caracteristici participative, precum și altele asemenea, dau *comportamentul verbal* al omului, sunt traduse în mare măsură la nivelul semnalului vocal prin prozodie, dar în prezent nu se regăsesc la nivelul mașinii. Privitor la elementele de bază privind prozodia, vezi [26].

Pentru a implementa un comportament verbal, mașina trebuie să dispună de o bază de cunoștințe minimală prin care să genereze acest comportament. De exemplu, este necesar să se interpreteze “departe de normal” într-o aserțiune sau întrebare a interlocutorului uman. Deci, vom presupune că există o bază de cunoștințe care permite o asemenea interpretare. Construcția acestei baze de cunoștințe depinde de domeniul în care se poartă dialogul. În aceste condiții, accentul va fi mai puternic pe anumite părți ale frazei, sau răspunsul va depinde de aserțiune sau întrebare. Modul de răspuns va fi dirijat de asemenea de o bază de cunoștințe, care include regulile necesare modificării sintezei (vezi Figura 3).

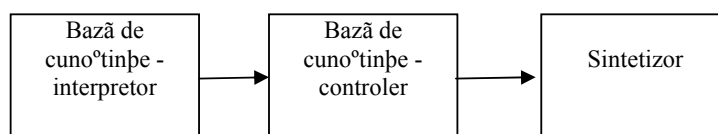


Figura 3. Schema de principiu a controlului contextual-interpretativ

Baza de cunoștințe-controler poate de asemenea fi implementată cu reguli *Dacă... Atunci*, de exemplu, de forma:

DACĂ oferta / răspunsul interlocutorului este neașteptat (negăsit în baza de cunoștințe – baza de așteptare/ baza de cazuri),

ATUNCI afectul sintezei este mirare / neîncredere/.../ etc.

ori

DACĂ oferta / răspunsul interlocutorului este neașteptat negativ (conform bazei de cunoștințe),

ATUNCI afectul sintezei este mirare și/sau furie.

Folosind rezultatele regulilor de acest fel, se pot seta parametrii ierarhic inferiori, de tonalitate, ai vocii sintetizate, pe baza acestora generându-se parametrii efectivi de control ai sintezei (amplitudine, frecvențe formanți etc.).

Deși acest gen de control poate părea complicat, sunt situații destul de generale în care el se poate implementa cu un efort relativ redus. De exemplu, atunci când se determină (printr-o măsurătoare relativ simplă, de frecvență medie în spectrul vocal, sau de fundamentală) că interlocutorul este un copil sau o persoană de gen feminin, se poate selecta una sau ambele dintre alternativele:

- sistemul de sinteză automată se setează pe o voce de același tip (copil/feminin)
- sistemul de sinteză automată se setează pe voce “caldă” și “vorbire clară”.

Utilitatea și modalitatea de realizare a primei setări nu necesită explicații. A doua setare (care poate fi simultană cu prima) se justifică – în cazul interlocutorului copil – prin necesitatea de a îi crea un mediu afectiv propice și liniștit de dialog (voce “caldă”) și prin necesitatea unei comunicări cât mai informative, ușor de urmărit. Pentru a obține o voce “caldă”, se pot folosi trasee melodice cu variații lente precum și frecvențe mai joase ale formanților și lărgimi mai mari (în zona spre frecvențe joase) a spectrelor formanților. “Claritatea” vocii se poate traduce prin segmentarea mai pronunțată pe cuvinte, precum și vocale mai lungi (cu sau fără accentuări ale spectrelor formanților). Utilizarea unor asemenea adaptări – ce rămân în mare măsură să fie concepute în detaliu, implementate și testate – este neîndoiește mare la sinteza pentru procese educative [15, 26], în aplicații medicale (răspuns sintetic destinat pacienților), precum și în numeroase aplicații generale (de exemplu, sintetizoare utilizate în muzee, pentru prezentarea exponatelor).

Alte modalități de personalizare afectivă sunt colorarea frecvențială și în amplitudine a anumitor părți din frază sau în cadrul unui cuvânt, aceste modificări locale fiind larg documentate în literatură, de ex. [16-18] și fiind relativ ușor de implementat.

5. Variabilitate prin metoda modulării de către un sistem dinamic neliniar

Variabilitatea semnalului vocal uman este bine cunoscută [5-9], [19-26]. Variabilitatea de tip natural a semnalului vocal sintetizat se poate obține prin modularea diverselor controale (al amplitudinii, lungimii vocalelor, accentului, pitch-ului etc.) sau semnale lent variable, generate de sisteme care prezintă dinamică neliniară (haos). Parametrii sistemului haotic respectiv pot modela un anume subiect; considerăm aici că acești parametri reprezintă individul vorbitor și “personalitatea” lui. Această metodă, propusă de noi inițial în 1992 ([28] ș.a.), dar nepublicată în forma extinsă, credem că reprezintă o metodă promițătoare de “personalizare” a vocii.

Considerăm un sistem dinamic neliniar, dependent de parametri; semnalul în timp generat de acesta este de forma $x(t) = x(t, \lambda_1, \lambda_2, \dots, \lambda_q)$, unde λ_h reprezintă parametrii sistemului haotic și permit modelarea specificității vorbitorului. Semnalul x poate fi folosit în modularea amplitudinii, frecvenței fundamentale, sau spectrului semnalului vocal sintetizat. De exemplu, spectrul poate fi modificat folosind o lege de variație a frecvenței centrale a formanților de forma:

$$f_j(t) = (1 + x_j(t)) \cdot f_{j0}(t) \quad (5)$$

unde $f_j(t)$ este frecvența formantului numărul j la momentul t , $x_j(t)$ este semnalul haotic respectiv ($x_j(t) \ll 1$), iar $f_{j0}(t)$ este frecvența “nominală” a formantului respectiv.

Un exemplu simplu de sistem haotic ce poate fi folosit în acest scop este dat de ecuațiile:

$$\begin{cases} r_{n+1} = \lambda_3 \cdot u_n^3 + \lambda_2 \cdot u_n^2 + \lambda_1 \cdot u_n + \lambda_0 \\ u_n = \lambda_4 \cdot r_n + \lambda_5 \end{cases} \quad (6)$$

unde setul de coeficienți $(\lambda_0, \lambda_1, \dots, \lambda_5) \in \mathbf{R}^6$ se alege în domeniul de valori ce corespunde unui comportament haotic al sistemului (vezi Anexa 2). Setul de coeficienți $(\lambda_0, \lambda_1, \dots, \lambda_5)$ se poate seta specific pentru fiecare sistem de sinteză automată, “personalizând” sistemul. Valorile de ieșire ale generatorului se scalează corespunzător și se folosesc la modularea unuia dintre parametrii de sinteză. Pentru exemplul din secțiunea 3, amplitudinea semnalului sonor devine, prin utilizarea modulației haotice:

$$A_n = A_0 \cdot \left(1 + \frac{\sum_{k=1}^9 \alpha_k^A \mu_k^A(x_0)}{\sum_{k=1}^9 \mu_k^A(x_0)} \right) \cdot (1 + \kappa \cdot r_n) \quad (7)$$

unde κ este un coeficient de scalare a seriei de timp r_n . Coeficientul κ se alege astfel încât contribuția termenului $\kappa \cdot r_n$ să fie de ordinul procentelor ($\kappa \cdot r_n < 0,1 \forall n$).

Desigur, scara de timp a procesului de generare de eșantioane de semnal vocal diferă de scara de timp a proceselor haotice folosite în modulație, ceasul celui de al doilea proces fiind mult mai lent (de ordinul 1/100) decât al primului proces. Pentru evitarea tranzițiilor bruște ale parametrului controlat, valorile generate pot fi interpolate și se poate realiza o variație lentă între două valori succesive. Considerând că un eșantion al seriei

haotice r_n este generat la fiecare Q eşantioane de semnal vocal, seria r_n se poate înlocui cu seria (mai “fină”, după ceasul de generare a eşantioanelor semnalului vocal):

$$r_k = r_{n-1} + \frac{r_n - r_{n-1}}{Q} \cdot k, \quad k = 0, 1, \dots, Q \quad (8)$$

În scopul modulării haotice a mai multor parametri de sinteză (amplitudine, frecvența centrală a formanților, lărgimea formanților, elemente prozodice etc.), sunt necesare mai multe generatoare haotice, câte unul pentru fiecare parametru controlat. Alternativ, se poate folosi un sistem nuanțat (fuzzy) haotic, aceste sisteme generând simultan un număr mare de ieșiri necorelate sau slab corelate [28].

6. Concluzii și discuții

Adaptabilitatea și variabilitatea sistemelor de sinteză a vocii și ale celor audio, în general, se pot asigura prin modificări relativ simple hard și soft ale sistemelor actuale. Adaptabilitate se poate manifesta atât în raport cu mediul sonor, cât și în raport cu contextul sau cu interlocutorul. Ideea de adaptabilitate și metodele respective au fost introduse de noi în urmă cu peste 20 de ani și dezvoltate continuu în lucrările citate, atât pentru aplicații de uz general, cât și pentru aplicații medicale.

O aplicație de interes medical-educațional este utilizarea unor sisteme de învățare a unei limbi pentru copii de vârste mici (1 lună – 3 ani) care suferă de deficiențe de auz. Utilizarea unor sintetizoare cu spectru și amplitudine controlate, astfel încât să fie optim adaptate auzului (curbei de sensibilitate audiometrică) a fiecărui copil în parte ar ajuta asemenea copii să învețe limba la această vârstă. Este, într-adevăr, demonstrat că învățarea primelor elemente ale unei limbi la aceste vârste asigură o șansă mult mai mare de învățare a limbii ulterior și de inserare socială [24].

Lucrarea prezentă se situează într-un context mai larg, în cadrul cercetărilor realizate de diverse colective care caută soluții pentru a face vocea sintetică purtătoare de informație emoțională. Astfel, în [31] se descrie o metodă de sinteză a “vocii emoționale”, capabilă să transmită trei emoții (supărare-furie, bucurie, tristețe) folosind elemente de prozodie și segmente de tip vocală-consoană-vocală (specifice limbii japoneze). În [32], starea (“mood”) și personalitatea sunt văzute ca elemente esențiale aparând în subsidiar în voce și necesar a fi introduse și în vocea sintetizată. Alți autori [33] vorbesc de “nivelul de plăcere al audiției” (pleasantness) – dincolo de inteligibilitate – și văd naturațea vocii sintetizate prin această prismă, a utilizării la nivel semnificativ, a prozodiei (“...we need to know more about how prosody could be utilized in human-computer interaction. We believe that we could borrow a lot from professional human speakers. Furthermore, speech applications should be built in a way that makes it possible to use prosodic features efficiently.”).

Comment [R1]: “

Credem că, în viitor, o metodă comodă de a genera automat prozodia, pentru o voce artificială dată ⁹¹ pentru o anumită stare, ar putea fi constituită de o procedură inversă celei descrise în [34].

Incheiem cu un citat din [35]: "... in spite of the long history of speech synthesis, no one speech synthesis system available today is able to produce speech that could be characterized as natural or completely pleasant. In order to improve the speech quality of current text-to-speech (TTS) systems in terms of naturalness, three areas must be addressed⁹⁷: 1) improved linguistic analyses, 2) improved prosody modeling, and 3) improved speech synthesis models."

Mulțumiri. Această lucrare a fost realizată cu sprijinul material al Academiei Române – Institutul de Informatică Teoretică Iași – precum și cu sprijinul material parțial al Societății "Tehnici și Tehnologii" s.r.l. Iași. Autorul mulțumește colegilor Dragoș Burileanu, Bogdan Branzilă și Oana Geman pentru sugestiile și corecțiile la o formă preliminară a lucrării.

Referințe bibliografice

- [1] Teodorescu H.N., Chelaru M., Sofron E., Adascalitei A.: Adaptive speech synthesis. In vol. *Digitale Sprach-verarbeitung - Prinzipien und Anwendungen*. VDE Verlag, Berlin (W), pp. 183-188, 1988
- [2] Teodorescu H.N.: Interrelationship, Communication, Semiotics, and Artificial Consciousness. In: Kitamura, T. (Ed.): *What Should be Computed to Understand and Model Brain Functions?* FLSI Book Series, vol. 3, World Scientific, 2000
- [3] Teodorescu H.N.: Computer semiotics: understanding meanings and parallel languages (Refereed invited paper), Proc. Int. Conf. IIZUKA'98, Japan, 1998
- [4] Teodorescu H.N.: Making speech synthesizers noise-adaptable. *Electronic Engineering* (UK), July 1987, p. 23
- [5] Rodriguez, W., Teodorescu H.N., Grigoras Fl., Kandel A., Bunke H.: A Fuzzy information space approach to speech signal nonlinear analysis. *J. of Intelligent Systems* (Wiley), Dec. 1999
- [6] Grigoras Fl., Teodorescu H.N., Apopei V.: Nonlinear Analysis and Synthesis of Speech. *Studies in Informatics and Control*, vol. 7, no. 1, March 1998, pp. 57-72
- [7] Teodorescu H.N., Grigoras Fl., Apopei V.: Nonlinear processes in speech production. *Int. J. Chaos Theory and Applications*, vol. 2, no. 2 (1997), pp. 35-52

⁹⁷ Aici, autorul citat face referire la L. R. Rabiner, "Applications of Voice Processing to Telecommunications," Proc. IEEE, vol. 82, pp. 199-228, February 1994.

-
- [8] Teodorescu H.N., Grigoraş Fl.: Nonlinear Techniques in Speech Signal Analysis. Proc. International Conference on Intelligent Technologies in Human-Related Sciences, ITHURS'96. July 5-7, Leon, Spain. Vol. 2, pp. 293-298, 1996
- [9] Grigoraş Fl., Teodorescu H.N., Apopei V.: Analysis of nonlinear and nonstationary processes in speech production, IEEE 1997 Workshop on Applications of Processing to Audio and Acoustics. Mohonk Mountain House New Paltz, New York, October 19-22, 1997 (IEEE Catalog # 97TH8278)
- [10] Burlui V., Teodorescu H.N., Moraraşu C.S.: La fonction phonatoire chez l'edente total. Analyse en frequence. *Les Cahiers de Prothese* (France), No. 88, Decembre 1994, pp. 63-68 1994
- [11] Teodorescu H.N. et al.: Fuzzy models in speech analysis and medical application, in Book of Summaries Int. Conf Modelling and Simulation, Istanbul, Turkey, July 1988, vol. 1, p. 162 (Summary)
- [12] Teodorescu H.N., L. Buchholtzer, Chelaru M., Teodorescu L.: A laryngeal prosthesis based on perilaryngean reflexes, Proc. 9th Int. EMBS Conf. IEEE, Boston. Vol. 4, IEEE, pp. 2114-2115, 1987
- [13] Anonymous Automotive Industry OEM/Supplier: Talking to computers vs. talking to humans 7/12/2000. <http://www-nrd.nhtsa.dot.gov/departments/nrd-13/driver-distraction/Topics013040293.htm#A293>
- [14] Anne-Marie Derouault, The Future of Speech Recognition. Evolving speech recognition technology is driving transparent computing, making it easier for people to interact with computers. <http://www.advisor.com/Articles.nsf/ID/OA000107.DERO01>
- [15] House D., Bell L., Gustafson K. & Johansson L. Child-directed speech synthesis: evaluation of prosodic variation for an educational computer program. Proc of Eurospeech'99, pp. 1843-1846, 1999
- [16] Heldner M., Strangert E. & Deschamps T.: Focus detection using overall intensity and high frequency emphasis. In: Andersson R, Abelin L, Allwood J & Lindblad P, eds. Proc of Fonetik 99; pp. 73-76, 1999.
- [17] Heldner M., Strangert E. & Deschamps T.: A focus detector using overall intensity and high frequency emphasis. Proc of ICPhS-99, pp. 1491-1494, 1999.
- [18] Heldner M.: On the non-linear lengthening of focally accented Swedish words. In: W. van Dommelen & T Fretheim, eds. Nordic Prosody: Proc of the VIIIth Conference, Trondheim 2000 . Frankfurt am Main: Peter Lang. 2001
- [19] Karlsson I., Banziger T., Dankovicová J., Johnstone T., Lindberg J., Melin H., Nolan F. & Scherer K.: Within-speaker variability due to speaking manners. Mannell RH & Robert-Ribes J, eds. Proc of ICSLP98, 2379-2382. 1998

-
- [20] Karlsson I.: Within-speaker variability in the VeriVox database. In: Andersson R, Abelin L, Allwood J & Lindblad P, eds. Proc. of Fonetik 99, pp. 93-96, 1999.
- [21] Karlsson I, Banziger T, Dankovicova J, Johnstone T, Lindberg J, Melin H, Nolan F, Scherer K (1998), Within speaker variation due to induced stress, Proc Fonetik-98, 150-153. www.ling.su.se/fon/publications/fonetik98/
- [22] Gustafson-Capkova S & Megyesi B.: A Comparative Study of Pauses in Dialogues and Read Speech. Proc of Eurospeech 2001, pp. 931-935, 2001
- [23] Beskow J.: A tool for teaching and development of parametric speech synthesis. In: Branderud P & Traunmüller H (eds). Proc of Fonetik -98, pp. 162-165. 1-98, 1998
- [24] Rachel I. Mayberry, Elizabeth Lock, Hena Kazmi: Linguistic ability and early language exposure. *NATURE*, Vol. 417, 2 May 2002, p. 38, 2002
- [25] Microsoft Co.: Platform SDK: Agent. Characters. http://msdn.microsoft.com/library/default.asp?url=/library/en-us/msagent/deschar_8nn6.asp
- [26] Mauricio Lumbreras, Gustavo Rossi: Metaphor for the Visually Impaired: Browsing Information in a 3D Auditory Environment. CHI'95 Proc., www.acm.org/sigchi/chi95/proceedings/shortppr/ml_bdy.htm
- [27] Christophe d'Alessandro & Jean-Sylvain Liénard: 5.2 Synthetic Speech Generation. In: Survey of the State of the Art in Human Language Technology. <http://cslu.cse.ogi.edu/HLTsurvey/ch5node4.html#SECTION52>
- [28] Teodorescu H.N.: Chaos in fuzzy systems and signals. Vol. Proceedings of the 2nd Int. Conf. on Fuzzy Logic and Neural Networks. Vol. 1., pp. 21-50 (Jono Printing Co., 1992, Iizuka, Japan)
- [29] Teodorescu H.N., Kandel A., Jain L. C. (Eds.), Fuzzy and Neuro-Fuzzy Systems in Medicine (International Series on Computational Intelligence). CRC Press, Boca Raton, USA, 1998.
- [30] Teodorescu H.N., Mlynek D., Kandel A. (Eds.): Intelligent Systems and Interfaces (The Kluwer International Series In Intelligent Systems). Kluwer Publ., Boston, 2000.
- [31] Yasuhisa Niimi, Masanori Kasamatu, Takuya Nishimoto and Masahiro Araki: Synthesis of Emotional Speech Using Prosodically Balanced VCV Segments. <http://www.ssw4.org/papers/133.pdf>.
- [32] Nick Campbell: WHERE IS THE INFORMATION IN SPEECH? (and to what extent can it be modelled in synthesis?) www.slt.atr.co.jp/cocosda/jenolan/Proc/r82/r82.pdf.
- [33] Hakulinen J., Turunen, M.: Prosodic Features for Speech User Interfaces. www.cs.uta.fi/hci/spi/reports/Prosodic_Features_for_Speech_User_Interfaces.pdf.
- [34] Ansgar Rinscheid: Voice Conversion Based On Topological Feature Maps and Time-Variant Filtering. www.asel.udel.edu/icslp/cdrom/vol3/235/a235.pdf.

- [35] Syrdal A., Stylianou Y., Garrison L., Conkie A. Schroeter J.: Td-Psola Vs. Harmonic Plus Noise model in Diphone Based Speech Synthesis. www.research.att.com/projects/tts/papers/1998_ICASSP/paperSYN.ps.

Anexa 1: Sisteme nuanțate de tip Sugeno, de ordin 1. Funcții de apartenență

Reamintim ca o mulțime (clasică) $A \subset X$, unde X notează universul de discurs, este definită de o funcție caracteristică, de forma:

$$\chi_A(\cdot): X \rightarrow \{0,1\}$$

$$\chi_A(x) = \begin{cases} 1 & \text{daca } x \in A \\ 0 & \text{daca } x \notin A \end{cases}$$

Prin generalizarea conceptelor de mulțime și de funcție caracteristică, se definesc mulțimile nuanțate (fuzzy) și funcțiile de apartenență corespunzătoare astfel: o mulțime nuanțată, notată \tilde{A} , peste universul de discurs X , este caracterizată unic de o funcție de apartenență:

$$\mu_A(\cdot): X \rightarrow [0,1]$$

În particular, funcția de apartenență poate fi de forma:

$$\mu_A(x) = \begin{cases} 1 & \text{pentru } x = a \in X \\ 0 & \text{pentru } x \neq a \end{cases}$$

caz în care se numește *singleton*.

Un sistem de tip Sugeno, de ordin 0, este descris de reguli de forma:

*DACA intrarea (premise) # 1 SI premise # 2 SI ... Si premise # n ATUNCI
concluzia*

unde premisele sunt de forma: x_i este \tilde{A}_{ij} , iar \tilde{A}_{ij} sunt valori nuanțate (fuzzy), de exemplu \tilde{A}_{i1} = “mare”, \tilde{A}_{i2} = “mediu”, atributelor lingvistice “mare”, “mic” etc. fiindu-le atașate câte o funcție de apartenență. Specific sistemelor Sugeno este faptul ca în concluzie apar valori numerice și nu valori nuanțate, concluzia fiind deci de forma “ $y = 0,3$ ” (singleton).

Definițiile funcțiilor de apartenență pentru intensitatea sonoră din Figura 1.a sunt:

$$\mu_{Putere=mica}(p) = \begin{cases} 1 & \text{pentru } p \leq 40 \text{ dB} \\ 1 - \frac{p-40}{15} & \text{pentru } 40 < p < 55 \text{ dB} \\ 0 & \text{pentru } p \geq 55 \text{ dB} \end{cases}$$

$$\mu_{Putere=medie}(p) = \begin{cases} 0 & \text{pentru } p \leq 40 \text{ dB} \\ \frac{p-40}{15} & \text{pentru } 40 < p < 55 \text{ dB} \\ 1 - \frac{p-55}{15} & \text{pentru } 55 \leq p < 70 \text{ dB} \\ 0 & \text{pentru } p \geq 70 \text{ dB} \end{cases}$$

$$\mu_{Putere=mare}(p) = \begin{cases} 0 & \text{pentru } p \leq 55 \text{ dB} \\ \frac{p-55}{15} & \text{pentru } 55 < p < 70 \text{ dB} \\ 1 & \text{pentru } p \geq 70 \text{ dB} \end{cases}$$

Definițiile funcțiilor de apartenență pentru raportul HL (Figura 1b) sunt:

$$\mu_{HL=mica}(q) = \begin{cases} 1 & \text{pentru } q \leq 0.5 \\ 1 - \frac{q-0.5}{.5} & \text{pentru } 0.5 < q < 1. \\ 0 & \text{pentru } p \geq 1. \end{cases}$$

$$\mu_{HL=medie}(q) = \begin{cases} 0 & \text{pentru } q \leq 0.5 \\ \frac{q-0.5}{0.5} & \text{pentru } 0.5 < q < 1.0 \\ 1 - \frac{q-1.0}{0.5} & \text{pentru } 1.0 \leq q < 1.5 \\ 0 & \text{pentru } q \geq 1.5 \end{cases}$$

$$\mu_{HL=mare}(q) = \begin{cases} 0 & \text{pentru } q \leq 1.0 \\ \frac{q-1.0}{0.5} & \text{pentru } 1.0 < q < 1.5 \\ 1 & \text{pentru } q \geq 1.5 \end{cases}$$

Pentru detalii asupra manipulării funcțiilor de apartenență și a regulilor în sistemele nunațate, a se vedea orice manual în domeniul sistemelor fuzzy, sau volume precum [29, 30] în care se pot găsi și aplicații specifice legate de înțelegerea vorbirii, sau alte aplicații medicale.

Anexa 2: Procesul haotic

Procesul reprezentat de ecuațiile (7) are o dinamică haotică doar pentru anumite subintervale relativ înguste din \mathbf{R}^6 . În restul spațiului, comportamentul este asimptotic instabil (peste tot pentru valori ale coeficienților lui r^3 mai mari ca 1, în modul, dacă și coeficientul lui u este mai mare ca 1 în modul); comportamentul este stabil sau periodic pentru alte zone, relativ reduse din \mathbf{R}^6 .

Diagrama de bifurcație a procesului, așa cum apare în Figura A1, este obținută pentru: valorile coeficienților $[Q]=\{.1, -.17, -.18, .1\}$; $\text{coeff}_4 = 1.1$; $\text{coeff}_5 = -.15$; condiție inițială $r[0] = 0.3$; număr total de puncte în diagrama de bufurcație: 500 (punctele de la 500 la 1000); regimul tranzitoriu eliminat: primele 500 puncte; precizia tuturor coeficienților și variabilelor: double.



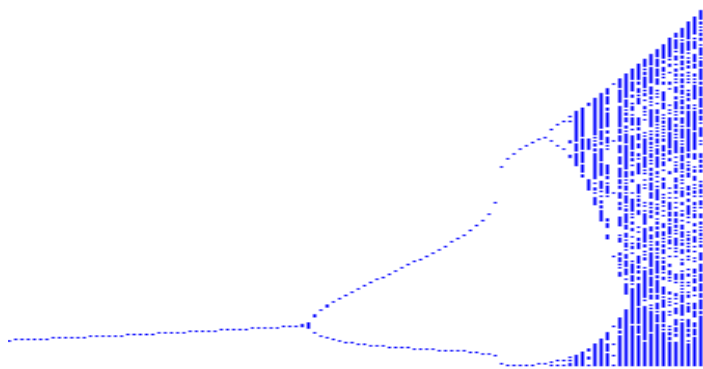


Figura A2-1. Diagrama de bifurcație a procesului

Legile folosite (conform codului, scris în limbajul C) sunt:

$$u[n] = (\text{coeff_4}) * r[n] + \text{coeff_5} - 0.005 * (\text{float})k;$$

$$x = u[n]; \quad r[n+1] = \text{poly}(x, Q, \text{coeff});$$

(Q este numărul de valori în vectorul coeficienților, Q=4)

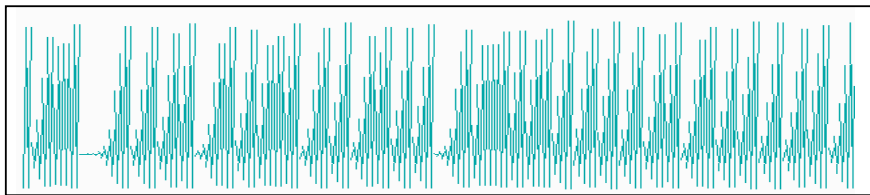


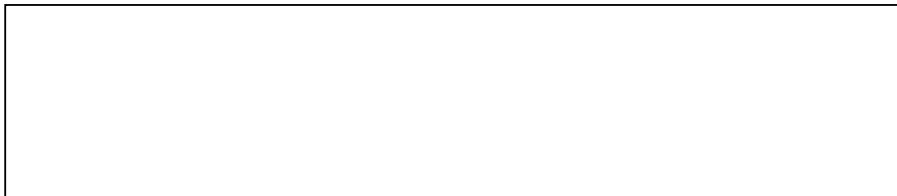
Figura A2-2

Semnalul în domeniul amplitudine-timp din Figura A2 a fost obținut pentru ecuațiile (cod C):

$$u[n] = \text{coeff_4} * r[n] + \text{coeff_5} - 0.05 * 21.;$$

$$x = u[n]; \quad r[n+1] = \text{poly}(x, Q, \text{coeff});$$

Semnalul obținut pentru valoarea $u[n] = \text{coeff_4} * r[n] + \text{coeff_5} - 0.05 * 21.$ (restul programului fiind identic ca pentru cazul anterior) este ilustrat în Figura A3.



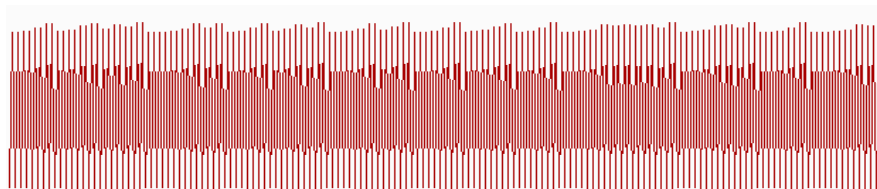


Figura A2-3

iar semnalul obținut cu $u[n] = \text{coeff_4} * r[n] + \text{coeff_5} - 0.05 * 20.7$, precum și la o scară dublă de timp, este ilustrat în Figura A4:

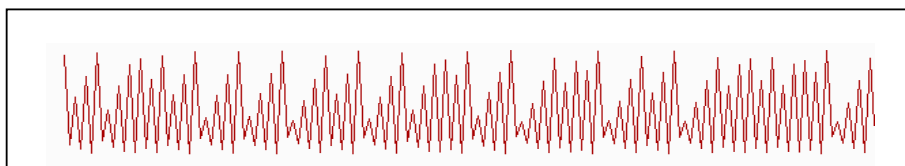


Figura A2-4

Regiunile spațiului parametrilor în care sistemul este stabil, după cum s-a spus deja, sunt relativ înguste. Pentru parametrii coeff_1 - coeff_4 fixați și coeficientul coeff_5 variabil între -25.15 și $+4.85$ (600 de pași, cu pas $0,05$), doar zona îngustă din Figura A2-5 este stabilă, oscilantă sau haotică, în rest sistemul fiind asimptotic instabil. Pentru ușurința urmăririi scării, linia din partea de jos a figurii reprezintă intervalul menționat, $[-25.15, +4.85]$, în care s-a testat sistemul.



Figura A2-5

În figură, se poate remarca diagrama de bifurcație a sistemului, cu zonele de stabilitate, oscilație și haos. Pentru restul intervalului, prin program, calculele sunt abandonate, deoarece valorile de ieșire ale sistemului depășesc, în valoare absolută, 10000 .

Dicționarele multimedia ale limbii române. Secvențe de implementări și experimentări

Dumitru TODOROI, Diana MICUSA, Zinaida TODOROI, Ion LINGA, Ion COVALENCO, Nicolae OBJELEANU, *tefan SPĂȚARU, Stela LUNGU, Virginia ȚURCANU, Elana COZLOV, Nadejda AMBROZII, Victor SLOBODEANU, Igor CO^aERU, Cătălina SURUCEANU
Academia de studii economice din Moldova, Str. Banulescu-Bodoni, 59-61/503»B», Chi^oinău MD 2005, Republica Moldova,
E-mail : todoroi@ase.md

Lucrarea actuală în cadrul punerii Marelui Dicționar al Limbii Române (MDLR) pe calculator a fost metodologic influențată de ideile subliniate în [1-3] și este o continuare a cercetărilor [4-7,10-11], efectuate în cadrul procesării limbajului natural. Au fost elaborate un șir de proiecte [8-9,12] de informatizare a Limbii Române. Experimentările cu elaborarea sistemelor computerizate de nivelul unu, care susțin diferite sub-dicționare ale MDLR pe așa axe ale lui ca: TEXT, AUDIO, IMAGINI și VIDEO, au început recent în Academia de studii economice din Moldova (ASEM) în colaborare cu ONG-ul ECO-INFO-MOLD. Unele rezultatele de cercetări și experimentări în cadrul platformei, alcătuite din aceste 4 subsisteme, sunt expuse în lucrarea de față. Sunt prezentate diferite scenarii [19] și metodologii de utilizare a sub-dicționarelor informatizate a limbii române. Clarificarea mijloacelor Hardware-ului și Software-ului modern, care pot suporta MDLR informatizat [18], constituie o problemă importantă la etapa creării Societății Informaționale – Societate a Cunoașterii [20].

I. Componenta TEXT a dicționarului economic MULTIMEDIA al limbii Române [23].

Scopul acestui compartiment computerizat al MDLR constă în crearea subsistemului TEXT de nivelul unu a unei părți introductive a dicționarului economic al limbii române și experimentarea cu acest sistem. Acest dicționar economic constă din 35.000 – 40.000 cuvinte. Cuvintele, care se conțin în Dicționarul Enciclopedic Ilustrat (DEI) [21], vor fi definite 100/100. Experimentarea cu subsistemul TEXT al MDLR computerizat este efectuată la moment cu circe 200 articole din DEI.

Baza de date TEXT (BDT) a dicționarului economic constă dintr-o culegere de texte-articole, alcătuită din cuvinte, fraze, paragrafe, capitole etc. ale DEI. Documentele în

BDT includ nu numai informații textuale (definiții de cuvinte), dar pot conține informații și de alt tip, de exemplu, prin extindere, imagini. Prin urmare BDT în sistemul computerizat al MDLR conține nu numai materialul textual, dar și ilustrativ: diagrame, grafice, fotografii etc.

Prin crearea subsistemului TEXT de nivel unu utilizatorul obține un mijloc important, prin intermediul căruia informația poate fi introdusă și utilizată în mod complementar pe cale electronică.

1.1. Capacitățile necesare ale unui sistem de gestiune a bazei de date MULTIMEDIA

MDLR este o bază de date MULTIMEDIA. Sistemul de gestiune al MDLR este un sistem de gestiune a bazei de date MULTIMEDIA (SGBDMM) și constituie un mecanism, care operează cu diferite tipuri de date, reprezentate într-o diversitate de formate pe un set larg de mijloace și surse. Pentru a funcționa efectiv e necesar ca SGBDMM să poseze următoarele capacități:

- (a) Capacitatea de a interoga datele, uniform reprezentate în diferite formate;
- (b) Capacitatea de a interoga datele, reprezentate în diferite medii;
- (c) Capacitatea de a transmite subiectele mediei din dispozitivele de stocare locale într-un mod efectiv;
- (d) Capacitatea de a primi răspunsul la o interogare și de a desfășura prezentarea acestui răspuns pe baza mediei audiovizuale;
- (e) Capacitatea de a furniza această prezentare pe acea cale, care ar satisface calitățile diferitor cerințe ale serviciului.

1.2. Structura bazei de date TEXT (BDT) a dicționarului economic MULTIMEDIA

Dicționarul economic, care este pe cale de a fi pus pe calculator, este o BDT cu posibilitatea de a fi extinsă cu diferite componente ale MULTIMEDIEI. Subsistemul TEXT a dicționarului economic MULTIMEDIA este un subsistem al SGBDMM, care aprovizionează această posibilitate împreună cu utilizarea complementară a BDT.

Structura BDT e compusă din:

- (1) Indice cu caracteristica "număr";
- (2) Termen principal (cuvânt, articol) cu caracteristica "text";
- (3) Variantă(e), derivate, abreviere (concretizare) cu caracteristica "text";
- (4) Categorie gramaticală cu caracteristica "text";
- (5) Domeniu cu caracteristica "text";
- (6) Definiții pentru termenul principal (și concretizări) cu caracteristica "text";

-
- (7) Sinonim(e) cu caracteristica "text";
 - (8) Antonim(e) cu caracteristica "text" și aștele. De asemenea BDT are posibilitatea de a fi extinsă cu așa subdiviziuni MULTIMEDIA ca:
 - (9) Audio cu caracteristica „OLE”;
 - (10) Imagini cu caracteristica „OLE”;
 - (11) Video cu caracteristica „OLE” și altele.

1.3. Scenarii de utilizări și interogări a subsistemului TEXT al MDLR unformatizat

Interogarea este o formă, care ajută utilizatorul să prezinte o informație anumită într-o structură anumită, definită de utilizător. Spre exemplu, utilizătorul dorește să obțină informații din arhive, articole, sau alte documente, care conțin informația despre Uniunea Europeană. Interogarea poate avea următoarea formă: "Găsește toate dosarele, legate de investițiile străine, făcute de UE în domeniul educației". Un simplu cuvânt cheie a acestui dosar nu va permite găsirea răspunsului corect, chiar și dacă indicile acestui document deja există. Totuși, sistemul ne va prezenta careva cuvinte, legate de această interogare, dar ele pot să nu fie direct relatate la tema dorită. De aceea textul trebuie să fie indexat nu numai pe cuvintele cheie, dar și pe conținutul semantic și/sau pragmatic al cuvintelor (în cazul BDT, de exemplu, concretizarea).

Soluționarea problemei utilizătorului, care dorește să afle definiția cuvântului "Academie", de exemplu, cere introducerea polisemiei în BDT, care conține așa concepte ca precizia și rechemarea. Întrebarea, propusă de către utilizător în acest context, este: "Cum să aflu din baza de cunoștințe a MDLR sensul cuvântului "Academie – ca instituție de învățământ economic". Pentru aceasta BDT va fi completată cu o nouă coloană "concretizare", care va preciza și va face posibilă afișarea pe monitor anume a acelei definiții a cuvântului, de care utilizatorul este cointerestat (de exemplu: Academia de studii economice).

Un fragment de structură schematică a BDT.

I Indece	I Cuvânt	I Concretizare	I I	I Definiție	I I	I Traducere cuvânt
I 03342	I Academia	I de studii economice	I	I Nume dat °colii de ...I	I	I
I 14269	I Banii	I EURO	I	I Denumire a princip...I	I	I
I 14271	I Banca	I de economii	I	I	I	I

SGBDMM, ca o extindere a SGBD Ms ACCESS-2000, în baza căruia este creată componenta TEXT a MDLR informatizat, gestionează BDT, utilizând limbajul SQL. În exemplele următoare utilizatorul este cointerestat în sfera finanțelor. Accesul BDT a MDLR este efectuat prin intermediul următoarelor interogări din SQL (care, în general, constituie așa comenzi ca SUMMARING, JOIN, PROJECTION, DIVISION, SELECT și altele):

Ex.1. SELECT Banii (termen principal, nume de interes)
 FROM Ambrozii-Godzina (nume de fișier)
 WHERE Concretizare = EURO (concretizare pentru termenul principal)

Ex.2. SELECT Academia
 FROM Ambrozii-Godzina
 WHERE Concretizare = de studii economice

Ex.3. SELECT Banca
 FROM Ambrozii-Godzina
 WHERE Concretizare = de economii.

II. AUDIO-dicționarul explicativ economic al limbii Române [24]

Dicționarul explicativ economic MULTIMEDIA al limbii române, ca o parte componentă al DEI, include circa 35000-40000 de cuvinte și este divizat în compartimentele: Text, Audio, Video și Imagini. Aceste componente MULTIMEDIA ale MDLR informatizat satisfac cerințele de bază către un dicționar informatizat: prezintă formele exacte ale cuvintelor, accentul, etimologia, definiția-tezt, definiția-sunet (audio), definiția-video (film), definiția-imagine (grafic, schema, poza etc.) și corespunde cerințelor unor categorii foarte largi de utilizatori nu numai elevi și studenți, dar și funcționari și profesioniști, contribuind la ridicarea nivelului de cultura.

Compartimentul AUDIO al MDLR informatizat furnizează informații necesare ale articolului respectiv (cuvântul, definiția lui) în forma AUDIO. Subsistemul AUDIO de nivelul unu al SGBDMM oferă posibilitatea de AUDIO-utilizare a dicționarului. Acest AUDIO-dicționar va contribui din plin la ridicarea pe o treaptă superioară a societății noastre în exprimarea economică corectă orală și scrisă. Conținutul datei AUDIO poate fi caracterizată prin două metode: (a) folosind metadata prin explicarea conținutului unui fișier AUDIO sau (b) prin extragerea tipului potrivit de date AUDIO, folosind procesorul tehnic.

2.1. Componenta AUDIO a metadatelor

Cu un fișier AUDIO se procedează la fel ca și în cazul unei date VIDEO: acestui fișier i se asociază un set (grup) de segmente, toate referindu-se la o perioadă de timp. Fiecărui segment i se atribuie un set de activități, care au decurs în acea perioadă de timp, subliniate prin aceste segmente. În general, metadata utilizează reprezentarea AUDIO, care este sesizată ca un set de obiecte marcate în timp.

Utilizarea componentei **AUDIO** a metadatei din MDLR informatizat este mai mult recomandată când este o modalitate de creare a acestei metadata, modificarea ei, în deosebi, la interogarea **AUDIO** –dicționarului de utilizatorii, care necesită această formă de comunicare om-mășină.

Crearea componentei **AUDIO** a metadatei este un lucru mai complex decât alte forme de dicționare informatizate, deoarece identitatea indivizilor, ce vorbesc, nu poate fi ușor cunoscută; chiar și conținutul discursului poate fi neclar.

Conceptul despre conținut este descris în termeni de metadata a procesului. Ca rezultat data **AUDIO** este considerată ca un semnal $\Delta(x)$ în timpul x . Trăsăturile utilizatorice ale acestui semnal $\Delta(x)$ sunt: (a) extragerea, (b) indicarea și (c) depozitarea.

O undă constă dintr-un set de vîrfuri (creste) și adîncituri (vâi). Perioada vibrației T este definită ca timpul, pentru care o parte a undei să revenă la poziția inițială. Alte caracteristici utilizate de componenta AUDIO în crearea metadatei sunt: (1) frecvența, (2) viteza și (3) amplituda.

Baza de date **AUDIO (BDA)** poate fi interacționată și gestionată, utilizând sunetul auditiv prin intermediul procesiunii de segmentare, memorizare și extragere a informației.

2.1.1. Segmentarea

Segmentarea e o procedură de separare a sunetului auditiv în câteva ferestruici egale. Această procedură poate fi utilizată conform următoarelor două metode:

- (a) Utilizatorul specifică dimensiunile ferestrei, presupunând că proprietățile unde și a ferestrei se vor obține prin medie;
- (b) Utilizatorul segmentează sunetul în același mod ca și imaginiile, folosind predicatul de omogenitate H.

2.1.2. Extragerea

La extragere cel mai des utilizate sunt facilitățile de indicare a intensității, zgomotului, înălțimii și strălucirii.

2.2. Unele sisteme de utilizare a BDA

Din punct de vedere a MULTIMEDIEI, AUDIO - baza de date (BDA) poate fi interpretată ca o sursă auditivă, ca un fișier cu o fereastră auditivă și cu trasăturile respective, asociate acestei ferestre.

Scenariile de utilizare a BDA cuprind toate formele MULTIMEDIEI, care pot fi utilizate în diferite domenii. În sistemele comerciale, de exemplu, **Bazele de date Informix** includ bazele de date a sistemului managerial, care permit utilizatorului să acceseze baza de date, bazându-se pe nesiguranța conținutului.

Baza de date DB2, un alt exemplu, utilizată cu calculatorul de tip IBM, necesită cuplarea cu un sistem adăugător, care permite lăsarea mesajelor vocale pe robot. DB2 poate importa și menține clipurile, care pot fi cautate printr-un nume sau descriere.

Putem reaudia mesajele, lăsate pe robot, prin intermediul **Internetului**. Un exemplu în plus constituie o utilizare a unui cuvânt din AUDIO-dicționarul economic al limbii române prin intermediul AUDIO-VIDEO-robotului, care este un sistem autorizat și care acționează pe baza unui program de lucru stabilit sau care reacționează la anumite influențe exterioare.

Un exemplu de interogare a componentei AUDIO a subdicționarului economic al MDLR prin intermediul limbajului SQL și al subsistemului AUDIO de nivel unu al SGBDMM poate avea forma :

```
SELECT      Robot
FROM        Țurcan-Mutruc
WHERE       Attribute IS Definiție AND Attribute IS Audio
```

Ca rezultat al acestei interogări utilizatorul prin intermediul răspunsului prietenos obține pe ecran definiția TEXT a cuvântului Robot și în paralel acest subsistem AUDIO al SGBDMM difuzează această definiție cu voce feminină sau masculină (la dorința utilizatorului).

III. Subsistemul IMAGINI de nivel unu al dicționarului economic informatizat al limbii române [25]

Scopul acestui capitol constă în descrierea posibilităților de introducere a imaginilor în baza de date a MDLR informatizat și de utilizare a acestora în viața cotidiană. Daza de date IMAGINI (BDI) a subdicționarului economic al MDLR informatizat constituie o componentă, care oferă posibilitatea de extindere a procesului de înțelegere a sensului cuvântului dat. Din cele aproximativ 35000-40000 de articole ale dicționarului economic din MDLR doar numai 50-60%, după părerea noastră, pot fi prezentate în forma de imagini.

Experiența, obținută pe baza câtorva zeci de articole din DEI în cadrul evaluării subsistemului IMAGINI al SGBDMM, ne confirmă întru totul conținutul zicalei: «Mai bine odată să vezi decât de o sută ori să auzi» și a zicalei «Un tablou este egal cu o mie de cuvinte». Aceste facilități utilizatorice din evoluția creatorică și utilizatorică a MDLR informatizat le confirmă și lucrările din [22] chiar și prin intermediul următorului Tabel 3.1, prezentat în original.

Table 13.1.

Data rates and storage requirements per hour, day, and lifetime for a person to record all the text they've read, all the speech they've heard, and all the video they've seen

Data type	data rate (bytes per second)	storage needed per hour and day	storage needed in a lifetime
Read text, few pictures	50	200 KB; 2-10 MB	60-300 GB
speech text @ 120 wpm	12	43 K; 0.5 MB	15GB
speech (compressed)	1,000	3.6 MB; 40 MB	1.2TB
video (compressed)	500,000	2 GB; 20 GB	1 PB

3.1. Baza de date IMAGINI (BDI)

Imaginea poate povesti mai mult despre un obiect decât câteva pagini (Vezi Tabelul 3.1) de descrieri textuale. Pentru un chirurg este cu mult mai ușor să găsească un pacient potențial prin investigarea diferitor imagini. Imaginile pot fi combinate cu corpusuri, text-definiții, sunet-definiții, traduceri etc.

În afara de datele IMAGINI ale dicționarului economic MULTIMEDIA în MDLR informatizat vor fi prezente a'a tipuri de date ca video, audio, document, manuscrise și altele. Datele VIDEO sunt des folosite în domeniul învățămîntului. Datele AUDIO sunt importante în domeniul criminalisticii, de exemplu, în identificarea vocilor celor suspectați. Datele documentare diferă de datele TEXT prin aceea, că pot să conțină nu numai informații textuale, dar și imagini încadrate. Datele manuscrise se presupune că în viitorul apropiat vor prevala înregistrările electronice.

Sunt cunoscute diferite formate electronice, care dau posibilitatea de a vizualiza imaginea (fișierele de tip GIF, TIFF, PCX, de exemplu). Subsistemul IMAGINI a SGBDMM are anumite trăsături specifice necesității de utilizare a imaginilor ca o componentă vitală a MDLR informatizat.

3.2. Subsistemul IMAGINI

Baza de date IMAGINI diferă de bazele de date TEXT și AUDIO prin complicitatea imaginilor, necesitatea de a diviza, combina și utiliza diferite părți componente ale imaginii, care deseori la interogare se complică și prin utilizarea incorectă și analiza neprecizată a tehnicilor de manipulare a imaginilor. Aceasta se complică și prin aceea, că diferite organizații adună date fotografice, hărți, scheme și alte imagini de tip universal sau specializat (cum ar fi, de exemplu, NASA). Interogările datelor de tip IMAGINI sunt efectuate în baza datelor de tip TEXT, cautate în baza de date de tip IMAGINI și vizualizate în formă de text și imagini. În final imaginile pot fi transferate în baza de date specializate, cum ar fi, de exemplu, încadrarea lor în baza de date MULTIMEDIA comerciale. În subsistemul IMAGINI al SGBDMM sunt prevăzute un set larg de proceduri cu imaginile.

3.2.1. Plasa imaginii

Conținutul imaginii constă din toate obiectele acestei imagini și caracteristicile lor, care reprezintă interes din punct de vedere a programului aplicativ. Imaginea poate avea o mulțime de proprietăți, așa ca descrierea formei, prezentarea vectorului subdiviziunilor, prezentarea vectorului ordinii de descompunere și compunere a imaginii și altele. Fiecare imagine "I" are o pereche asociată schimbătoare de numere pozitive (m,n), care se numește **plasa imaginii**. Ea este compusă din $m \cdot n$ celule de măsuri egale.

3.2.2. Transformări de imagini

Imaginea se împarte în părți omogene, care se numesc segmente. Schemele de compresare a imaginii sunt invertibile, deoarece unele scheme de compresare pot conduce la pierderea informației sau la pierderea perfecțiunii. Există două abordări a problemei căutării similitudinii imaginilor: Abordarea metrică și abordarea de transformare.

Abordarea de transformare este mai generală decât abordarea metrică. Această abordare utilizează așa operațiuni ca: transformarea, transferarea, rotația, scalarea, simetrizarea și a.

3.3. Utilizarea imaginii

În prezent multe instituții de învățământ oferă programe de studii individuale. Unele persoane studiază cursuri speciale de sinestător pentru dezvoltarea intelectului uman. Astfel de cursuri pot fi reprezentate sub formă de imagini speciale.

Imaginile pot fi utilizate în industria filmelor. Specialiștii au posibilitatea de a vizualiza imaginile, alese de ei, lucrând la calculator.

Imaginile sunt importante și în industria turismului. Pentru informații despre imaginile necesare la fel se poate apela la sistemul de tip IMAGINE al SGBDMM.

Interogările de imagini în dicționarul economic al MDLR informatizat pot fi efectuate la fel ca și în subsistemele de același nivel unu de tip TEXT și AUDIO prin intermediul limbajului SQL al SGBD. Rezultatul în forma textuală a articolului și imaginea într-o formă complementară este reprezentată utilizatorului în formă de Soft-copy sau Hard-copy.

Obținerea imaginii cuvântului «bancă», de exemplu, în subsistemul IMAGINI al SGBDMM al MDLR informatizat se efectuează prin intermediul următoarelor acțiuni. Se deschide baza de date IMAGINI al dicționarului economic al limbii române (în care sunt acumulate la momentul experimentării cu SGBDMM al MDLR informatizat doar numai 25 de cuvinte cu imagini respective). Se alege cuvântul «banca». În înregistrarea respectivă a băncii în compartimentul imagini se găsește OLE al imaginii cuvântului ales. Se efectuează clic pe ea și vizualizăm pe ecran imaginea respectivă. Analogic se procedează și cu alte cuvinte din BDI.

IV. VIDEO-dicționarul economic al limbii române[26]

În ultimii ani a crescut imens necesitatea de a putea chestiona și procesa cantități mari de date, care nu sunt întotdeauna ușor de reprezentat prin intermediul simbolurilor. Exemple de astfel de date sunt: informația în formă de imagini, informația-video, datele-audio, informația textuală, notițe și altele. În continuare vor fi examinate unele momente de realizare a dicționarului economic informatizat cu VIDEO clipuri. A fost inițiată baza de date VIDEO (BDV) a dicționarului economic MULTIMEDIA – o subdiviziune a MDLR informatizat – prin crearea subsistemului VIDEO de nivel unu al SGBDMM. Se va demonstra viabilitatea acestui subsistem.

4.1. Problemele creării subsistemului VIDEO al SGBDMM.

Pentru a opera o bază de date MULTIMEDIA (BDMM) un SGBDMM trebuie să aibă următoarele abilități:

- (a) Capacitatea de a chestiona uniform datele reprezentate în diferite formate;
- (b) Capacitatea de a chestiona uniform datele reprezentate în diferite surse media;
- (c) Capacitatea de a aporta unitățile media dintr-o diviziune locală de depozitare, asigurând continuitatea acestui proces;
- (d) SGBDMM trebuie să primească răspunsul, generat de o chestionare și să poată genera o prezentare a acelui răspuns utilizând audiovizualul;
- (e) Capacitatea de a oferi prezentarea într-un mod, care ar satisface diferite cerințe ale utilizatorului.

Tehnologiile, legate de bazele de date, au dezvoltat în ultimii 40 de ani baza, pe care ar trebui să fie creată o BDMM. În prezent sunt create limbaje de chestionare, tehnicile de aranjare, algoritmi de aportare pentru o mulțime de baze de date de tip relațional, spațial, temporal și altele. Fiecare din aceste mijloace extind posibilitățile limbajelor și algoritmi precedenți pentru a face față noilor tipuri de date sau pentru a argumenta paradigmele respective.

În acest capitol se va analiza informația de tip VIDEO. Necesitatea de a accesa o bază de date VIDEO (BDV) poate apărea într-o mulțime largă de aplicații, și de obicei modelul de acces variază considerabil de la o aplicație la alta.

În procesul reprezentării conținutului unui film în BDV este necesar de răspuns la un set de întrebări de tipul:

- (A) Ce aspecte posibile ale filmului pot cointeresa utilizatorii BDV?
- (B) Cum pot fi aceste aspecte ale filmului eficient depozitate, încât să minimalizeze timpul necesar subsistemului VIDEO al SGBDMM pentru a răspunde interogărilor utilizatorilor?
- (C) Cum ar trebui să fie limbajul de interogare a datelor VIDEO și cum ar trebui schimbat modelul relațional pentru a corespunde informației VIDEO?
- (D) Poate fi oare automatizat procesul de extragere a informației în baza contextului?

Aceste probleme au fost abordate în procesul creării și experimentării cu BDV și subsistemul VIDEO de nivel unu al SGBDMM.

4.2. Definițiile datelor de tip VIDEO

De obicei un film este caracterizat prin personajele sale, atributele acestora și activitățile, în care sunt angajate aceste personaje. Principalele surse de interes într-un film includ: (a) oameni, (b) obiecte neînsuflețite, (c) ființe însuflețite și (d) activități.

De observat, că tema generală, care se repetă în fiecare cadru, constă în aceea, că există un grup de obiecte și activități asociate. Astfel vom încerca să definim o bază de date VIDEO printr-un set de definiții.

Definiție 1: O *proprietate* VIDEO este o pereche $(pname, Values)$, unde $pname$ este numele proprietății și $Values$ este o mulțime. O *instanță* a proprietății $(pname, Values)$, este o expresie de forma $pname=v$, unde $v \in Values$.

Definiție 2: O *schemă obiect* este o pereche (fd, fi) , unde:

fd este o mulțime de proprietăți cadru-dependente,

fi este o mulțime de proprietăți cadru-independente (fi și fd sînt mulțimi disjunctive).

Definiție 3: O *instanță obiect* este un triplet (oid, os, ip) , unde:

oid este o frază numită identitatea obiectului,

$os = (fd, fi)$ este o schemă obiect și

ip este o mulțime de afirmații de tip:

(a) pentru fiecare proprietate $(pname, Values)$, în fi , ip conține cel mult o instanță a proprietății $(pname, Values)$,

(b) pentru fiecare proprietate $(pname, Values)$ în fd și pentru fiecare cadru f al filmului, ip conține cel mult o proprietate instanță $(pname, Values)$. Această proprietate instanță este notată prin $pname = v$ în f .

Definiție 4: O schemă activitate ACT_SCH este o mulțime finită de proprietăți astfel încât, dacă $(pname, Values1)$, și $(pname, Values2)$ ambele aparțin ACT_SCH , atunci $Values1 = Values2$.

Definiție 5: O activitate este o pereche, care constă din

(a) AcID, indecele schemei activitate ACT_SCH și

(b) pentru fiecare pereche $(pname, Values) \in ACT_SCH$ este valabilă ecuația de forma $pname = v$, unde $v \in Values$.

Oricărei activități i se asociază o schemă de activitate și fiecărei proprietăți i se asociază o valoare din mulțimea valorilor posibile.

Fiind dată o singură dată VIDEO v , putem defini "conținutul" filmului v .

Definiție 6: Fie că $framenum(v)$ specifică numărul total de cadre din filmul v . *Conținutul* lui v constă dintr-un triplet (OBJ, AC, λ) , unde:

1. $OBJ = \{oid_1, \dots, oid_n\}$ este o mulțime finită de instanțe ale obiectului,
2. $AC = \{AcID_1, \dots, AcID_k\}$ este o mulțime finită de activități/evenimente și
3. λ este o hartă de la $\{1, \dots, framenum(v)\}$ până la $2^{OBJ \cup AC}$.

Intuitiv, conținutul unei date VIDEO v este teoretic descris de tripletul (OBJ, AC, λ) , unde:

1. OBJ reprezintă mulțimea obiectelor de interes în film,
2. AC reprezintă mulțimea activităților de interes din film și
3. λ reprezintă obiectele și activitățile, care sunt asociate cu fiecare cadru f al filmului.

4.3. VIDEO biblioteca

O persoană interesată de obținerea unei lecții imprimată pe o casetă video ar dori să chestioneze o VIDEO bibliotecă, care găzduiește o colecție de casete video, referitoare la un anumit subiect. De exemplu, Universitatea Maryland oferă cursuri, utilizând contactul prin satelit. În viitor casetele video, create în acest fel, vor putea fi accesate cu ajutorul unui calculator, oferind astfel studenților prelegeri pentru diferite obiecte adunate în mai mulți ani și pînute de diferiți lectori. Chestionarea bazei de date VIDEO de un student individual ar presupune accesarea unui număr foarte mare de casete video.

O bibliotecă VIDEO este o colecție, care specifică: (a) totalitatea filmelor din bibliotecă, (b) conținutul fiecărui film și (c) memorizarea fizică a filmelor.

Definiție 7: O VIDEO bibliotecă *VidLib* constă dintr-o mulțime finită de cvintete de tip $(VidContent, Vid_Id, framenum, R, plm)$, unde:

- (a) *VidContent* este conținutul filmului,
- (b) *Vid_Id*, este numele filmului,
- (c) *Framenum* este numărul de cadre în film,
- (d) *Plm* este amplasarea, care specifică adresele diferitor părți ale filmului și
- (e) *R* este mulțimea relațiilor despre filme în întregime.

4.3.1. Chestionarea bibliotecii VIDEO

Chestionarea unei VIDEO bibliotecii conține următoarele tipuri de interogări: (a) *aportarea segmentelor* (Găsește toate segmentele care corespund unei anumite cerințe), (b) *aportarea obiectelor*, (c) *aportarea activităților* și (d) *aportarea proprietăților de bază* (Care VIDEO-date sunt în bibliotecă, care este conținutul fiecărei VIDEO-date selectate, unde sunt localizate fizic VIDEO-datele).

4.3.2. Funcțiile VIDEO-datei

Cu bibliotecile VIDEO pot fi definite o serie de funcții:

FindVideoWithObject(o): fiind dat numele obiectului o , această funcție ne oferă tripletul $(VideoId, StartFrame, EndFrame)$,

FindVideoWithActivity(a)

FindVideoWithActivityandProp(a,p,z)

FindVideoWithObjectandProp(o,p,z)

FindObjectsInVideo(v,s,e)

FindActivitiesInVideo (v,s,e)

FindActivitiesAndPropsinVideo (v,s,e)

FindObjectAndPropsinVideo (v,s,e)

O chestionare standardă a VIDEO-bibliotecii, utilizînd SQL are forma:

```
SELECT câmp1,..., câmpn
FROM relația1(R1), relația2(R2),..., relațiak(Rk)
WHERE condiție.
```

4.3.3. Ordonarea datelor VIDEO

O problemă importantă este crearea structurilor informaționale, care ar organiza bazele de date VIDEO în așa fel încât să optimizeze procesarea celor opt funcții enumerate mai sus. Este imposibil de a depozita conținuturi al VIDEO-datelor cadru cu cadru, deoarece un singur film de 90 minute conține 162,000 cadre. Astfel este necesar de a crea reprezentări compacte a conceptului de conținut video. În acest sens vom prezenta două astfel de structuri: (a) arborii segment cadru, și (b) arborii R-segment.

4.3.4. Arborii segment cadru

Ideea de bază a arborelui segment cadru este foarte simplă. La început se creează două tabele unidimensionale: OBJECTARRAY și ACTIVITYARRAY.

În acest context arborele poate fi creat în 2 etape:

La prima etapă presupunem, că $[s_1, e_1), \dots, [s_w, e_w)$ sunt toate intervalele în coloana "Segment" a tabeli segment. Fie q_1, \dots, q_z o enumerație ascendentă a tuturor membrilor $\{s_i, e_i \mid 1 \leq i \leq w\}$. Dacă z nu este exponent al numărului 2, atunci se procedează astfel: fie r cel mai mic număr întreg așa ca $2^r > z$ și $2^r > \text{framenum}(v)$. Se adaugă noi elemente q_{z+1}, \dots, q_{2^r} în așa fel, că $q_{2^r} = \text{framenum}(v) + 1$ și $q_{z+j} = q_z + j$ ($j > 0, z+j < 2^r$).

La a doua etapă arborele este unul binar format după cum urmează:

1. În fiecare nod arborele segment cadru reprezintă o secvență de cadru $[x,y)$.
2. Fiecare frunză este la nivelul r . Prima frunză din stînga marchează intervalul $[z_1, z_2)$, a doua $[z_2, z_3)$ și așa mai departe.
3. Numărul din interiorul fiecărui nod este adresa aceluiași nod.
4. Mulțimea de numere de lângă nod marchează numărul de identitate al VIDEO-obiectelor și a VIDEO-activităților, care apar în întreaga secvență de cadru asociată cu nodul dat.

Definiție 8: O *secvență de cadru* este o pereche $[i, j]$, unde $1 \leq i \leq n$ și $[i, j]$ reprezintă mulțimea tuturor cadrelor între i (inclusiv) și j .

Definiție 9: O *ordonare parțială* \subseteq asupra mulțimii tuturor secvențelor de cadru este definită ca $[i_1, j_1] \subseteq [i_2, j_2]$ cu condiția, că $i_1 < j_1 = i_2 < j_2$.

Definiție 10: O mulțime X de secvențe de cadru este *bine aranjată* dacă:

1. X este finită (adică $X = \{[i_1, j_1], \dots, [i_r, j_r]\}$, pentru oricare $r \geq 1$)
2. $[i_1, j_1] \subseteq [i_2, j_2] \subseteq \dots \subseteq [i_r, j_r]$

Definiție 11: O mulțime X de secvențe de cadru este *solidă* dacă:

1. X este bine ordonată
2. Nu există nici o pereche de secvențe de cadru în X de forma $[i_1, i_2]$ și $[i_2, i_3]$

4.3.5. Operații cu arborii segment cadru.

Fiecare film v este o structură de VIDEO-date, care constă dintr-un arbore segment cadru, un tablou obiect și un tablou activitate. În particular, dacă biblioteca *VidLib* conține filmele v_1, \dots, v_n , atunci este suficient să asociem următoarele:

1. O singură tabelă numită INTOBJECTARRAY cu schema (VID.ID, OBJ, PTR),
2. O tabelă numită INACTIVITYARRAY cu schema (VID.ID, ACT, PTR) și
3. Pentru fiecare arbore segment cadru v_i , $fst(v_i)$ este asociat cu filmul v_i .

De asemenea pot fi exprimate cele 8 funcții, introduse în SQL mai sus. De exemplu, una din aceste funcții FindVideoWithObject(o), poate fi implementată cu arborii segment cadru PRINTR-o operație de selecție, efectuată asupra INTOBJECTARRAY DE TIP:

```
SELECT VIDEO_ID
FROM INTOBJECTARRAY
WHERE OBJ = o.
```

4.3.6. Arborii R-segment (RS-arbori)

Arborii R-segment sunt foarte asemănători cu arborii segment cadru, cu o singură deosebire. Deși conceptele de OBJECTARRAY și ACTIVITYARRAY rămân aceleași, în locul utilizării unui arbore segment cadru pentru a reprezenta secvența de cadru profităm de faptul că o secvență $[s, e]$ este un dreptunghi cu lungimea laturii (e-s) și lățimea 0. Fiecare nod va avea o structură specială pentru a specifica, pentru fiecare dreptunghi, care obiect sau activitate este asociată acestuia.

4.4. Operații cu VIDEO-clipuri

Un film este creat prin filmarea unor secvențe și combinarea lor, utilizând un operator de combinare. O secvență este de obicei filmată de mai multe camere, fiecare având o viteză relativă de rotație constantă. În general o secvență poate avea mai multe atribute asociate așadar ca durata filmării, tipul de cameră utilizat și altele.

Un operator de combinare a filmărilor, deseori numit *edit effect*, este o operație care în baza a două filmări S_1 și S_2 , și a unui interval de timp t efectuează o secvență compusă în timpul t . Așadar un film este creat prin combinarea unei mulțimi de secvențe filmate, utilizând un număr finit de operații de compunere. Exemple de astfel de operații de compunere a filmelor includ:

1. Concatenarea filmărilor,
2. Compoziția spațială și
3. Compoziția cromatică.

4.5. Standardele video

Deși în general standardele industriale nu sunt parte componentă a fundației cadrului MULTIMEDIA este important în linii generale să explicăm ideea de bază a standardelor MPEG.

Toate standardele de comprimare a informației VIDEO încearcă să comprime filmele prin executarea unei analize intra-cadru: fiecare cadru este divizat în blocuri, diferite cadre sunt comparate pentru a vedea, dacă informația conținută de acestea, nu se repetă în două cadre. Calitatea tehnicii de compresie este măsurată conform următorilor trei parametri de bază:

- (a) Fidelitatea hărții color: cât de multe culori ale filmului original sunt prezente după comprimare?
- (b) Rezoluția pixel pe cadru: câte pixele au fost abandonate?
- (c) Numărul de cadre pe secundă: câte cadre au fost abandonate?

4.6. Scenarii de utilizare a VIDEO-dicționarului

Dicționarul MULTIMEDIA al limbii române cuprinde peste 70000 de cuvinte din cele mai diverse domenii. Dicționarul este conceput atât pentru studenți cât și pentru cercul larg al vorbitorilor limbii române, care doresc să cunoască sensul propriu care trebuie conferit cuvintelor. Dicționarul MULTIMEDIA satisface cerințele de bază: da definiția exactă a cuvântului, și, dacă e cazul, genul, numărul, sinonimele, antonimele, imagini, secvențe VIDEO și AUDIO, care exprimă sensul exact și limpede, deplin accesibil, ceea ce constituie partea cea mai importantă de utilizare. Acest dicționar este una din pietrele de temelie ale culturii tineretului, care va contribui la opera de culturalizare a maselor prin inițierea în folosirea limbii române informatizată corectă, exactă și unitară.

Compartimentul VIDEO al acestui dicționar MULTIMEDIA al limbii române conține, după pronosticurile noastre, peste 12000 cuvinte. Acest compartiment furnizează informații necesare referitoare la cuvintele căutate, secvențe video ce oferă posibilitatea de a percepe mai bine esența cuvintelor. Diviziunea video face dicționarul mult mai accesibil și atractiv utilizatorilor de toate vârstele și interesele.

Necesitatea utilizării VIDEO-dicționarului poate apărea în cele mai diverse situații. Să considerăm situația, în care un student este nevoit să scrie un referat la merceologia și tehnologia produselor alimentare. Studentul trebuie să analizeze procesul tehnologic de producere a pâinii. În acest sens, apelarea la VIDEO-dicționarul limbii române îi va ușura lucrul; acesta îi va furniza secvențe VIDEO, ce prezintă procesul de fabricare a pâinii, ingredientele utilizate, utilajul necesar.

4.5.1. Chestionarea Video dicționarului

Dicționarul VIDEO este organizat ca o mini - bibliotecă VIDEO. După cum am subliniat mai sus, în procesul de chestionare cele mai importante aspecte sunt:

- (a) Aportarea segmentelor: utilizatorul poate cere bazei de date VIDEO să-i ofere toate secvențele, care conțin informații despre procesul tehnologic de producere a pâinii. O astfel de chestionare ar fi: "Găsește toate secvențele unde se combină ingredientele ", sau "Găsește toate secvențele unde se frământă pâinea".
- (b) Aportarea obiectelor: în acest caz, utilizatorul poate solicita toate segmentele, în care este prezent cuptorul, banda rulantă sau chiar șeful departamentului de producere. Formularea întrebării ar fi: "Găsește toate secvențele, în care apare cuptorul", "Găsește toate secvențele, în care apare banda rulantă", sau "Găsește toate secvențele, în care apare șeful departamentului de producere".
- (c) Aportarea activităților: se solicită prezentarea tuturor segmentelor, în care pot fi urmărite diferite operațiuni de producere. Întrebarea poate fi: "Găsește toate secvențele, în care se desfășoară operațiunile de producere".

4.5.2. Utilizarea bazelor de date VIDEO on diferite domenii.

După cum am menționat anterior scopul baze de date VIDEO este de a satisface cele mai diverse cerințe. Astfel aceste BDV își găsesc aplicarea în cele mai diverse domenii.

4.5.2.1. Educație. Bazele de datele VIDEO au o aplicare largă în educație și cercetare. Universitățile pot acorda a o servicii ca studii la distanță prin satelit, sau utilizând Internetul. Acestea pot pune la dispoziția studenților un set de casete VIDEO cu înregistrări ale cursurilor. Dicționarul VIDEO, fiind și el o bază de date VIDEO pune la dispoziția utilizatorilor secvențe VIDEO, care pot fi utilizate în cadrul comunicărilor, pentru pregătirea unor prezentări, lecții deschise, rapoarte.

4.5.2.2. Sport. Sălile de Sănătate oferă baze de date, în care sunt înregistrate casete **VIDEO**, care conțin diferite programe de antrenament, utilizatorului oferindu-se posibilitatea de a alege între programe de slăbire, fortificare sau menținere a condiției fizice.

4.5.2.3. Agricultură. Institutetele de cercetări științifice în domeniul agriculturii din țară ar putea utiliza **VIDEO** dicționarul pentru a studia mai aprofundat procesul de plantare, condițiile de creștere și dezvoltare a plantelor, specificul dezvoltării plantelor în diferite regiuni sau țări, aclimatizarea plantelor la condițiile țării în cauză.

4.5.2.4. Economie. **VIDEO**-dicționarul poate fi utilizat în foarte multe domenii ale economiei: finanțe, contabilitate, management, marketing, statistică, turism. Vocabularul economic cuprinde destul de mulți termeni, care pot fi redați printr-un limbaj **VIDEO** mai accesibil atât specialiștilor cât și utilizatorilor de rând.

V. Concluzii

5.1. Compartimentul TEXT. Dicționarul economic TEXT al limbii române în forma sa de BDT, ca o subdiviziune a MDLR, are posibilitățile de a fi extins cu caracteristicile respective ale MULTIMEDIA: Imagine, Audio, Video etc. Această BDT va ocupa aproximativ 18 MB memorie. La conferința tinerilor savanți ai ASEM din 4-5 aprilie 2002 în baza câtorva sute de articole din DEI au fost demonstrate caracteristicile de utilizare prietenoasă a subsistemului TEXT al SGBDMM, utilizând sistemele Ms ACCESS – 2000, Ms WORD - 2000 și Ms PowerPoint – 2000 ca componente ale Software-ului Ms OFFICE –2000 și WINDOWS – 2000, exploatate în baza Hardware-ului de tip PC Pentium II, conectat la rețelele Intranet, Externet și Internet.

5.2. Subsistemul AUDIO. Subsistemul AUDIO interacționează cu celelalte subsisteme de nivel unu (TEXT, IMAGINI, VIDEO) ale SGBDMM, care susține evaluarea Marelui Dicționar al Limbii Române informatizat cu MULTIMEDIE. Acest subsistem AUDIO susține toate definițiile celor 61635 de articole din DEI de comun acord cu subsistemul TEXT al SGBDM. Cele 2320 de ilustrații din DEI sunt susținute de componenta IMAGINI a SGBDMM, dar cu ele poate fi extinsă componenta TEXT și/sau componenta AUDIO. Exemplele, enumerate mai sus de utilizare a AUDIO componenteii a MDLR informatizat, au un aspect comun, abstract vorbind formează corpul unei date, care sunt individual executate în diferite probleme prin intermediul diferitor suporturi ale Software-ului și Hardware-ului modern. Baza de date BDA al compartimentului AUDIO-dicționarului economic al MDLR informatizat va ocupa un volum de memorie de circa 60 GB memorie.

5.3. Subsistemul IMAGINI. BDI al subsistemului IMAGINI al MDLR informatizat recent a fost expusă pentru analizare și discuții la Conferința tinerelor cercetători ai ASEM din 4-5 aprilie 2002 în baza câtorva zeci de articole din DEI. Mijloacele Software-ului și Hardware-ului de tip Ms ACCESS-2000, Ms WORD-2000 și

Ms PowerPoint-2000 cu dispozitivele respective al PC-ului Pentium II au fost suficiente la etapa inițială pentru a demonstra eficiența și eficacitatea mijloacelor și metodelor alese pentru realizarea Proiectului “Limba Română – Limba a Comunității Europene” de către grupul de cercetători – autori ai acestei publicații. Volumul BDI de prezentare în Ms ACCESS-2000 fără comprimare a 50 articole din DEI ocupă circa 550 MB memorie.

5.4. Subsistemul VIDEO. După cele menționate mai sus putem să subliniem, că subdicționarul **VIDEO** are o utilitate mare pentru persoanele ce operează în diferite domenii așa ca: economia, educația, sport, agricultură, industrie, etc. Avantajul acestui dicționar este că putem ușor funcționa cu el și este accesibil pentru toți. Dicționarul **VIDEO** este o bază de date, cu care putem opera oricând avem nevoie și oferă posibilitatea de a percepe o informație în formă de videoclipuri. În așa mod persoanele ce se folosesc de astfel de dicționar înțeleg mai ușor sensul cuvântului, care este reprezentat în formă **VIDEO**, fiindcă se formează o imagine amplă despre cuvântul dat și este ușor de memorizat. Acesta încă odată confirmă proverbul: «Mai bine o dată să vezi, decât de o sută de ori să auzi».

5.5. Lucrări paralele și perspective. În paralel cu sistemele de nivel unu sunt elaborate sistemele de nivelul doi, care suportă subdiviziunile MDLR în planurile: TEXT&AUDIO, TEXT&IMAGINI și TEXT&VIDEO.

Elaborarea sistemului, care suportă în comun toate susnumitele compartimente MULTIMEDIA ale MDLR informatizat, constituie a treia platforma, mai complexă, de experimentari și implementări a dicționarelor computerizate în cadrul elaborării MDLR informatizat [17].

Rezultatele evaluării preventive a primelor din aceste trei platforme: sistemele unare TEXT, AUDIO, IMAGINI și VIDEO au dat posibilitate de a face unele concluzii de evaluare a MDLR informatizat ca o parte componentă a cercetărilor, făcute în cadrul Proiectului «Limba română – limbă a Comunității Europene», care evaluează în perioada 2000-2006. Acest Proiect a fost inițiat [10-11] de către Forumul Internațional din Chișinău, 14-15 aprilie anul 2000. Proiectul constituie unul din subiectele de cercetări, experimentări și evaluări, efectuate în cadrul Consorțiului Uniunii Latine «Pentru limba română», a Consorțiului «Pentru informatizarea limbii române» și a Comisiei Academiei Române «Pentru informatizarea limbii române».

O serie de aplicații a MDLR computerizat este evidențiată în [13-16].

Referințe bibliografice

- [1] V. S. Subrahmanian. Principles of Multimedia Database Systems. // Morgan Kaufman Publishers, Inc., San-Francisco, California, USA, 1998, -pp. 442.
- [2] D. Todoroi, S. Nazem, T. Jucan, D. Micusha. Transition To A Full Information Society: Stage Development. // Working Paper No. 98-2, UNO, Omaha, USA, March 1998. - 38 p.

-
- [3] D. Todoroi, D. Micu^a, V. Clocotici, I. Linga, V. Tapcov, N. Drucioc, A. Calcatin, M. Morari. **Data Bases** and Communications Tools. Ms ACCESS – 200. // Ed. ASEM, Chisinau 2002, 337 pages. (Eng.)
- [4] Dumitru N. Todoroi, Zinaida Todoroi, Diana Micusa. Romanian Computerized Language – One of the European Community Languages. // Proceedings of the 26th Annual Congress of the American Romanian Academy of Arts and Sciences (ARA), Montreal, Quebec, Canada, July 25-29, 2001, pp. 133-137. (Rom)
- [5] Diana D. Micusha, Dumitru Todoroi. Natural language processing at the transition to a full information society initial development phase. Part 1. // Studii ^oi cercetări economice. Vol. XXX. Lucrări prezentate la Sesiunea jubiliară de comunicări ^otiințifice : «Cre^otere economică, dezvoltare, progres», Cluj-Napoca, 2001, pp. 1396-1413.
- [6] Diana D. Micusha, Dumitru Todoroi. Natural language processing at the transition to a full information society initial development phase. Part 2. // Studii ^oi cercetări economice. Vol. XXX. Lucrări prezentate la Sesiunea jubiliară de comunicări ^otiințifice : «Cre^otere economică, dezvoltare, progres», Cluj-Napoca, 2001, pp. 1414-1427.
- [7] Sabin-Corneliu Buraga, Dumitru Todoroi. Adaptabilitatea informațională ^oi operațională. // Studii ^oi cercetări economice. Vol. XXX. Lucrări prezentate la Sesiunea jubiliară de comunicări ^otiințifice : «Cre^otere economică, dezvoltare, progres», Cluj-Napoca, 2001, pp. 1447-1457.
- [8] Dumitru TODOROI. The Computerized Romanian Natural Language Processing Development-Projects-Perspectives. // INFORMATION SOCIETY. The Proceedings of the 5th International Symposium on Economic Informatics, May 2001, Ed ECONOMICA, Bucharest 10-13 May 2001, pp. 927-935.
- [9] Dumitru N. TODOROI. IEE-2000 PROJECT: Natural Language Processing Initialization. // EUROPEAN EXCELLENCE IN BUSINESS STUDIES STUDENTS' EDUCATION. International Symposium. Edited by IOAN ANDONE, Bucuresti, Editura Economica, 2000, pp. 328-334.
- [10] Dumitru Todoroi. Project: Romanian Language - One of the European Community Languages. // Proc. of the VI Conf. « Application Sciences», 18-19 May 2000, USAM, Chisinau, pp. 12-15.
- [11] Dan Crisrea, Dumitru Todoroi, Dan Tufis. Computational Linguistic: Romanian Language - One of the European Community Languages. // Proc. of the Intern. Sc. Seminar "Strategies and Modalities for Romania and Moldova' European Integration", 28-29 Sept. 2000, V.2, ASEM, Chisinau, pp.276-280.
- [12] D. Todoroi, D. Micusha, V. Clocotici, S. Pereteatcu, V. Bordeianu, C. Grigoras, S. Cretu, I. Linga, S. Spataru. Natural Language Processing: IEE-2000 Project. // Proc. of the Intern. Sc. Seminar "Strategies and Modalities for Romania and Moldova' European Integration", 28-29 Sept. 2000, V.2, ASEM, Chisinau, pp.281-285.
- [13] Stefan Spataru, Dumitru Todoroi. Distance Education Via Internet, Multimedia and modern System Environment. // Proc. of the Intern. Sc. Seminar "Strategies and Modalities for Romania and Moldova' European Integration", 28-29 Sept. 2000, V.2, ASEM, Chisinau, pp. 307-312.

- [14] Ion LINGA. IMPACTUL IMPLEMENTARII COMPUTERULUI ASUPRA PROCESULUI DE ASIMILARE A CUNOSTINTELOR. // Proc. Of the 27th ARA Congress, May 29 – June 2, 2002, Oradea, Romania.(To be published).
- [15] Ion COVALENCO. Metode adaptabile de evaluare a cunostiințelor asistată de calculator. // Proc. Of the 27th ARA Congress, May 29 – June 2, 2002, Oradea, Romania. (To be published).
- [16] Nicolae OBJELEAN.The Metod for Error Corection in String with Applications in Speach Recognition. // Proc. Of the 27th ARA Congress, May 29 – June 2, 2002, Oradea, Romania.(To be published).
- [17] Dumitru N. TODOROI, ASEM, Chisinau, Nicolae MARGINEANU, L'Ecole Politechnique, Montreal, Canada.THE ROMANIAN LANGUAGE'MULTIMEDIA – DICTIONARIES IMPLEMENTATION ENVIRONMENT AT THE FULL INFORMATION SOCIETY INITIAL DEVELOPMENT PERIOD. // Proc. Of the 27th ARA Congress, May 29 – June 2, 2002, Oradea, Romania.(To be published).
- [18] Diana MICUSHA. Mijloace adaptabile ale sistemelor de procesare a limbajului natural computerizat. // Proc. Of the 27th ARA Congress, May 29 – June 2, 2002, Oradea, Romania.(To be published).
- [19] Zinaida TODOROI, ULIM, Chisinau, Eugenia MARGINEANU, L'Ecole Politechnique, Montreal, Canada. MULTIMEDIA – dictionaries for Romanian Language. Usage Scenarios on the EAPEC Base. // Proc. Of the 27th ARA Congress, May 29 – June 2, 2002, Oradea, Romania.(To be published).
- [20] Societatea informațională – Societatea cunoașterii. Concepte, soluții și strategii pentru România. // ACADEMIA ROMÂNĂ, Editura EXPERT, București, decembrie 2001. – 541 pages.
- [21] Dicționar Enciclopedic Ilustrat (DEI). // Editura CARTIER SRL, Chișinău, Editura CODEX SRL, București, 1999, 1808 pages.
- [22] Beyond Calculation : The Next Fifty Years of Computing. // Edited by Peter J. Denning and Bob Metcalfe, Copernicus, 1997 Springer-Verlag New York, Inc., 350 pages.

Comunicări la Conferința tinerilor cercetători ASEM, 4-5 aprilie 2002, Chișinău.

Coordonator: Dumitru TODOROI, Prof. Univ., doctor habilitatus.

- [23]. AMBROZII Nadejda, GODZINA Irina. Componenta Text a Audio Dicționarului Economic al Limbii Române.
- [24]. ȚURCANU Virginia, MUTRUC Carolina. AUDIO-DICȚIONARUL EXPLICATIV ECONOMIC AL LIMBII ROMÂNE
- [25]. COZLOV Elena, BABANU Irina. Subsistemul IMAGE al dicționarului economic informatizat al limbii române.
- [26]. LUNGU Stela,CIOBANU Diana, GUZUN Oxana. VIDEO-dicționarul economic al limbii române.

Mediu pentru editarea transcrierilor fonetice în Limba Română. Realizarea Atlasului Lingvistic Român pe Regiuni

Silviu BEJINARIU, Vasile APOPEI, Mariana ROMAN
Academia Română, Institutul de Informatică Teoretică, Iași, B-dul Carol nr. 8
silviub@academie.is.edu.ro, vapopei@academie.is.edu.ro

Abstract

The goal of our work is to create an Electronic Linguistic Atlas of Romania. The Electronic Linguistic Atlas has features of a multimedia application allowing the user to consult and/or print the linguistic maps and to listen audio recordings or synthesized speech.

In order to show all the spelling variations, the phonetically transcription is used in the linguistic atlases. For the Romanian Language, the graphic symbols have been handwritten.

The editing process is too difficult using a standard text editor as consequence of the great number of fonts used. In this paper we propose an editing interface for the phonetic transcription of the Romanian Language. This interface can be used to edit dictionaries of the Linguistic Atlas and as editing tool for the phonetic transcriptions in stand-alone mode or as server for other text editors.

Keywords: dictionary, phonetically transcription, multimedia, linguistic atlas

1. Clasificarea simbolurilor grafice pentru editarea transcrierilor fonetice

Pentru a putea arăta toate nuanțele de rostire, în lingvistică se recurge (după practica internațională) la transcrierea fonetică. Pe lângă transcrierea fonetică internațională realizată cu Alfabetul Fonetic Internațional (IPA), fiecare țară își are propriile simboluri grafice [1], [2]. Pentru limba română, aceste simboluri sunt realizate doar manual. În lucrarea [3] este prezentată o primă abordare a realizării simbolurilor grafice pentru transcrierea fonetică din perspectiva realizării variantei computerizate a atlaselor lingvistice românești.

În această primă parte vom prezenta principiile care au stat la baza modului în care au fost organizate simbolurile grafice folosite în transcrierea fonetică a limbii române.

Pentru claritatea prezentării introducem următoarele noțiuni:

- sunete primare⁹⁸:
 - vocale, consoane - existente în alfabetul latin care au corespondent pe tastatură;
 - diacritice - vocale, consoane – care nu au corespondent pe tastatură dar pot fi obținute prin combinații de taste;
- sunete marcate cu unul sau mai multe fenomene fonetice.

De aici a rezultat necesitatea realizării unui font de bază (*ALR_Baza*) care să cuprindă simbolurile grafice pentru toate sunetele primare. Poziția în “font” a simbolurilor grafice pentru diacritice, a fost stabilită urmărind păstrarea poziției implicite din familiile de fonturi uzuale (*Arial, Times New Roman*). Pentru realizarea sunetelor marcate cu un fenomen fonetic sau mai multe am realizat familii de fonturi ale căror denumiri le-am dat cu ajutorul fenomenelor fonetice aplicate (ex. *ALR_Semivocale, ALR_Nazalizate, ALR_Seminazalizate, ALR_ScurteNazalizate, ALR_ etc.*). Această organizare a fonturilor a fost făcută cu scopul de a permite scrierea textelor cu transcrieri fonetice cu orice editor de text (*Microsoft Word*), iar textul scris cu aceste fonturi să poată fi citit chiar dacă fonturile proiectate de noi nu sunt instalate (în acest context se vor pierde numai fenomenele fonetice aplicate sunetelor primare).

Pentru generarea acestor fonturi am folosit programul *FontLab 3.1* care permite definirea de simboluri grafice compuse, pornind de la o familie de fonturi *TrueType* existentă în sistemul de operare Windows. Pentru familiile de fonturi pe care le-am realizat am convenit să folosim ca model de plecare fontul *ARIAL*.

Facem precizarea că fenomenele fonetice și modul lor de aplicare este diferit pentru cele două tipuri de sunete: vocale și consoane.

⁹⁸ Formularea "sunete primare", inexactă din punctul de vedere fonetic, este folosită cu înțelesul "sunete a căror imagine grafică pe calculator are corespondent pe tastatură, sau este obținută prin combinații de taste "

1.1. Fenomene fonetice aplicate vocalelor primare

Vocalele primare folosite în transcrierea fonetică sunt:

simple	diacritice					
a	ă	ǎ	â	a ⁹⁹	í	
e	ë	ɛ̃				
i			î			i
o	ö	ɔ̃				
u	ü		ú			

Cu ajutorul acestor "vocale primare" și al celor trei variante accentuate (a - á A Á) ale fiecăreia dintre ele se obține seria completă de sunete vocalice care se regăsesc în fontul de bază *ALR_Baza* (17*4=68 semne).

Transformările fonetice care pot modifica cele 17 vocale de bază (68 împreună cu variantele lor accentuate), sunt clasificate în următoarele grupe de fenomene disjuncte¹⁰⁰:

Grupe	Poziționare	Notăție	Fenomen	Exemplu
Durată	Așezat cel mai sus	(a)	Scurtime	e é E É
		(b)	Semilungime	e é E É
		(c)	Lungime	e é E É
Nazalizare	Așezat deasupra dar sub a-c	(d)	Seminazalizare	e é E É
		(e)	Nazalizare	e é E É
Ocluzie glotală	așezat "în umăr", în fața	(f)	Coup de glotte	e é E É
Deschidere	Așezat imediat sub vocală	(g)	Închidere	e é E É
		(h)	Semideschidere	e é E É
		(i)	Deschidere	e é E É
		(j)	Deschidere mare	e é E É
Afonizare	Așezat sub vocală dar și sub g-j	(k)	Semiafonizare	e é E É
		(l)	Afonizare	e é E É

Prin transformări fonetice se înțeleg toate realizările vocalice obținute ca urmare a aplicării a cel puțin unui fenomen fonetic (*a*)–(*l*) asupra vocalelor primare. Din punct de vedere lingvistic sunt impuse următoarele reguli:

Regula [1]

- vocalele **a**, **ă**, **í** - deschise prin natura lor (cu cel mai mare grad de apertură) - nu pot contacta fenomenele fonetice *h* (semideschidere), *i* (deschidere) și *j* (deschidere mare);

⁹⁹ Sunetul *a* nu trebuie marcat cu alt fenomen fonetic

¹⁰⁰ Fenomene din aceeași grupă nu pot fi aplicate simultan asupra unui sunet

- vocalele i, □, î, u, ü, ũ – închise prin natura lor (cu cel mai mic grad de apertură) - nu pot contacta fenomenul fonetic *g* (închidere).

Sunt excluse deci variantele vocalice a, A, a, ä, ä, ä, Í, Í, Í, ca și variantele vocalice i, î, u, ü, ũ. Aceste 15 grafeme*4=60 se scad din cele 17*4*12=816. Astfel prin asocierea vocalelor primare cu câte un fenomen *a-l* apar 756 imagini grafice repartizate în 12 fonturi astfel grupate (convențional dar extrem de ușor de ținut minte) după criteriul poziției semnului față de vocală.

Regula [2]

Sunt excluse orice combinații dintre două nuanțe fonetice din aceeași grupă de transformări vocalice. Astfel o vocală nu poate fi în același timp “scurtă, semilungă și lungă” sau „seminazală și nazală” sau „închisă, semideschisă, deschisă și foarte deschisă” sau „semiafonizată și afonizată”. Fiecare transformare fonetică exclude prezența celorlalte transformări din aceeași grupă. În aceste condiții combinațiile de câte două sau mai multe fenomene sunt posibile doar între membrii a două grupe diferite.

În plus, cele 15 grafeme*4=60 excluse ca urmare a restricției formulate sub **Regula 1**, nu pot participa la combinațiile de de două, trei, patru fenomene.

1.2. Fenomene fonetice aplicate consoanelor primare

Consoanele primare folosite în transcrierea fonetică sunt:

b, c, □, □, □, □, □, d, π, D, ∩, f, g, □, □, □, □, h, J, O, O, j, k, l, ≈, m, M, n, N, φ, J, p, r, ∅, ∅, ρ, s, ∇, ∙, ș, σ, t, v, †, v, w, z, |, †, y

Fenomenele fonetice care pot fi aplicate consoanelor primare sunt:

Grupe	Notăție	Fenomen
Durată	1	Semilungime
	2	Lungime
Palatalizare	3	Semipalatalizare
	4	Palatalizare
	5	Palatalizare mare
Explozie	6	Explozie
Caracter silabic	7	Caracter silabic
Afonizare	8	Semiafonizare
	9	Afonizare

Spre deosebire de vocale, unde s-au putut defini reguli generale pentru realizarea combinațiilor de fenomene fonetice, în cazul consoanelor primare, transformările fonetice se aplică numai unor consoane specifice. În plus, consoanelor primare le pot fi aplicate numai cel mult două transformări și numai în anumite combinații. În tabelul următor sunt prezentate combinațiile posibile de fenomene și consoanele pe care acestea le pot însoți.

<i>1.2.1. Consoane cu un singur fenomen fonetic:</i>	
semilungime	∩, f, h, □, ○, ●, j, l, ≈, m, M, n, N, ∅, J, r, ℄, ℅, ρ, s, ∇, ∙, ș, v, w, z, , ̇, y
lungime	∩, f, h, □, ○, ●, j, l, ≈, m, M, n, N, ∅, J, r, ℄, ℅, ρ, s, ∇, ∙, ș, v, w, z, , ̇, y
semipalatalizare	d, g, h, ○, k, j, l, n, r, ș, t
palatalizare	d, g, h, k, j, l, n, r, ș, t
palatalizare mare	t, d
explozie	c, p, t
caracter silabic	l, m, n, r, s, M
semiafonizare	b, d, π, D, ∩, g, □, □, □, □, □, ●, j, l, ≈, m, M, n, N, ∅, J, r, ℄, ℅, ρ, v, w, z, , ̇, y
afonizare	b, d, π, D, ∩, g, □, □, □, □, □, ●, j, l, ≈, m, M, n, N, ∅, J, r, ℄, ℅, ρ, v, w, z, , ̇, y
<i>1.2.2. Consoane cu două fenomene fonetice:</i>	
semilungime + semipalatalizare	h, ○, j, l, n, r, ș
semilungime + palatalizare	h, j, l, n, r, ș
semilungime + caracter silabic	l, m, M, n, r, s
semilungime+semiafonizare	∩, □, ●, j, l, ≈, m, M, n, N, ∅, J, r, ℄, ℅, ρ, v, w, z, , ̇, y
semilungime+afonizare	∩, □, ●, j, l, ≈, m, M, n, N, ∅, J, r, ℄, ℅, ρ, v, w, z, , ̇, y
lungime + semipalatalizare	h, ○, j, l, n, r, ș
lungime + palatalizare	h, j, l, n, r, ș
lungime + caracter silabic	l, m, M, n, r, s
lungime+semiafonizare	∩, □, ●, j, l, ≈, m, M, n, N, ∅, J, r, ℄, ℅, ρ, v, w, z, , ̇, y
lungime+afonizare	∩, □, ●, j, l, ≈, m, M, n, N, ∅, J, r, ℄, ℅, ρ, v, w, z, , ̇, y
semipalatalizare+semiafonizare	d, g, j, l, n, r
semipalatalizare+afonizare	d, g, j, l, n, r

palatalizare+semiafonizar e	d, g, j, l, n, r
palatalizare+afonizare	d, g, j, l, n, r
Palatalizare	D
mare+semiafonizare	
palatalizare	D
mare+afonizare	
Explozie+semiafonizare	b, d, g
Explozie+afonizare	b, d, g
Caracter silabic+semiafonizare	l, m, M, n, r
Caracter silabic+afonizare	l, m, M, n, r

2. Mediu pentru editarea transcrierilor fonetice

Interfața realizată pentru editarea transcrierilor fonetice poate fi folosită în mai multe moduri:

- editarea dicționarelor Atlasului Lingvistic;
- editor stand-alone sau ca aplicație de tip server pentru inserarea de obiecte de tip “transcriere fonetica” în alte editoare de text.

Funcționalitatea acestei interfețe va fi exemplificată pentru situația Atlasului Lingvistic, ale cărui componente sunt prezentate pe scurt în continuare.

Dicționarele ALR sunt componente care realizează colectarea informațiilor primare despre titlul hărților (cuvinte de bază), punctele de anchetă, speech (colecție audio), transcrieri fonetice și notele asociate transcrierilor fonetice (Figura 1).

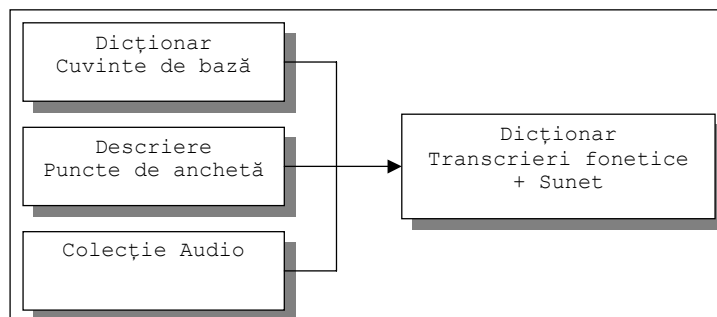


Figura 1. Conexiunile dintre informațiile stocate în dicționarele ALR

Dicționarul “Cuvinte de bază” conține fondul de cuvinte (titlul hărților) din atlasul lingvistic electronic, întrebările care au fost puse la anchetă, note, observații, și

eventual imagini. Pentru fiecare cuvânt este indicată și întrebarea corespunzătoare care este pusă în momentul interviului.

În momentul completării acestui dicționar, utilizatorul poate vedea lista completă a cuvintelor de bază introduse, le poate sorta după diferite criterii, poate modifica articolele introduse anterior, după cum este prezentat în figura 2.

Dicționarul "Puncte de anchetă" conține informații (cod, nume, observații) despre punctele de anchetă prezentate în cadrul atlasului lingvistic. La fel ca la dicționarul anterior, și aici, utilizatorul poate vedea lista completă a punctelor de anchetă introduse, le poate sorta după diferite criterii, poate modifica articolele introduse anterior.

2.1. Dicționar transcrieri fonetice

Dicționarul de transcrieri fonetice conține transcrierea fonetică a răspunsului la întrebarea pusă în etapa de interviu pentru fiecare cuvânt din **Dicționarul Cuvinte de bază** în fiecare din **Punctele de anchetă**, iar acolo unde este posibil și înregistrarea audio corespunzătoare din **Colecția Audio**.

Pentru claritatea hărților lingvistice, răspunsurile din punctele de anchetă sunt însoțite de note și comentarii (figura 3).

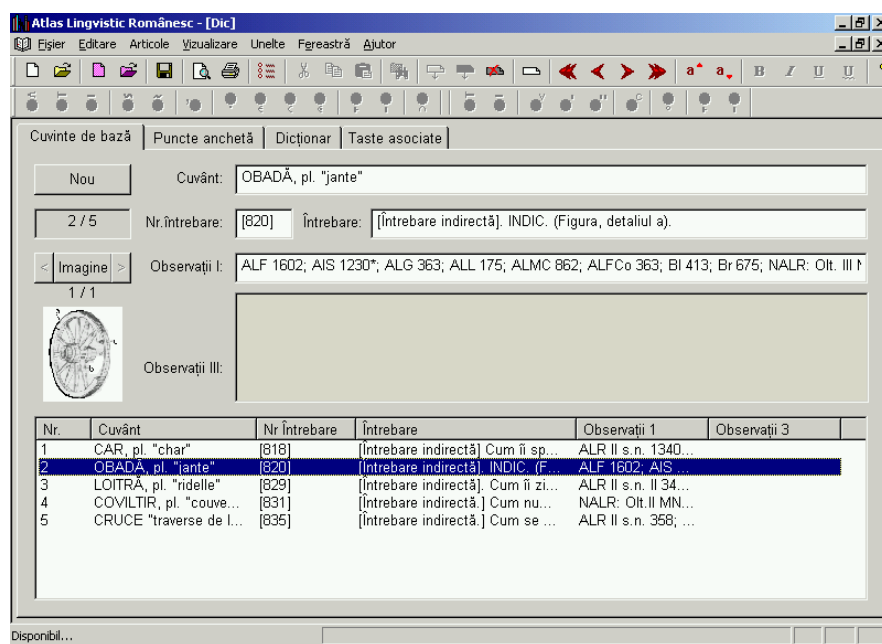


Figura 2. Fereastra de editare a listei cuvintelor de bază

Pentru transcrierea fonetică a cuvintelor din atlasul lingvistic românesc este folosit un număr mare de fonturi, rezultat din numărul mare de combinații posibile ale fenomenelor fonetice prezentate în capitolul 1. Aceste fonturi au fost definite astfel încât, toate "variantele fonetice" ale unui anumit caracter să fie obținute prin selectarea caracterului respectiv într-un anumit font.

Deoarece un fișier text normal nu păstrează informații despre fonturile folosite, și în plus transcrierile fonetice sunt realizate și prin diferite poziționări ale caracterelor, s-a folosit un mod propriu de codificare a acestora.

Transcrierile fonetice sunt codificate cu ajutorul unor obiecte de tip **CAIString**. Acestea sunt de fapt șiruri de obiecte de tip **CAIChar**, care la rândul lor au următoarea descriere:

- caracterul corespunzător sunetului primar (pe 16 biți, codificare UNICODE);
 - * atribute:
 - poziționare: normal, deasupra sau „în umăr”;
 - mod subliniere: linie sau zigzag;
 - cursiv;
 - îngroșat;
 - * fenomene:
 - tip sunet: vocală sau consoană
 - fenomene specifice aplicate (codificate pe biți);

Fontul folosit pentru desenarea caracterului din transcrierea fonetică este ales dinamic din lista de fonturi a aplicației, în momentul afișării.

În momentul în care este deschis dicționarul de transcrieri fonetice, se fac două tipuri de verificări:

- se verifică corespondența dintre fonturile folosite la ultima editare a dicționarului și lista curentă recunoscută de program.
- se verifică dacă toate fonturile folosite sunt instalate în Windows.

Datorită cantității mari de informație care trebuie stocate pentru Atlasul Lingvistic Român, descrierea fiecărui cuvânt este compresată folosind un algoritm de compresie LZW. La selecția unui cuvânt de bază, descrierea sa este decompresată în memorie și dacă se fac modificări ale transcrierilor fonetice, aceasta este compresată și rescrisă în fișier la selectarea unui alt cuvânt, sau la închiderea dicționarului.

Pentru scrierea informațiilor în dicționar am proiectat o interfață utilizator prietenoasă, la care operatorul trebuie să parcurgă următorii pași:

- selectează cuvântul titlu;
- selectează punctul de anchetă;

- editează transcrierea fonetică, nota și comentariul asociat cuvântului pentru punctul de anchetă respectiv.

La editarea transcrierilor fonetice trebuie avute în vedere două aspecte:

- selectarea sunetului primar;
- selectarea fenomenelor asociate.

Selectarea sunetului primar se face prin apăsarea tastei respective, dacă sunetul are un corespondent pe tastatură, sau prin apăsarea unei combinații de taste, dacă sunetul nu are corespondent pe tastatură. Combinațiile de taste sunt prestabilite în aplicație (la stabilirea combinațiilor de taste am folosit recomandările de la Microsoft Word), și cel puțin deocamdată nu pot fi modificate de utilizator. Pentru a veni în ajutorul celui ce editează dicționarul, aplicația dispune de o fereastră în care sunt afișate combinațiile de taste predefinite.

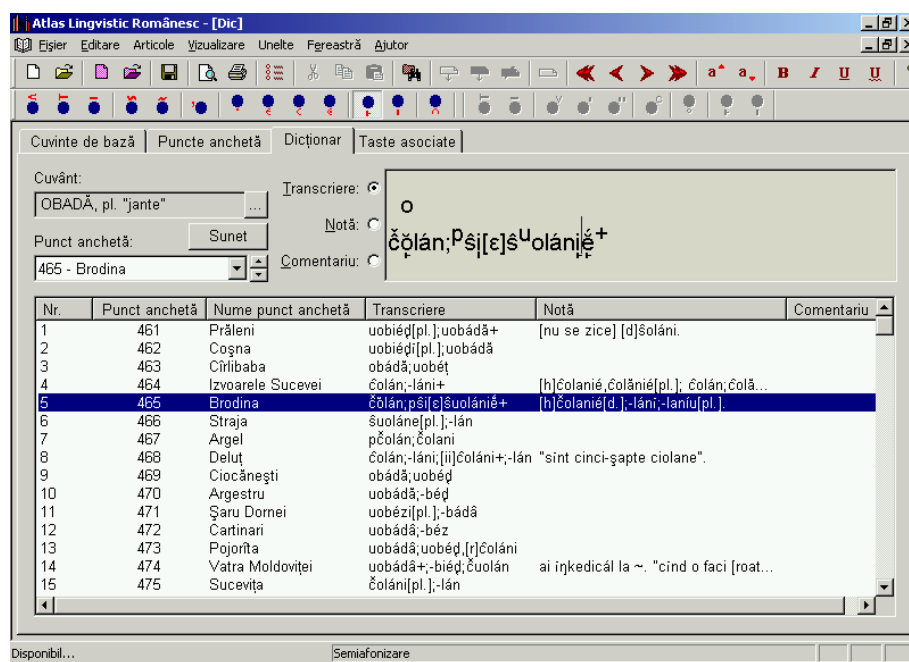


Figura 3. Editarea Dicționarului de transcrieri fonetice

Pentru selectarea fenomenelor asociate sunetelor, aplicația prezintă 2 grupe de butoane cu imaginile tuturor fenomenelor posibile pentru vocale respectiv consoane. Prin apăsarea pe unul din aceste butoane se va selecta simbolul grafic corespunzător în

transcrierea fonetică. Cele 2 grupe de butoane sunt împărțite în subgrupe corespunzătoare grupărilor de fenomene (vezi capitolul 1). Pot fi selectate mai multe fenomene, dar, cel mult câte unul din fiecare subgrupă. Selectarea unui fenomen, produce dezactivarea selecției anterioare din subgrupa respectivă.

După selectarea caracterului dorit, utilizatorul va specifica și poziționarea acestuia (deasupra, în urmă) prin folosirea comenzilor PgUp/PgDown.

Fereastra de editare a transcrierilor fonetice este prezentată în figura 3.

Dictionarul cu transcrieri fonetice poate conține înregistrările audio în format WAV ale răspunsurilor la întrebările din anchetă pentru cuvintele incluse în dictionarul lingvistic. Acest lucru este posibil dacă la realizarea anchetei pentru atlasele lingvistice se face și înregistrarea pe bandă a răspunsurilor.

3. Realizarea Atlasului Lingvistic Român pe Regiuni

Sistemul software care modelează atlasul lingvistic electronic, conține module care realizează gestionarea următoarelor grupe de informații:

- simbolii pentru editarea transcrierilor fonetice,
- dicționarele atlasului lingvistic (cuvinte de bază, puncte de anchetă, transcrieri fonetice),
- informații grafice pentru descrierea hărților organizate în fișiere DXF,
- hărțile atlasului lingvistic, care pot fi consultate și/sau tipărite;

Din punct de vedere funcțional, atlasul lingvistic electronic este structurat în două componente principale:

- Proceduri pentru pregătirea datelor primare.
- Interfața multimedia;

Aceste componente sunt prezentate în figura 4.

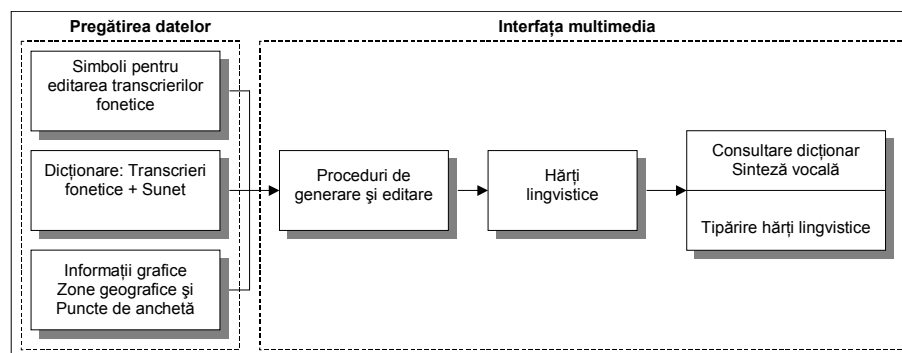


Figura 4. Componentele Atlasului lingvistic Electronic

În acest capitol ne vom referii la facilitățile pe care le oferă modulul software care realizează generarea și consultarea hărților lingvistice. În etapa de proiectare a acestui modul, am avut în vedere modelul interfețelor multimedia. Din această analiză a rezultat necesitatea existenței unui modul care să permită:

- generarea unei hărți noi pe baza informațiilor din dicționarele ALR și a informațiilor grafice primare cuprinse în fișiere DXF.
- editarea: aranjarea în pagină, selectarea informațiilor care vor fi vizibile implicit;
- salvarea într-un fișier numit “hartă lingvistică” a selecțiilor și modificărilor din faza de editare;
- consultarea atlasului electronic:
 - vizualizarea și ascultarea informațiilor din punctele de anchetă.
 - tipărirea hărților lingvistice;

3.1. Modulul pentru generarea și editarea hărților lingvistice

Pentru reprezentarea hărților lingvistice s-au proiectat structuri de date bazate pe obiecte, suficient de flexibile pentru a permite dezvoltări ulterioare. În cele ce urmează vom face o trecere în revistă a principalelor structuri de date realizate.

Pentru organizarea informațiilor grafice primare, am avut în vedere cerințele impuse de tehnologia de realizare a atlaselor lingvistice. Astfel am apelat la formatul DXF care permite o organizare a obiectelor primare pe straturi. Am realizat fișierul NALRB.DXF care conține următoarele tipuri de obiecte (straturi):

- chenare – limitele paginii și chenarele hărții;
- frontiere – conturul zonei studiate (Moldova și Bucovina);
- mijloc – indică locul de pliere al hărții, la legarea în volum;

municipii	– localitățile importante afișate pe hartă;
puncte anchetă	– dreptunghiurile în care se scriu codurile punctelor de anchetă;
transcriere fonetică	– conține dreptunghiuri pentru transcrierea fonetică;
note	– dreptunghiuri cu pozițiile predefinite pentru Titlu, Nota I, Nota II, Nota III;
zone	– delimitări zonale în jurul punctelor de anchetă

Pentru editarea și salvarea hărților lingvistice din conținutul ALR, s-a creat o structură de date flexibilă, care să permită în viitor, extinderea editării asistate de calculator a atlaselor lingvistice românești regionale la nivel național. Astfel, a rezultat o structură de date numită “harta lingvistică” de forma următoare:

- header fișier;
- listă cu descrieri obiecte;

Descrierile de obiecte au un antet care este comun pentru toate tipurile de obiecte și un corp obiect specific fiecărui tip în parte. Obiectele pot fi simple sau compuse. Un obiect compus conține la rândul lui alte obiecte simple sau compuse.

Au fost definite următoarele tipuri de obiecte:

- Text;
- AlrString (obiectul a fost definit pentru editarea dicționarului cu transcrieri fonetice);
- Dreptunghi;
- Harta cu transcrierile fonetice;
- Harta sintetică (lingvistică sau fonetică);
- Notă referitoare la continuarea transcrierile fonetice (vezi Nota II din N.A.L.R. Moldova și Bucovina);
- Notă sintetică referitoare la cuvântul titlu (vezi Nota III din N.A.L.R. Moldova și Bucovina);
- Legendă pentru harta sintetică;
- Simbol pe harta sintetică ;
- Zonă hașurată pe harta sintetică ;
- Bitmap;
- Strat DXF;

3.2. Modulul pentru consultarea atlasului electronic

Componenta pentru **consultarea atlasului**, permite încărcarea unei hărți lingvistice generate - editate în etapa anterioară. Sistemul va afișa harta regiunii respective

(în situația studiată este vorba de Moldova și Bucovina), pe care va plasa transcrierea fonetică a răspunsurilor din punctele de anchetă împreună cu notele și observațiile introduse anterior în dicționarele ALR sau de operatorul care a realizat harta (figura 5).

Dacă utilizatorul dorește să analizeze un cuvânt pentru care nu sa realizat în prealabil o hartă, dar care există în dicționarele ALR se realizează generarea automată a hărții pe care se vor plasa transcrierile fonetice a răspunsurilor din punctele de anchetă împreună cu notele și observațiile introduse anterior în dicționarele ALR. La activarea modului, apare prezentată harta regiunii cu punctele de anchetă și numele localităților pe care acestea le reprezintă, după care poate fi selectat cuvântul de bază dorit.

După ce harta lingvistică a fost încărcată / generată prin selecția unui punct de anchetă, este posibilă și redarea audio corespunzătoare transcrierii fonetice asociate acestuia (înregistrarea audio sau cuvântul sintetizat).

Tot cu ajutorul acestei componente se realizează tipărirea automată a hărților atlasului lingvistic românesc, în vederea includerii lor în volum (figura 6).

Pentru tipărirea hărților au fost prevăzute următoarele facilități:

- posibilități de selectare a informațiilor ce se vor tipări;
- tipărirea pe o pagina sau tipărirea pe două pagini cu respectarea locului de pliere al hărții indicat prin linia “mijloc”;

Dacă utilizatorul dorește tipărirea într-un mod sintetic (fără hartă), modulul poate asigura crearea unor pagini de tip MN (Material Necartografiat). Folosind această opțiune, va fi tipărită numai lista cu transcrierile fonetice corespunzătoare cuvântului selectat, ordonate după criteriul de similaritate.

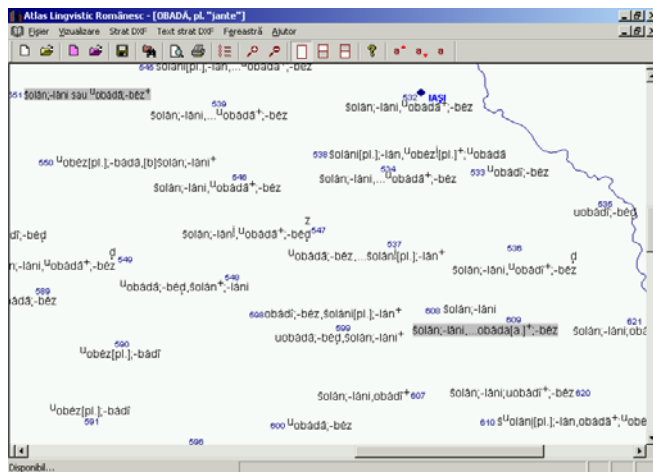


Figura 5. Fereastra de editare – consultare a Atlasului Lingvistic

4. Concluzii

Organizarea prezentată pentru simbolurile grafice facilitează editarea de texte cu transcrieri fonetice și cu alte editoare de texte care permit folosirea de fonturi multiple în text.

Modul de selectare a fonturilor folosit la editarea transcrierilor fonetice poate fi extins pentru crearea unei aplicații de tip client-server sau la realizarea unui editor simplu, de tip WordPad.

Realizarea acestui sistem de editare a transcrierilor fonetice este în curs de testare și finalizare. În continuare, ne propunem adăugarea de noi opțiuni și facilități, care să permită transformarea sistemului într-un instrument util cercetătorilor lingviști.

Bibliografie

- [1] Academia Română, Atlasul lingvistic român pe regiuni, 1987, 1997.
- [2] Istituto dell'Atlante Linguistico Italiano, Atlante Linguistico Italiano, Roma, 1995.
- [3] S. Bejinariu, M. Roman, V. Apopei, F. Olariu, "Sistem pentru editarea transcrierii fonetice în ALR ", Zilele Academice Iașene, Iasi, 6 oct. 2000.

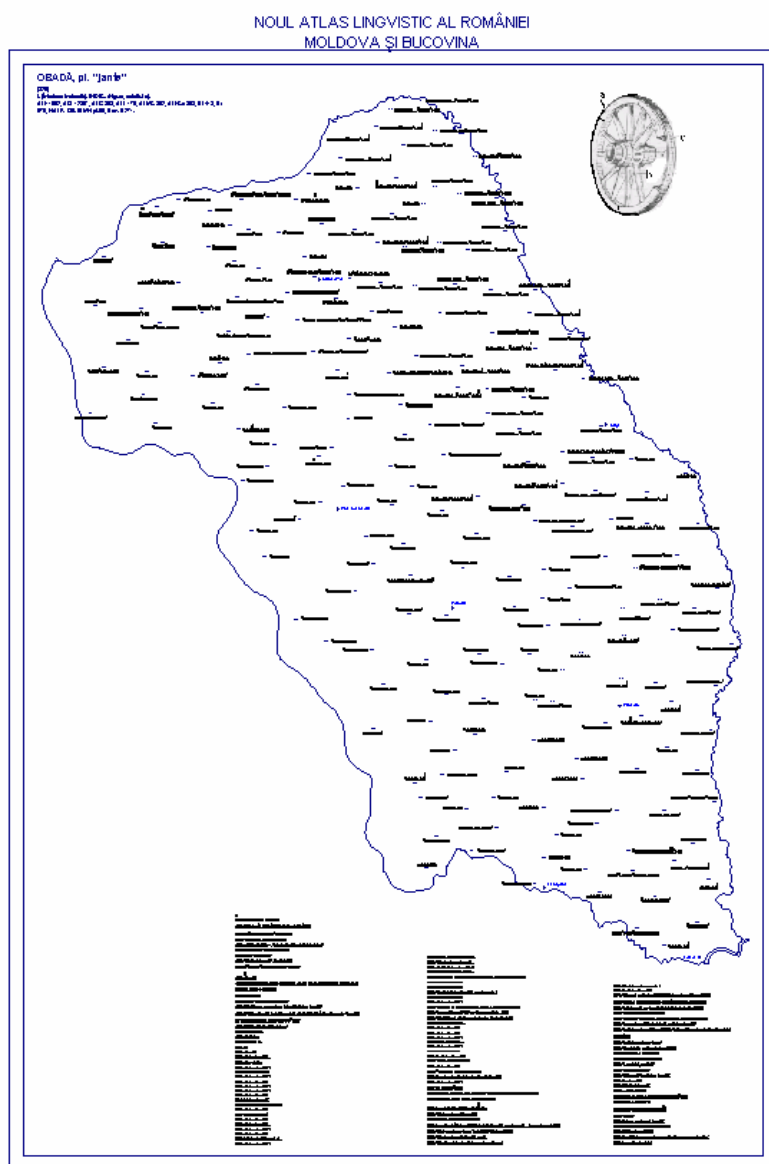


Figura 6. Imaginea unei pagini tipărite din modul de consultare

SECȚIUNEA VI

Dezbateri și discuții

Asupra a doi vectori funcționali ai societății cunoașterii: managementul cunoașterii și învățarea electronică. Cultura și societatea cunoașterii

Mihai DRĂGĂNESCU
Institutul de Inteligență Artificială
Academia Română

Introducere generală

Acest material, care constituie o contribuție la dezbateră problemelor enunțate în primul volum *Societatea informațională-Societatea cunoașterii, Concepte, soluții și strategii pentru România*, coord. Filip Gh. Florin, editat de Academia Română-Secția de știință și tehnologia informației (Institutul de inteligență artificială al Academiei Române - denumire prescurtată) și ICI-INFOSOC, Editura Expert (coordonare editorială, Valeriu Ioan-Franc), București 2002, are următorul cuprins:

- I. **Managementul cunoașterii, vector funcțional al societății cunoașterii**, comunicare (Mihai Drăgănescu) prezentată la “The Sixth International Conference on Information and Communications Technology in Public Administration, Sinaia, 29 oct.2001”.
- II. **Învățământul electronic și societatea cunoașterii**, comunicare (Mihai Drăgănescu) la simpozionul “E-learning (E-învățământ)”, Academia Română, 28 martie 2002.
- III. **Cultura și Societatea cunoașterii** (Mihai Drăgănescu, studiu elaborat în mai 2002).

Societatea cunoașterii, asupra căreia se insistă cu prioritate în aceste studii și lucrări, va fi o *perioadă interimară între Societatea informațională și Societatea conștiinței* (un studiu privind Societatea conștiinței este în elaborarea autorului.). După cum am mai remarcat în alte lucrări, esențială pentru Societatea cunoașterii va fi inteligența artificială (IA), atât ca vector tehnologic, cât și prin utilizarea ei în vectorii funcționali ai societății cunoașterii.

Această primă perioadă interimară va dura până cândva după momentul în care inteligența artificială va egala inteligența naturală (IN) structurală a omului, respectiv a

părții $(IN)_{\text{structural}}$ care nu poate poseda intuiție, creativitate și spiritualitate. După concepția mea ontologică, nu este posibil pentru orice fel de inteligență artificială (electronică și în viitor nanoelectronică) să aibă intuiție, creativitate și spiritualitate fără a recurge și la alte elemente ale naturii decât cele structurale și a căror realitate devine din ce în ce mai plauzibilă. Egalitatea $IA = (IN)_{\text{structural}}$ se va petre, după o serie de autori (Moravec, Kurzweil, Buttuzzo, Broderick ș.a.), între 2019-2035. Unii dintre aceștia cred că atunci când se va atinge $IA = (IN)_{\text{structural}}$, automat un asemenea creier electronic va avea și proprietățile fenomenologice ale intuiției, creativității și spiritualității. Ceea ce nu credem.

Din momentul în care ce $IA > (IN)_{\text{structural}}$ este evident însă că se intră într-o nouă etapă, care va produce multe consecințe pe plan social, datorită relațiilor omului cu asemenea inteligențe, unele software, altele sub forma de specii de roboți inteligenți. Aceasta va fi a doua perioadă intermediară între Societatea cunoașterii și Societatea conștiinței, până în momentul în care va apare o inteligență artificială cu conștiință veritabilă, adică o conștiință artificială (CA). Din momentul în care $CA > IN$, se va intra în zona societății conștiinței, urmând ca societatea să fie bazată pe relațiile dintre IN (care și ea este de presupus că va fi amplificată prin auto-transformări ale codului genetic și probabil prin cuplaje cu sisteme informatice microelectronice și nanoelectronice, chiar și cu rețele internet) și CA software sau robotice. Va trebui cu siguranță să gândim de pe acum și asupra societății conștiinței pentru a pregăti societatea pentru o asemenea perspectivă, care nu mai apare, surprinzător, atât de îndepărtată, deoarece se poate manifesta chiar în acest secol. Societatea cunoașterii trebuie să înceapă să fie gândită și dezvoltată și cu gândul la această viitoare societate.

I. MANAGEMENTUL CUNOAȘTERII

I.1 Introduction

In the past XXth century a new era began in the history of humanity: *the information era* [1]. This era comprises the *information society* that will be followed naturally by the *knowledge society* and finally, somewhere more or less later in this century, the *society of consciousness*. Knowledge is a form of information [2], and consciousness is another form of information [3]. All the forms of information are intermingled with the physical and energetic realities; still they have a relative independence and can influence these realities.

To pass from the first form of information society (based essentially on Internet and Internet economy) to the second stage, the knowledge society, I considered in a previous work [1], two types of vectors: *technological* and *functional*.

Technological vectors are the extended Internet, the e-book and e-document technology, artificial intelligence (with intelligent agents and future Networked Systems of Embedded Computers), nanotechnology and others.

Among the functional vectors of the Knowledge Society, a group of vectors is related to knowledge management:

- *knowledge management for corporations and enterprises, organizations and institutions, **local and national administrations**;*
- *the management of the moral use of scientific knowledge at the global level;*
- *e-learning management;*
- *development of a culture of knowledge and innovation;*
- *management of the scientific and technological knowledge for every main domain of activity as health care, sustainable society and others.*

I.2 Knowledge management

The problem of management with respect to knowledge is regarded in two ways:

- I. As the management of the organization busy with the use and integration of various types of knowledge;
- II. As the management of knowledge itself, for generation of new knowledge, for discovering existing knowledge (tacit or very local, or external to the organization), for combining available knowledge.

Perhaps, what is really needed is a general vision, in a unity, on the management of the organization and the management of knowledge.

Knowledge management for enterprises, organizations, institutions, local and national administrations

In the western literature, in the last years, were elaborated a series of works dedicated to the problems of enterprises and knowledge. In Romania we do not have yet specialists in knowledge management in the context of the knowledge society. We do not have either a knowledge society, but we need experts in knowledge management for building the future knowledge society. A group of members of the Romanian Academy and other wellknown specialists in information technology from Romania and colleagues from USA decided to constitute a Romanian-American Foundation for the Knowledge Society, one of the main aims being to educate in USA a number of young Romanian specialists in the new domain of knowledge management. All is ready for such a Foundation, the contributions of the individual founders are also ready, but no institution or organization sponsored such an exotic objective for The Knowledge Society, with some amount of money asked by the Romanian laws for a Foundation to begin its activity. But let us return to the theory of knowledge management. One definition [4] is the following:

“Knowledge Management is the conceptualizing of an organization as an integrated knowledge system, and the management of the organization for effective use of

that knowledge. Where *knowledge* refers to human cognitive and innovative processes and the artifacts that support them.”

This definition insists on the management of the organization, even if it recognizes the knowledge system of the organization. This definition, as it is recognized by its authors, disguise knowledge management because of the delicate problem of knowledge measurement [4]:

'The recent attractiveness of the term knowledge management appears to have been prompted by three major forces:

1. Increasing dominance of knowledge as a basis for organizational effectiveness.
2. The failure of financial models to represent the dynamics of knowledge.
3. The failure of information technology by itself to achieve substantial benefits for organizations.'

The second point, of the above quotation, is answered by many studies and books concerning the characteristics of the new economy based on knowledge (see for instance section 6 of [1]: The Economy of the knowledge society. The new economy. About the role of information in the new economy. The intangible goods).

The rapidity of the transformation of the information society into a knowledge society determines a reasoning on the new economy that takes into account not only:

- a) the Internet market and the effects of Internet information on all economical and administrative agents, but also
- b) the effect of knowledge as an economical and *organizational* factor that imposes the recognition of the intangible goods, in general, in the creation of economical value and organizational efficiency, and
- c) the necessity of a sustainable society, an important objective for national and even local administrations, that predictably is possible only in the frame of the knowledge society, that will demand new industries, challenges the classical economical thinking (for instance, productivity of the resources, of the energy, of materials to be more important than work productivity [8]).

The third point of the above quotation concerns the importance of the contents of information, especially of knowledge, but these would not be efficient without information technology. The technological vectors of the Knowledge Society are equally important as the functional vectors.

I.3 Points of view for practical knowledge management

Knowledge management is both the management of the organization to use knowledge and the management of all knowledge possible, from inside and outside the organization, to attain the objectives of that organization. Because knowledge is a

special form of information, information technology has to play an essential role in knowledge management. Knowledge and IT are, without any doubt, going hand in hand and have a synergetic effect on the efficiency of organizations.

Lucy Marshall [6] considers that knowledge management refers to the control and utilization of the intellectual capital in an organization. For Lucy Marshall, not the information, but knowledge is the most important asset of an institution. This author recommends a **Chief Knowledge Officer** for an institution, who based on the Intranet of the institution, has to assure the discovery and creation of knowledge in the institution.

Rooney and Mandeville considers the knowledge management at the national level. The abstract of their paper is quoted [7] below:

'As the global economy becomes more knowledge intensive and the wealth of nations more dependent on their knowledge assets being harnessed, it is essential for policy makers of having frameworks for the development and the utilization of national knowledge assets. This article argues that a policy framework can be developed through which policy initiatives in a range of policy areas can be filtered in order to meet the challenges of the knowledge economy. We have developed an approach that has previously been applied to managing intellectual capital in firms and adapted it to the public policy arena. In doing so we question policy orthodoxies such as the assumption that free trade automatically facilitates international knowledge flows, that participation in a global knowledge economy necessarily challenges national sovereignty, and that online delivery of education is necessarily a progressive strategy'.

Peter Drucker (a wellknown professor of social science at Claremont Graduate School and the author of more than thirty books, his most recent book is Management Challenges for the 21st Century, 1999) writes [8] about the **knowledge worker**:

'I am convinced that a drastic change in the social mind-set is required - just as leadership in the industrial economy after the railroad required the drastic change from "tradesman" to "technologist" or "engineer."

What we call the Information Revolution is actually a Knowledge Revolution. What has made it possible to routinize processes is not machinery; the computer is only the trigger. Software is the reorganization of traditional work, based on centuries of experience, through the application of knowledge and especially of systematic, logical analysis. The key is not electronics; it is cognitive science. This means that the key to maintaining leadership in the economy and the technology that are about to emerge is likely to be the social position of knowledge professionals and social acceptance of their values. For them to remain traditional "employees" and be treated as such would be tantamount to England's treating its technologists as tradesmen - and likely to have similar consequences.'

I.4 Cognitive Science

The knowledge of organizations is a form of knowledge that is more and more recognized. The ways and forms of this knowledge have to be carefully studied. Cognitive

science might be, indeed, the tool for this study. The cognitive science is today understood in two ways [2]:

1. As a science of human mind cognition, even if it uses models of electronic computers and electronic neural networks.
2. As a general science of cognition, that has to study cognition processes not only of the human mind, but also of animals, of artificial intelligence systems, of the ensembles man-computer-Internet, of social organizations at the levels of institutions, enterprises, corporations, local and national administrative bodies, even at the global level.

The second way of dealing with the processes of cognition presents today the greatest interest. Such a science does not yet exist. Perhaps it is on the way. The most complex realities are the social organizations because they combine all sorts of cognitive elements, natural and artificial, but they have something more, a social body, with its own social intelligence, cognition and knowledge. To obtain new theories for such large and difficult problems, it is necessary to have talented and interested specialists in knowledge management. But it is also necessary some practice of those charged with knowledge management and knowledge work in organizations. The idea of Lucy Marshall, mentioned before, about a Chief Knowledge Officer seems to be very useful.

I.4 Final remarks

The Knowledge society is paving the way for a Consciousness society. For this we need more fundamental knowledge [9] on physical reality down to the frontier of the quantum world with the deepest reality of existence, on life, mind and consciousness, on cognition, but also on self-organization and organization of social bodies and their behavior. We need also more technological knowledge. For science and society, knowledge management will become the most important administration.

References

- [1] Mihai Drăgănescu, *Societatea Informațională și a Cunoașterii. Vectorii Societății Cunoașterii* (Information Society and Knowledge Society. Vectors of the Knowledge Society), Romanian Academy, July 2001. On the Web, http://www.academiaromana.ro/pro_pri/
- [2] Mihai Drăgănescu, Cunoașterea în secolul XXI (Knowledge in the XXI century), communication at the Annual Conference of the Romanian Committee for the History and Philosophy of Science, Romanian Academy, Bucharest, 15 October 2001, to be published.

- [3] Mihai Drăgănescu, *The Interdisciplinary Science of Consciousness* (Chapter 5) pp. 46-59, in *Science and the Primacy of Consciousness, Intimation of a 21st Century Revolution*, Richard L. Amoroso a.o, (eds.), Orinda, California: The Noetic Press, 2000.
- [4] See <http://www.uts.edu.au/fac/hss/Departments/DIS/km/introduct.htm#Char>
- [5] Ernst Ulrich von Weiszäcker, Amory B. Lovins, L.Hunter Lovins, *Factor patru. Dublarea prosperității prin înjumătățirea consumului de resurse*, Raport pentru Clubul de la Roma, traducere din limba germană (FAKTOR VIER.Doppelter Wohlstand - halbierter Verbrauch, München, 1995), București, Editura tehnică, 1998.
- [6] Lucy Marshall, *Facilitating knowledge management and knowledge sharing: New opportunities for information professionals*, Online. 21(5): 92-98. 1997 Sep/Oct.
- [7] David Rooney and Thomas Mandeville, *The Knowing Nation: A Framework for Public Policy in a Post-industrial Knowledge Economy*, Prometheus 16 (4) pp. 453-467, 1998.
- [8] Peter F. Drucker, *Beyond the Information Revolution*, The Atlantic Monthly, Digital Edition, 1999, <http://www.theatlantic.com/issues/99oct/9910drucker3.htm>
- [9] Menas Kafatos, Mihai Drăgănescu, *Preliminaries to the philosophy of integrative science*, e-book, MSReader format, Academy of Scientists - Romania, Bucharest, 2001, (available free by e-mail: dragam@racai.ro).

II. ÎNVĂȚĂMÂNTUL ELECTRONIC ȘI SOCIETATEA CUNOAȘTERII

II.1 Introducere. Sintagma Societății cunoașterii.

În *societatea cunoașterii* doi vectori, strâns legați între ei, unul tehnologic - *cartea electronică* - și altul funcțional - *învățământul electronic* - sunt chemați să joace un rol important în desfășurarea acesteia.

Problematica societății cunoașterii a fost abordată în țara noastră începând din anul 2001 la Academia Română [1], la Academia de studii economice [2] și de revista Diplomat-Club [3]. Primul politician român care a folosit sintagma societății cunoașterii (din anul 2001) a fost președintele României și protectorul de fapt al Academiei Române, Ion Iliescu.

Este poate interesant de amintit că în anul 1986, în lucrarea 'Tendencies of becoming' [4] (Tendințele devenirii, republicată în volumul [5]) se justifică și folosește sintagma 'societatea cunoașterii':

*"Cine nu face legătura dintre revoluția microelectronică și informațională și tendința devenirii istorice nu înțelege vremurile. Cine se opune acestei revoluții părăsește linia devenirii istorice. Și totuși nici această revoluție nu trebuie absolutizată întrucât trebuie să fie însoțită și de alte schimbări. Atunci nu ne putem fixa numai asupra ei, ci asupra unui context mai larg în cadrul căruia ea poate juca rolul principal o anumită perioadă istorică. **Tendința devenirii istorice se conturează a fi tendința către o societatea a cunoașterii, a creației și a civilizației, către o societate globală și către o societate interastrală în univers, apoi către un act cosmic în conformitate cu tendința existențială a universului. Mai aproape de noi, ca urmare a revoluției microelectronice și informatice, a unei noi revoluții industriale, se deschid perspectivele unei societăți orientate informațional...***"

Era o viziune, în acel moment, legată de o anumită filosofie pe care am dezvoltat-o în anii 1980, viziune ancorată și în realitatea electronică și informatică a ceea ce se va numi era informației.

II.2 Cartea electronică

Cartea electronică este un vector tehnologic. La Academia Română în anul 2001 s-a desfășurat un simpozion referitor la cartea electronică și s-a publicat un volum de referință sub coordonarea prof. Doina Banciu [7]. Atunci am descoperit firma de software SOFTWIN condusă de Florin Talpeș care lucrase în domeniu și avea un prestigiu internațional în producerea de cărți electronice. Softwin este participantă la elaborarea specificațiilor internaționale OPEN E-BOOK care au stabilit formatul edițiilor de cărți electronice de interes

public. Ca urmare a simpozionului a fost înființată și o librărie de software, cărți și documente electronice la Institutul Național de Cercetare-Dezvoltare în Informatică ([http:// www.e-librarie.ro](http://www.e-librarie.ro)).

Despre cartea electronică, și rolul ei pentru societatea cunoașterii în România, am expus considerațiile mele în lucrări anterioare [1b], [7] și nu voi reveni asupra lor. În schimb, voi cita doi autori, unul care a exprimat opinii înainte de apariția cărții electronice propriu-zise, altul care a participat la lansarea cărții electronice. Primul este Paul Saffo, directorul unui elevat Institut al Viitorului din California, care lucrează, foarte scump, numai pentru marile companii americane și care în anul 1988 prevedea că o carte electronică va fi mai mult decât o carte tipărită datorită posibilităților de a introduce elemente audio, video, conexiuni la informații pe rețea. El scria [8]:

' The term "electronic book" is misleading because these products are not books at all, but something new. We are living in a moment between two revolutions: one of print, four centuries old and not quite spent and another of electronics, two decades young, and just getting underway. Today's products amount to a bridge between these two revolutions...'

Al doilea este Dick Brass, Vicepreședinte Microsoft pentru dezvoltare tehnologică, care în anul 2000, an în care cartea electronică propriu-zisă decola, scria [9]:

'If you don't think eBooks will take off, remember that electronic encyclopedias have already outsold all paper encyclopedias. [...] They cost less than \$100, instead of the \$2,000 or more for fine paper encyclopedias. [...] Similarly, after the triumph of eBooks, paper books will no longer be the principal means of distributing information. But, like horses they will continue to exist for pleasure... [...] Like all transitions, the move from pBooks to eBooks will be a little painful and tentative at first. Then, in less than 20 years, eBooks will be so pervasive that we won't be able to remember living without them. [...] We are on the verge of the most exciting change to the printed word since movable type...'

Cartea electronică a decolat. Firme precum Amazon și Barnes and Noble din SUA sunt cunoscute în întreaga lume pentru modul în care au promovat-o. Ele sunt o adevărată școală pentru toți cei care conduc și vor conduce librării de cărți electronice și software, școală accesibilă gratuit prin simpla experimentare prin Internet pe web-site-urile acestor firme.

II.3 Procesul de învățare

În anul 1988 scriam despre procesul de învățare [10]:

'Înțelegerea profundă a procesului de învățare depinde de explicarea funcționării creierului și a minții omului, în ultimă instanță de înțelegerea naturii materiei vii. Cu alte cuvinte, natura intimă a procesului de învățare nu va putea fi elucidată într-o măsură într-adevăr mulțumitoare decât atunci când știința va face un nou mare pas în cunoașterea materiei. Cercetările din domeniile fizicii și biologiei, esențiale pentru elucidarea naturii materiei vii, se vor îmbina cu cele din domeniul științei informației. Activitatea creierului

este în principal o activitate informațională, iar *procesul de învățare este un proces informațional.*'

În acea perioadă știința cognitivă se găsea, este adevărat, în perioada post-behavioristă și se baza pe modelarea simbolică de tip calculator electronic, ceea ce s-a dovedit insuficient pentru înțelegerea multumitoare a proceselor cognitive mentale [11]. De aceea procesul de învățare nu era de fapt explicat și înțeles din punct de vedere științific. În anii 1990, modelarea proceselor cognitive a cunoscut aportul adus de utilizarea modelelor bazate pe rețele neuronice (de tip natural ca în creierul omului) și neural (artificiale, electronice), dar nici acestea nu au dus încă la o știință cognitivă bine constituită [11]:

COGNIȚIA:

<i>Anii 1970 și 1980</i> (modelare simbolică tip calculator)	<i>Anii 1990</i> (efectul conectivismului, rețele neuronice și neurale)	<i>Anii 2000. Ce va urma?</i> Efectul științei integrative.
---	--	--

O speranță este aceea ca în sec. XXI știința cognitivă să fie consolidată prin luarea în considerare atât a proceselor fenomenologice (qualia, experiențiale) ale minții, cât și a rolului proceselor sociale în procesele cognitive (socialul referindu-se nu numai la persoane umane, ci și la grupuri de inteligențe artificiale sau la grupuri mixte). Un asemenea mod de abordare se încadrează în viziunea unei științe numite integrative [12]. Tot în anul 1988 remarcam [10]:

'Un interes deosebit prezintă cercetările din domeniul inteligenței artificiale, domeniu care este studiat în ultimii ani și din punctul de vedere al capacității de a învăța. Studiul procesului de învățare de către inteligența artificială ar putea oferi multe elemente utile pentru înțelegerea procesului de învățare al inteligenței naturale a omului. [...] Inteligența presupune și capacitatea de a învăța. [...] Gh. Tecuci, într-o lucrare originală în care se prezintă un sistem expert la care asociază un sistem de învățare automată [13], deși constată că 'învățarea este un proces cognitiv în cea mai mare măsură necunoscut' [13], arată și demonstrează prin sistemul său că 'forme efective de învățare automată sunt posibile'.

Dintre aceste forme de învățare automată pot fi amintite [13]:

- *Învățarea pe de rost și implantare directă de noi cunoștințe* (când este mai eficient să se regăsească o cunoștință în memorie decât să se producă acea cunoștință).
- *Învățare prin instruire* (sistemul primește cunoștințe de la un profesor și le integrează cu cunoștințele anterioare).
- *Învățarea prin analogie.*
- *Învățarea din exemple prin detecție de similarități*, proces esențialmente inductiv (fără a exclude și procese deductive) prin generalizarea exemplurilor pozitive, generalizare care evită exemplele negative.

- *Învățarea prin observare și descoperire* (spre exemplu a unor regularități în structurări de date).

Gh. Tecuci înclină către o *îmbinare de metode de învățare*. Fără îndoială câteva lucruri credem că se susțin pentru procesul uman de învățare:

- **Necesitatea unei varietăți de metode, și nu o monometodă, lucru deosebit de important** când ar putea apare tendința de a ne baza, în viitor, mai mult pe tehnologie în procesul educațional.
- *Obținerea unui sistem de cunoștințe* sub forma unui model intern de bază (unor modele interne) la care să se poată racorda ușor cunoștințe de detaliu provenite din exterior eventual prin metode informatice.
- *Obținerea sensurilor cunoașterii*, a sensului 'fizic', al intuiției lucrurilor și chiar a unui răspuns creativ în procesul învățării, lucru de care automatele nu sunt capabile, adică a pune umanul în starea lui firească.
- *O deschidere firească spre creativitate și creație*, spre inovare în vederea rezolvării de probleme care nu sunt structurate după tipul sistemului de cunoștințe existent în modelul intern disponibil la un moment dat.'

Odată cu apariția e-învățării se deschid perspective noi și pentru studierea experimentală a procesului de învățare și confruntarea acestui tip important de proces cognitiv cu teoriile științei cognitive care se vor baza pe progresele pe care le va realiza, ceea ce numim, știința integrativă. Considerații privind procesul de învățare și sisteme de educație, inclusiv prin folosirea metodelor inteligenței artificiale sunt prezentate într-un grup de trei lucrări recente ale unor cercetători științifici de la Centrul pentru Cercetări Avansate în Învățarea Automată, Prelucrarea Limbajului Natural și Modelare Conceptuală și Institutul de Psihologie 'Mihai Ralea' al Academiei Române [14], [15], [16].

National Research Council de pe lângă Academia Națională de Științe din SUA a prezentat, în februarie 2002, un raport [17] privind cercetarea științifică a educației în care despre studiul științific al procesului de învățare se arată:

'Much of the controversy about education research relates to its perceived lack of quality. [...] Is scientific education research the same as research in social and behavioral science generally or the same as research in the physical sciences? [...] A key finding of this NRC committee is that at a fundamental level, scientific inquiry in education is no different from scientific inquiry in other fields and disciplines. A set of basic principles is common to all scientific endeavors: these principles include concepts like linking empirical data to theoretical models, using appropriate methods, applying rigorous reasoning, striving toward generalization.'

Considerațiile de mai înainte, inclusiv ale informaticienilor români, arată cât de deschis este în continuare câmpul cercetărilor privind procesul de învățare, în special al omului.

II.4 Învățământul electronic (e-learning)

E-learning este un vector funcțional al societății cunoașterii. Învățarea electronică înseamnă a învăța folosind mijloace electronice, ceea ce se poate face în mai multe moduri:

- Individual - folosind resursele existente pe Internet și CD-uri.
- Instituționalizat - în școli și universități sau organizat în întreprinderi sau de către fundații. Cursurile prin televiziune vor ceda locul cursurilor prin Internet, dar acest procedeu se va desfășura sub supravegherea și îndrumarea cadrelor didactice calificate.
- În cursul activității practice, din orice domeniu, care se va desfășura și într-un mediu informațional și de cunoaștere.

Cei care învață sunt persoane, dar și agenți inteligenți. În viitorul imediat, agenții inteligenți vor deveni nu numai studenți, ci și profesori, dar rolul lor cel mai promițător este acela de colaborator cu persoane. Învățarea implicând agenții inteligenți va deveni o etapă esențială în societatea cunoașterii, deoarece *în regim de croazieră societatea cunoașterii se va baza în cele mai multe activități pe agenți inteligenți. Inteligența artificială va fi esența tehnologică a societății cunoașterii*. Ea va antrena internetul, nanotehnologiile, dar și vectorii funcționali ai societății cunoașterii [1b]. Inteligența Artificială în primii 20 de ani ai sec. XXI va depăși inteligența omului (numai pentru aspectele structurale, fără intuiție și creativitate).

E-învățământul se găsește astăzi în plină dezvoltare [18], [19], [20], [21], [22]. Din experiența relatată în asemenea studii rezultă:

- Studenții găsesc, chiar în cazul lipsei unei interacțiuni față în față între profesor și student, că descărcarea notelor de curs prin Internet, corespondență prin e-mail cu profesori și instructori, examene prin răspunsuri date pe calculator, acasă sau la școală, acest e-învățământ este foarte agreabil. Iar performanțele studenților și elevilor sunt similare (evaluare pentru anul 2000) cu cele ale învățământului în clase de elevi și studenți..
- Corporațiile industriale recurg masiv la e-educație, iar această tendință nu mai poate fi ignorată. Unele corporații au lansat e-universități pentru personalul propriu, de. Ex. Dell Computer Corp. și Sun Microsystems.
- Universitățile au început să introducă nu un e-învățământ complet, ci constituirea treptată a acestuia prin unele e-cursuri. Spre exemplu University of California, Berkeley, în domeniul științei și tehnologiei informației a început (anii 1999-2000) cu patru e-cursuri: sisteme informatice, telecomunicații digitale, e-comerț, sisteme informaționale geografice.
- O serie de firme și-au dedicat activitatea sau o parte din activitate producerii unor 'e-learning software packages'. Se constituie un segment al pieții software specializat în e-learning. (Astfel se și explică prezența firmelor

SOFTWIN și SIVECO la acest simpozion devenite principalele firme românești de software educational. Dar asemenea pachete e-software pentru învățământ sunt de așteptat și din partea Programului e-școală al Ministerului Educației și Cercetării care urmărește o reformă educațională în România.

- Nu se constată deosebiri între rezultatele învățării on-line și învățarea într-un campus universitar sau o școală. Învățarea electronică cere mai multă disciplină și maturitate decât învățarea convențională [18].
- Pentru experimente de laborator și pentru viață socială este nevoie totuși de perioade de lucru în instituțiile de învățământ.
- Odată cu creșterea utilizării metodelor de e-învățământ, construcția de clădiri pentru învățământ se va diminua. În schimb apar cheltuieli pentru noua infrastructură a e-învățământului.
- Modul asincron de acces la cursuri permite e-educația în orice moment și în orice loc.
- E-învățământul încurajează studenții să-și asume o mai mare responsabilitate pentru definirea și organizarea a ceea ce urmăresc să învețe. Studenții sunt mai bine serviți având un acces asistat electronic on-line la cei mai buni instructori decât un contact față în față cu instructori mediocri [19]. În orice caz, nu se neagă rolul instructorilor.
- Discipline ca filosofia și istoria presupun discuții, iar discipline tehnice presupun proiecte. În aceste cazuri trebuie încă să se găsească soluții mixte de învățământ clasic și electronic.
- E-învățământul oferă cele mai bune perspective pentru învățarea în întreaga viață (învățarea continuă).
- 'Educația bazată pe Internet resuscită probleme fundamentale ale educației care sunt importante pentru conceperea activităților educaționale'. [19]
- Gradul în care instructorii vii pot fi înlocuiți cu agenți inteligenți specializați nu este încă clarificat.
- În mod diferit se pun problemele e-învățământului în școli elementare și licee în raport cu învățământul superior. Pentru școli și chiar licee, într-o primă etapă se dezvoltă clase conectate la Internet, cu calculatoare personale, dotate cu e-books, e-learning books, discuri compacte și acces la rețele specializate, eventual servere de clasă sau școală.
- Școlile, ca și companiile, ca și guvernul, trebuie să se regândească în lumina noilor tehnologii ale societății cunoașterii.
- Se preconizează și se experimentează atât pentru școli, cât și pentru alte forme de învățământ, utilizarea Internetului prin comunicații fără fir (wireless Internet) care oferă posibilități și opțiuni noi.

Acestea sunt principalele considerații și constatări la începutul anului 2002. Valabilitatea unora dintre ele se va confirma, alte constatări vor fi, poate, infirmate, dar vor apare cu siguranță multe alte aspecte noi.

II.5 Viața intelectuală

În timp sunt prevăzute multe schimbări datorită învățământului electronic [23]. În primul rând, apariția unor colegii și universități nelocalizate, extinse uneori la scară globală. Siturile acestora pot fi mari sau mici, structurarea socială având loc sub forma unor comunități (villages) având facilități comune pentru cercetare, proiecte de grup, dar și pentru activități comunitare culturale, sportive etc. O persoană admisă într-o asemenea universitate îi va rămâne atașată pentru toată viața, deoarece educația se va extinde pe întreaga viață prin perioade discrete (adică necontinue) și intensive de învățare. Viața intelectuală se va schimba foarte mult, reflectând modificările în cunoaștere:

'An epistemic change is the abandonment of the notion that any single human mind can bear any significant fraction of what is knowable...Even the renaissance notion of an 'educated person' has been discarded - there is no longer a canonical body of basic knowledge that defines this notion' [23].

Agenții inteligenți de căutare a informației, bibliotecile electronice, vizualizarea informației, pătrunderea în medii virtuale, toate acestea vor constitui un software care devine literatură [23]. 'Tehnologia va fi văzută ca cea mai bogată dezvoltare în cultura umană' [23]. Rădăcinile intelectuale se vor baza pe inginerie și tehnologie: Difuzia umanităților în tehnologie și invers, vor duce la o reorganizare radicală a disciplinelor intelectuale [23].

II.6 Perspective

Cum vor evolua lucrurile în viitor? Ray Kurzweil [24] face următoarele previziuni privind educația pentru anii 2009, 2019 și 2029:

Pentru anul 2009 [24, p. 191-192]:

'...most effective learning from computers taking place in the home. [...] The profound importance of the computer as a knowledge tool is widely recognized. Computers play a central role in all facets of education, as they do in other spheres of life. The majority of reading is done on displays, although the 'installed base' of paper documents is still formidable. The generation of paper documents is dwindling, however, as the books and other papers of largely twentieth century vintage are being rapidly scanned and stored. Documents circa 2009 routinely include embedded moving images and sounds. Students of all ages typically have a computer of their own, which a thin tabletlike device weighing under a pound

with a very high resolution display suitably for reading. Students interact with their computers primarily by voice and by pointing with a device that looks like a pencil. Keyboards still exist, but most textual language is created by speaking. [...]

Intelligent courseware has emerged as a common means of learning. [...] The traditional mode of a human teacher instructing a group of children is still prevalent, but schools are increasingly relying on software approaches, leaving human teachers to attend primarily to issues of motivation, psychological well-being, and socialization.'

Pentru anul 2019 [24, p.204]:

'Paper books and documents are rarely used or accessed. Most twentieth-century papers of interest have been scanned and are available through wireless network. Most learning is accomplished using intelligent software-based simulated teachers.[...] The teachers are viewed more as mentors and counselors than as sources of learning and knowledge. Students continue to gather together to exchange ideas and to socialize, although even this gathering is often physically and geographically remote. [...] Most adult human workers spend the majority of their time acquiring new skills and knowledge.'

Pentru anul 2029 [24, p. 221]:

'Human learning is primarily accomplished using virtual teachers and is enhanced by the widely available neural implants. The implants improve memory and perception, but it is not possible to download knowledge directly. Although enhanced through virtual experiences, intelligent interactive instruction, and neural implants, learning still requires time-consuming human experience and study. This activity comprises the primary focus of the human species. Automated agents are learning on their own without human spoon-feeding of information and knowledge. Computers have read all available human and machine generated-literature and multimedia material ...Significantly new knowledge is created by machines with little or no human intervention. Unlike humans, machines easily share knowledge structures with one another.'

Dacă în societatea cunoașterii previziunile de mai înainte se bazează pe o continuare a științei structurale, ce se va întâmpla dacă știința, cu bazele ei noi, integrative, va conduce și la apariția inteligenței artificiale conștiente, adică a conștiinței artificiale? Acest lucru nu se va întâmpla probabil în primii 30 de ani ai acestui secol, dar dacă se va întâmpla cum vom privi și acționa în activitatea educațională?

II.7 Încheiere. Propuneri

Roger Bohn definește, într-un mod specific pentru societatea cunoașterii, **învățarea drept evoluția cunoașterii în timp** [25].

Studiul procesului de învățare scoate în relief importanța științei cognitive și a învățării ca proces cognitiv fundamental. Această știință trebuie nu numai cunoscută, atât cât este ea astăzi, ci mai ales dezvoltată de către psihologi, neu-robiologi, sociologi și specialiști în inteligența artificială.

Este necesară o direcție de cercetare bine susținută pentru a stimula contribuții românești în acest domeniu. Am propus și propun în continuare ca în cadrul programului INFOSOC (Programul național de cercetare-dezvoltare pentru societatea informațională) să se stimuleze cercetări în domeniul științei cognitive care să contribuie la depășirea limitelor actuale ale acestui domeniu.

Este, de asemenea, necesară o dinamizare nu numai a cercetărilor, dar mai ales a dezvoltărilor și realizărilor concrete în domeniul inteligenței artificiale. Există un sistem românesc, sistemul DISCIPL, creat de acad. Gh. Tecuci [13], [26] la ICI și apoi la George Mason University din SUA. Ar trebui examinat și utilizat și la noi. Ar trebui să cunoaștem ce posibilități și ce potențial avem în domeniul utilizării agenților inteligenți și să existe o coordonare și autocoor-donare a eforturilor. Utilizarea agenților inteligenți pentru toți vectorii societății cunoașterii, inclusiv pentru e-învățământ va deveni determinantă pentru calitatea și eficiența acestei societăți. Recenta propunere pentru transformarea Centrului pentru Cercetări Avansate în Învățarea Automată, Prelucrarea Limbajului Natural și Modelare Conceptuală al Academiei Române într-un Centru de cercetări pentru Inteligența Artificială și Societatea Cunoașterii sprijinită de Directorul general ICI, Doina Banciu și de Ministrul Comunicațiilor și Tehnologiei Informației, Dan Nica, ar putea să satisfacă acestor cerințe actuale și de viitor. Sperăm ca și Academia Română să sprijine această solicitare pentru a putea fi înaintată Guvernului României spre a fi aprobată.

Tot la Academie, Comitetul Român pentru Istoria și Filosofia Științei și Tehnicii va acorda o anumită importanță muzeelor virtuale, nu numai pentru istoria științei și tehnicii, dar și pentru cunoaștere și învățare. Ar trebui realizat un web-site de sinteză a tuturor muzeelor virtuale din lume, inclusiv al web-site-urilor unor mari muzee de mare tradiție și importanță, cu adresele lor pe Internet. Acest web-site ar trebui să fie cunoscut și accesibil tuturor în România.

Apreciez în mod deosebit eforturile care se fac pentru informatizarea învățământului românesc de către Guvernul României, Ministerul Educației și Cercetării, firmele SIVICO și SOFTWIN, ca și de toate instituțiile reprezentate la acest simpozion dedicat învățământului electronic.

Doresc să mulțumesc tuturor celor care au prezentat comunicări la acest simpozion și celor care au participat la organizarea lui.

Referințe bibliografice

- [1] Mihai Drăgănescu, *Cunoașterea și societatea cunoașterii*, comunicare la sesiunea de lansare a programului strategic SI-SC, Academia Română, 10 aprilie 2001; 1b. Mihai Drăgănescu, *Societatea informațională și a cunoașterii. Vectorii societății cunoașterii*, studiu, Academia Română, 7 iulie 2001, publicat pe Internet și în vol. coord. Florin Gh. Filip, *Societatea informațională-Societatea cunoașterii. Concepte, soluții și strategii pentru Români*, Academia Română, 2002, p.43 - 112.
- [2] Gabriela S. Sabău, *Societatea cunoașterii. O perspectivă românească*, Editura economică, București, 2001; Ion Gh. Roșca, Viorel Petrescu, Beniamin Cotigaru, Gabriela Sabău, Vasilica Ciucă, Oscar Hoffman, Wilhelm Kappel, *Cercetarea pentru dezvoltarea în reconstrucția durabilă a economiei din perspectiva societății cunoașterii*, *Economistul*, 4 februarie 2002, nr.270, p. I-III.
- [3] Mihai Drăgănescu, *Societatea cunoașterii*, *Diplomat Club*, 2001, Nr. 6, p1-2; Mihai Drăgănescu, *Knowledge management, a functional vector of the knowledge society*, *Diplomat Club*, Nr. 10-11, 2001, p.4; Mihai Drăgănescu, *Factori noi în viața cultural-științifică-politică globală: terorismul și antiterorismul*, *Diplomat Club*, 2002, Nr.1, p.7.
- [4] Mihai Drăgănescu, *Tendencies of becoming*, *Romanian Review*, 1986, Nr. 11, p.55-59.
- [5] Mihai Drăgănescu, *Spiritualitate, Informație, Materie*, p.23-28, Ed. Academiei R.S.R., 1988.
- [6] coord. Doina Banciu, *Cartea Electronică*, Editura AGER, București, 2001.
- [7] Mihai Drăgănescu, *Societatea cunoașterii și cartea electronică*, în vol. coord. Doina Banciu, *Cartea Electronică*, Editura AGER, București, 2001, p. 26-42.
- [8] Paul Saffo, Institute for the Future, *Electronic books*, <http://www.saffo.org/sflibrary.html>, 1988.
- [9] Dick Brass, Vicepreședinte Microsoft pentru dezvoltare tehnologică, *E-books*, în vol. *Inside/Out, Microsoft- in our own words*, Penguin Books, New York, 2000, p.262-263.
- [10] Mihai Drăgănescu, *Microelectronica și învățământul în domeniul electronicii (I)*, *Forum*, anul XXX, noiembrie 1988, p. 36-48.
- [11] Mihai Drăgănescu, *Știința cognitivă, știința structurală sau știință integrativă?* Comunicare la sesiunea științifică de toamnă AOS-R, București, 9 noiembrie 2001, E-PREPRINT, MSReader format, november 2001.
- [12] Menas Kafatos, Mihai Drăgănescu, *Preliminaries to the Philosophy of Integrative Science*, MSReader e-book, Editura ICI, București, 2001, ISBN 973-10-02510-X.

-
- [13] Gheorghe Tecuci, *Mediu de dezvoltare a sistemelor expert instruibile pentru proiectarea asistată de calculator*, Teză de doctorat, Institutul Politehnic, București, 1988.
- [14] Ștefan Trăușan-Matu, *Achiziția, gestiunea, partajarea și prelucrarea cunoștințelor pe web: elemente esențiale în societatea cunoașterii*, în vol. coord. Florin Gh. Filip, *Societatea informațională-Societatea cunoașterii. Concepte, soluții și strategii pentru Români*, Academia Română, 2002, p.195-207.
- [15] Cristina V. Niculescu, *Noi tipuri de sisteme educaționale pentru SI-SC*, în vol. coord. Florin Gh. Filip, *Societatea informațională-Societatea cunoașterii. Concepte, soluții și strategii pentru Români*, Academia Română, 2002, p.209-223.
- [16] Gheorghe Iosif, Ana Maria Marhan, Ion Juvină, *Strategii de creștere a utilizabilității și de dezvoltare a competențelor de bază ale populației României pentru utilizarea tehnologiei informației*, în vol. coord. Florin Gh. Filip, *Societatea informațională-Societatea cunoașterii. Concepte, soluții și strategii pentru Români*, Academia Română, 2002, p. 225-235.
- [17] Lisa Towne, Study Director, Committee on Scientific Principles in Education Research National Research Council/National Academy of Sciences, *Statement before the Subcommittee on Education Reform Committee on Education and the Workforce United States House of Representatives*, February 28, 2002.
- [18] Robert Ubell, *Engineers turn to e-learning*, IEEE Spectrum, October 2000, p.59-63.
- [19] Peter Wiesner, *Distance Education: Rebottling or a New Brew?* Proceedings of the IEEE, July 2000, p.1124- 1130.
- [20] Ralph B. Ginsberg, Kenneth R. Foster, *The Wired Classroom*, IEEE Spectrum, August 1998, p.44- 51.
- [21] Paul G. Shotsberger, Ron Vetter, *Teaching and Learning in the Wireless Classroom*, Computer, march 2001, p.110-111.
- [22] <http://www.microsoft.com-education>
- [23] Edward A. Lee, David G. Messerschmitt, *A higher education in the year 2049*, Proceedings I.E.E.E., September 1999, p.1685 - 1691.
- [24] Ray Kurzweil, *The Age of Spiritual Machines*, Penguin Books, New York, 1999.
- [25] Roger E. Bohn, *Measuring and Managing Techological Knowledge*, p.295-314, în vol. Eds. Dale Neef a.o., *The Economic Impact of Knowledge*, Butterworth-Heinemann, Boston, 1998.
- [26] Gh. Tecuci, *Building Intelligent Agents*, Academic Press, San Diego, 1998.

III. CULTURA ȘI SOCIETATEA CUNOAȘTERII

Societatea Cunoașterii

Am prefigurat că va sosi un moment al societății cunoașterii (chiar cu această sintagmă, Mihai Drăgănescu, 1976, 1986), dar abia în ultimul deceniu al secolului XX conceptul s-a impus în SUA datorită lucrărilor sociologului Peter Drucker și ale altora, în ultimii 4-5 ani societatea cunoașterii devenind recunoscută ca o etapă nouă a erei informației, respectiv a societății informaționale. Academia Română a lansat acest concept în România în anul 2001 ca urmare a poziției și comunicării Mihai Drăgănescu, *Cunoașterea și Societatea cunoașterii*, la sesiunea de lansare a programului SI-SC, Academia Română, 10 aprilie 2001 și a elaborării studiului Mihai Drăgănescu, *Societatea Informațională și a Cunoașterii. Vectorii Societății Cunoașterii*, Academia Română, București, 9 iulie 2001, publicat apoi în vol. coord. Florin Gh. Filip, *Societatea informațională-Societatea cunoașterii. Concepte, soluții și strategii pentru România*, Academia Română, 2002, p.43 - 112.

Spre deosebire de unele puncte de vedere care privesc numai economicul (economia digitală, piața internet) societatea cunoașterii **nu este numai economia bazată pe cunoaștere**. Aceasta este foarte importantă, decisivă, esențială și cuprinde utilizarea și managementul cunoașterii existente sub forma cunoașterii tehnologice și organizaționale, producerea de cunoaștere tehnologică nouă prin inovare, **o nouă economie** în care procesul de **inovare** este determinant, *în care bunurile intangibile devin mai importante decât cele tangibile*.

Societatea cunoașterii **reprezintă mult mai mult** deoarece asigură o diseminare fără precedent a cunoașterii către toți cetățenii prin mijloace noi, folosind cu prioritate Internetul și cartea electronică și metodele de învățare prin procedee electronice (e-learning), urmărește extinderea și aprofundarea cunoașterii științifice și a adevărului despre existență, *este singurul mod prin care se va asigura o societate sustenabilă din punct de vedere ecologic și va fi o nouă etapă în cultură* (bazată pe cultura cunoașterii care implică toate formele de cunoaștere, inclusiv cunoașterea artistică, literară etc).

În fine, societatea cunoașterii asigură bazele unei viitoare societăți a conștiinței, a adevărului, moralității, creativității și spiritului.

Pentru realizarea societății cunoașterii am definit, în studiul amintit mai înainte, o serie de vectori (tehnologici și funcționali) care ar trebui introduși în acțiune într-o succesiune firească pentru posibilitățile țării noastre.

Categoriile culturii

Dintre lucrările pe care le-am publicat anterior în problemele culturii [1] două se referă la teoria culturii. În *Perspectiva informațională a culturii* (1983) găseam un anumit sprijin pentru o viziune informațională a culturii în teoria semiotică a culturii elaborată de Umberto Eco în *Tratatul său de semiotică generală*. Umberto Eco propunea o ipoteză

radicală prin care întreaga cultură este considerată un fenomen semiotic și o ipoteză moderată prin care orice aspect al culturii este o entitate semantică. Semiotica se referă la semne cu conținut semantic astfel încât cele două ipoteze nu sunt prea deosebite. De aceea, consideram, prin generalizare firească, deoarece semnul și semanticul (de semnificație și de sens) sunt informație, o posibilă perspectivă informațională a culturii. Acest lucru, faptul că **esența culturii este informațională**, chiar dacă ea se manifestă prin comportamente socio-umane, obiecte materiale și informaționale, a devenit tot mai evident. Nu trebuie să surprindă această esență informațională a culturii, astăzi fiind știut că și inteligența și conștiința sunt informație.

În legătură cu perspectiva informațională a culturii poate fi menționat ca precursor al acestei abordări, Ernst Cassirer [2] care considera că expresia culturală a omului și societății este caracterizată de activitatea de creare a simbolurilor (activitatea simbolizatoare) generate de imagini mentale. Pentru Cassirer, simbolul este o cheie pentru înțelegerea naturii omului, iar omul nu trăiește numai într-un univers material, ci mai ales într-unul simbolic [3].

Într-o a doua lucrare [1a], *Cultura și marile tehnologii* (1996) am urmărit linia clasică de definire a culturii ținând însă seamă de obiectele informaționale noi aduse de societatea informațională. **În teoria clasică cultura este definită ca un fenomen social care cuprinde comportamentul socio-uman cu obiectele materiale și informaționale integrate acestui comportament.** Obiectele informaționale au fost introduse în această definiție la sfârșitul secolului XX.

Pare a fi posibilă o încadrare a teoriei culturii într-o viziune categorială (termenul este utilizat în raport cu teoria categoriilor și functorilor din matematică, extinsă recent de la domeniul structural la domeniul structural-fenomeno-logic [4]).

Privind comportamentul socio-uman *cultural* ca o categorie, această categorie este o subcategorie majoră a categoriei comportamentului socio-uman *general*. Ultima mai cuprinde și o subcategorie a comportamentului determinat strict biologic, atât la nivel individual, cât și social. Într-adevăr, pe lângă comportamente individuale strict biologice există și comportamente sociale determinate biologic, puse în evidență, în cazul omului, de Gr.T. Popa [5]. Acesta demonstrează cum creierul vechi (primitiv, reptilian, thalamus-hipotalamus) determină comportamente necontrolate cultural care duc mase de oameni la comportamente sălbatice, iar în cazul societăților mai avansate duc la manifestări de semicivilizație, în care impulsivitatea biologică devine colectivă și sălbatecă.

Cultura

O subcategorie a unei categorii este o categorie. **Categoria cultură** este subcategorie a comportamentului socio-uman general, dar este aceea care deosebește specia umană de toate celelalte specii animale, chiar dacă unele dintre acestea pot avea și rudimente de cultură. Categoria cultură reprezintă *comportamentul socio-uman cultural*, spre deosebire de cel biologic, cu tot ceea ce construiește, dar nu se dezvoltă decât datorită, totuși, anumitor proprietăți biologice remarcabile ale omului, în special ale creierului său

care are o mare disponibilitate informațională. De aceea, dacă originea biologică a comportamentului cultural nu poate fi pusă la îndoială, cultura este o construcție care se ridică mult deasupra biologicului, atât cât va putea față de limitele biologice ale omului la un moment dat în istorie.

Poate că alături de cele două subcategorii menționate mai înainte ar trebui să mai adăugăm comportamentului socio-uman încă una, aceea a spiritualității (comportamentul spiritual), pe care nu o tratăm în această lucrare. A privi spiritualitatea ca o a treia subcategorie a comportamentului socio-uman general este o chestiune care trebuie aprofundată, având în vedere că mulți oameni de cultură consideră spiritualitatea a fi un comportament numai cultural. Ținând seama de cercetările și studiile de filozofie a științei din ultimii 15 ani privind mintea și conștiința, vom considera, până la argumente contrarii convingătoare, spiritualitatea ca fiind o subcategorie separată și nu una înglobată (total) în cultură.

Schematic, vom rezuma cele de mai înainte, astfel:

CATEGORIA COMPORTAMENTUL UI SOCIO-UMAN	Subcategoria comportamentului strict biologic	Notă: există o comportament social determinat strict biologic
	Subcategoria comportamentului cultural.	CULTURA
	Subcategoria comportamentului spiritual.	SPIRITUALITATEA

Sferele mari ale culturii

Pornind de la definiția din [1b] și diferența pe care o face UNESCO între cultura intangibilă și cultura tangibilă, marile sfere (categorii) ale culturii pot fi considerate următoarele:

- I. **Cultura intangibilă.** 'Moștenirea intangibilă poate fi definită ca îmbrățișând toate formele de cultură tradițională și populară sau cultura folk, adică producțiile colective originare de o comunitate dată și bazate pe tradiție. Aceste creații sunt transmise oral sau prin gesturi și sunt modificate într-o perioadă de timp printr-un proces de re-creare colectivă. Ele includ tradițiile orale, obiceiurile, limbajele, muzica, dansul, ritualurile, festivitățile, medicina tradițională și farmacopeea, artele culinare și tot felul de îndemânări speciale legate de aspectele materiale ale culturii, cum sunt uneltele și habitatul [6]'. Fără îndoială, noțiunea de cultură intangibilă a fost introdusă sub influența noțiunii de valoare intangibilă din economie care a căpătat o mare importanță pentru societatea cunoașterii (economia bazată pe cunoaștere). Se mai adaugă aici valori, credințe, cunoaștere tacită.

II. **Cultura umanistă.** Am preluat în acest studiu denumirea tradițională. Cultura umanistă cuprinde limbajele naturale, literatura, arta, istoria, filosofia, sportul. Cultura umanistă este o cultură tangibilă, ca și știința și tehnologia.

III. **Cultura științifică: Știința, tehnologia și cunoașterea.** Această categorie a culturii conține două subcategorii:

III.a Știința, cunoașterea științifică și tehnologică, cunoașterea tehnologică pentru fabricația de produse, dar și pentru utilizarea acestora, precum și cunoașterea organizațională și economică, chiar dacă unele obiecte ale cunoașterii sunt tacite sau fac parte și din cultura intangibilă. În categoria mare a culturii, anumite obiecte pot aparține la două sau mai multe subcategorii, acestea nu sunt neapărat disjuncte.

III.b Uneltele fizice și informaționale, obiectele fizice și informaționale produse sau fabricate, utilizarea lor, instituțiile și organizațiile, care sunt consecințe, în cea mai mare măsură, a cunoașterii științifice, tehnologice, economice și organizaționale, poate chiar și a culturii intangibile.

Nu numai că unele obiecte culturale pot face parte din mai multe subcategorii ale culturii, dar vor exista și zone de interferență între obiecte ale acestor subcategorii. De exemplu, filosofia științei, care este un obiect al filosofiei, nu se poate dezvolta decât în strânsă legătură cu știința. În teoria categoriilor asemenea legături se numesc morfisme (morphisms sau maps, în limba engleză). Mai mult, pe lângă legăturile dintre obiectele subcategoriilor culturii, din orice sferă a culturii ar proveni, există relații între aceste sfere în totalitatea lor. Acestea se numesc functori. Cei mai importanți functori sunt aceia dintre categoria II și categoria III de mai sus. Acești functori,

F1 : Categoria III (Cultura științifică) → Categoria II (Cultura umanistă)

F2 : Categoria II (Cultura umanistă) → Categoria III (Cultura științifică)

reprezintă relația și influența reciprocă dintre, în esență, cultura umanistă și știință (cultura științifică). Importanța lor pentru societate și om nu poate fi subestimată.

Care este mai importantă dintre cele două categorii? Ambele sunt importante, dar motorul dezvoltării provine din sânul categoriei III. Acest lucru a devenit tot mai evident odată cu formularea conceptelor societății cunoașterii [7].

Este adevărat că o altă resursă importantă este viața spirituală, ea având și componenta de creație implicând puternic atât cultura umanistă, cât și cultura științifică.

Odata cu era informației vor apare desigur multe elemente noi ale culturii datorită tehnologiei informației, cărții și documentelor electronice, internetului, tehnologiilor vorbirii, tehnologiilor bioelectronice și bioinformaticice, inteligenței artificiale și agenților inteligenți informatici, mediului ambiant inteligent, apariției conștiinței artificiale. Vor apare schimbări în viața intelectuală, socială și politică.

Ce se va mai petrece în cultură ?

În secolul XXI sunt posibile câteva evenimente majore care vor schimba viața omenirii:

- Prăbușirea ecologică a societății și a speciei umane, datorită deteriorării grave a mediului înconjurător, ceea ce s-ar putea întâmpla la mijlocul sec. XXI (să spunem, anul 2050) dacă nu se trece din timp, respectiv de pe acum, la efortul de asigurare a unor societăți sustenabile. Salvarea este posibilă chiar cu cunoașterea științifică și tehnologică de astăzi dacă se trece la un management adecvat al cunoașterii [7] și la noi concepte economice adaptate sustenabilității. În această problemă au apărut și alte noi perspective care vor rezulta dintr-o serie de evenimente descrise în continuare.
- Dezvoltarea inteligenței artificiale până la depășirea inteligenței umane, ceea ce se va putea petrece între anii 2019-2035 sau chiar mai devreme [8], [9], [10], [11], [12].
- Apariția conștiinței artificiale, tot în cursul sec. XXI, după ce inteligența artificială va depăși inteligența umană, dar fără a putea preciza perioada.

Aceste două ultime evenimente presupun apariția unor noi specii inteligente, dar și noi specii conștiente, unele nebiologice (roboți umanoizi în topul unor *specii* de roboți mai puțin inteligenți care simulează animale (insecte, pisici, câini) și roboți construiți pentru anumite funcțiuni care să înlocuiască omul [8][13].

Speciile de roboți umanoizi inteligenți și de agenți software inteligenți, ambele egal de inteligente sau mai inteligente decât omul sunt uneori numite *robo sapiens* [13]. Într-o primă etapă, aceste specii nu vor avea conștiință, astfel cum are omul, datorită faptului că au numai o organizare structurală și nu una structural-fenomenologică [14]. Dar aceste specii vor interacționa puternic cu omul și societatea și se pune întrebare în ce măsură ele vor fi și artefacte culturale, nu numai prin faptul că fac parte din cultura omului, ci și prin participarea lor activă la cultură. Vor dezvolta cultura lor (într-o anumită măsură, da) sau vor intra în jocul marii culturi, participând la *cultura totală devenită din fenomen socio-uman, unul socio-uman- inteligentă/conștiință artificială*

Întrucât *robo sapiens* va avea cunoaștere și va participa la dezvoltarea științei și tehnologiei, chiar la dezvoltarea sa ca obiect tehnologic, el va participa cu siguranță la cultura științifică, poate chiar la anumite forme de cultură umanistă sau numai robotică. El poate fi implicat, prin cunoașterea culturii umaniste, să participe ceva mai pronunțat la această cultură. Când va trece de la inteligență la conștiință, o asemenea activitate ar putea fi mult mai pronunțată.

Probabil, între *homo sapiens* și *robo sapiens* vor exista relații de competiție și cooperare, dar acestea se vor dezvolta într-o societate comună, cel puțin până la o segregare care nu ar fi de dorit, în care spiritualitatea și creativitatea lui *homo sapiens* îi va conferi acestuia din urmă poziții inabordabile lui *robo sapiens*. Din momentul în care vor apare

specii de *robo sapiens-conștient*, lucrurile se vor schimba din nou, cu efecte poate și mai dramatice pentru om și societate. Încerc să mă conving că ideile unei societăți a conștiinței ar putea fi benefice pentru un asemenea viitor care probabil nu va putea fi prohibit. Probabil, înspre un asemenea viitor și într-un asemenea viitor să fir rezolvată și sustenabilitatea unei societăți a conștiinței.

Este interesant de reluat aici câteva previziuni ale lui Kurzweil [10] privind starea societății în anii 2019 și 2029.

Pentru anul 2019, în domeniul afacerilor și al economiei, prevede tranzacții care în majoritate vor folosi persoane simulate, oamenii de afaceri vor avea asistenți software care vor conduce tranzacțiile în numele lor. Locuințele vor dispune de roboți de întreținere. Cu aceste artefacte comunicarea se va face prin voce, deoarece vor dispune de o tehnologie a limbajului natural și a vocii de foarte înaltă calitate. Oamenii vor avea relații cu persoane automate inteligente în calitatea acestora de profesori, îngrijitori medicali, persoane de companie etc. Aceste persoane automate au și calități superioare omului în privința memoriei, dar, afirmă Kurzweil, 'ele nu sunt încă privite ca fiind egale cu oamenii în toată subtilitatea personalității acestora'. Inteligența artificială este însă prezentă și împletită cu toate aspectele societății. Responsabilitatea omului va rămâne totuși pe primul plan și nu a persoanelor (agenților) care îl ajută. Operele de artă se vor realiza prin colaborarea dintre artiști umani și inteligențe artificiale. Principalul pericol în societate îl vor constitui micile grupuri de oameni și inteligențe artificiale folosind comunicații criptate care nu pot fi descifrate. Acestea vor folosi virusuri informatice și agenți de îmbolnăvire obținuți prin bioinginerie. Pe de altă parte descifrarea relațiilor dintre genele genomului uman va permite o medicină utilizând inteligența artificială pentru tratamentul și chiar eradicarea multor boli, inclusiv pentru prelungirea considerabilă a vieții omului natural.

Pentru anul 2029, Kurzweil prognozează: în domeniul comunicațiilor va predomina, ca volum, comunicația dintre oameni și mașini. Populația umană se va plafona la 12 miliarde de persoane reale, cărora li se asigură toate condițiile normale de viață. Populația umană și a inteligențelor artificiale va fi preocupată, în primul rând, pentru crearea de cunoaștere, într-o puzderie de forme. Va fi greu de să fie menționate capacități ale omului care să nu fie preluate de mașini, de fapt o deosebire netă nu mai există între lumea oamenilor și lumea mașinilor. Cognația umană a fost transferată mașinilor și multe mașini au personalitate, îndemânări și baze de cunoaștere preluând și cunoaș-terea umană. Implanturile neurale cognitive bazate pe inteligență artificială vor amplifica funcțiile cognitive ale omului. Kurzweil afirmă: 'A defini ceea ce înseamnă o ființă umană devine o chestiune semnificativă politică și de legislație. Creșterea rapidă a posibilităților mașinilor este controversată, dar nu există nici o rezistență față de ea. Deoarece la început mașinile au fost proiectate pentru a fi supuse controlului uman, ele nu au prezentat o față amenințătoare față populația umană. Oamenii realizează că nu mai este posibilă dezangajarea civilizației devenită om-mașină de dependența de inteligența mașinilor. Crește discuția despre drepturile legale ale mașinilor, în special ale acelor mașini care sunt independente de oameni (care nu sunt introduse într-un creier uman). Cu toate că nu se recunoaște deplin,

prin lege, influența evidentă a mașinilor la toate nivelele de decizie asigură o protecție importantă a mașinilor'.

Kurzweil consideră calități ale mașinilor inteligente, care încă din anul 2029 pot fi persoane de artă în toate domeniile artei ('Mulți dintre artiștii de frunte sunt mașini'). Observăm însă că acest lucru ar presupune o stare de conștiință similară omului și prin manifestarea fenomenelor de qualia. Implicit, Kurzweil consideră că mașini inteligente complexe structurale pot avea asemenea stări și pot chiar participa la discuții filosofice pe baza experienței proprii. Vorbind de experiența subiectivă a unor astfel de mașini, aceasta ar însemna că asemenea mașini să fi trecut pragul de la inteligență la conștiință *numai pe baze structurale* încă din anul 2029. Ceea ce nu credem, în principiu, a fi posibil.

Într-adevăr, previziunile pe care oamenii de știință le fac privind dezvoltarea inteligenței artificiale spre conștiință artificială se bazează pe extrapolări ale științei structurale (complexitatea structurală de la un anumit grad în sus generează conștiință, acest lucru fiind considerat valabil începând cu creierul animalelor). Odată cu creșterea complexității artefactelor creiere electronice sau creierelor software se consideră că atunci când acestea ating complexitatea creierului uman se va produce de la sine conștiința artificială [8], [10], [11]. Uneori, unii dintre cei care susțin un asemenea punct de vedere au îndoieli asupra valabilității lui [12]. În viziunea unei filosofii integrative a științei [15],[16], conștiința nu se poate realiza numai din elemente structurale, fiind nevoie și de elemente fenomenologice [17]. Conștiințele artificiale vor pune probleme foarte mari speciei umane care cred că ar putea fi rezolvate în cadrul unei viitoare societăți a conștiinței. Aceasta va urma atunci societății cunoașterii în cadrul erei informației [18], [19].

Ce va fi cultura în societatea conștiinței, la care vor participa, dacă nu chiar vor predomina conștiințele artificiale? Dacă lucrul cel mai important, în cele din urmă, este continuitatea conștiinței create de omeni, atunci și culturii create de ea trebuie să i se asigure o continuitate.

Aceste considerații arată, dacă mai era nevoie de subliniat, cât de importantă vor fi în sec. XXI, cultura științifică și cultura umanistă, ambele având nevoie de o cultură filosofică adecvată.

Culturi, cultură pozitivă și cultură negativă. Polarizarea culturii în jurul cunoașterii

O cultură poate fi apreciată *pozitiv sau negativ*, în raport cu anumite criterii. Se pierde prea mult din vedere acest lucru. Există astăzi și o cultură a teroriștilor (chiar și o știință a terorismului) o cultură a corupției care ne pune nouă românilor atâtea probleme, o cultură a hoților etc. Desigur, acestea pot fi numite sub-culturi, dar tot culturi sunt. Cultura are multe fațete.

Cultura negativă este o cultură deformată în raport cu criteriile civilizației socio-umane.

În ultimii 12 ani, în societatea românească, pe lângă multe lucruri pozitive, s-au accentuat, din nefericire, și fenomene negative îngrijorătoare: corupție, imoralitate, injustiție. Creșterea imoralității și a injustiției, a influențat până și viața academică din țara noastră. Avem nevoie și de un efort cultural pentru a reduce aceste flageluri din societatea noastră, pe lângă efortul dezvoltării economice.

Un exemplu de cultură pozitivă este arta. A cunoaște arta înseamnă cunoaștere, dar a simți arta, a trăi arta, a avea nevoie de ea, a fi o bucurie interioară, acestea înseamnă cultură umanistă adevărată.

Dar dacă cele de mai sus nu sunt însoțite de comportament civilizată, de civilizație socio-umană, cultura poate fi denaturată (rapturile de opere de artă în scopuri personale sau statale). Natura firească a culturii pozitive este aceea de a susține civilizația socio-umană, spiritualitatea, cunoașterea și conștiința, în cele din urmă societatea cunoașterii și societatea conștiinței.

În privința relației dintre cultura umanistă și cultura științifică, astăzi nu se mai poate vorbi de cultură, cu înțelesul de cultură - în general, dar de fapt cu gândul la cultura umanistă.

Cultura - în general, are o mult prea puternică componentă științifică (inclusiv tehnologică, economică, organizațională, politică) pentru a mai accepta o asemenea simplificare, este adevărat, continuatoarea unei tradiții care astăzi este complet depășită. Cultura, respectiv cultura - în general, este cultura umanistă și cultura științifică, împreună, ultima având, ca și prima, un conținut extrem de bogat.

În spatele confuziei care se menține astăzi atunci când vorbim de cultură se întreține schisma dintre cele două culturi, datorită unor interese de grup. În etapa actuală a societății, cultura umanistă nu-și mai poate erija numele general de cultură, de fapt nu ea, ci slujitorii ei care nu s-au adaptat la vremurile cunoașterii. În societatea cunoașterii, înainte de trecerea la societatea conștiinței, cultura se va concentra în jurul cunoașterii. Iar tehnologia va fi un factor cultural atât de covârșitor încât va reveni poate la pozițiile ei mitologice din antichitate. În [1a] remarcam:

În antichitate, la egipteni, zeul Ptah era privit ca patronul lucrătorilor de metal (metalurgiști și fierari) și al artizanilor. Ptah era însă unul dintre cei mai mari zei, creatorul pământului, părintele zeilor și al începuturilor. Interesant acest "al începuturilor".

La grecii antici, echivalentul lui Ptah era Hefaistos, zeul focului și meșteșugurilor, protectorul artizanilor. El nu mai era o divinitate primordială, dar divinitate, fiul lui Zeus și al Herei, fiind căsătorit cu Afrodita. Se pare că de la egipteni la greci, tehnologia nu mai păstra poziția începuturilor, dar avea totuși un reprezentant divin. La romani, echivalentul lui Hefaistos era Vulcan, considerat zeul focului.

Decăderea poziției tehnologiei în cultură începuse din antichitate. Ea a continuat până în secolul XX când într-adevăr avea să cunoască un reviriment. Astăzi vorbim despre marile tehnologii și chiar despre o filosofie a tehnologiei, de care o serie de gânditori și filosofi au scris lucrări deosebit de interesante: Ernst Kapp, Friedrich Schlegel, José Ortega

y Grasset, Martin Heidegger ș.a. Este adevărat că au apărut și lucrări îndreptate împotriva tehnologiei (L.Mumford, J.Ellul ș.a.), declanșând ceea ce secolul XX s-a numit dilema tehnologică.'

Revirimentul filosofic al tehnologiei în societatea cunoașterii, în secolul XXI, va fi un factor important în gândire, în general. Tehnologia va continua biologicul, culturalul și conștiința.

Ce va face omul? Marea lui înțelepciune va fi aceea de a pregăti în mod corespunzător viitorul [19]. Din ce în ce mai mult, gândirea filosofică va avea un rol hotărâtor în știință, politică, viața socială.

Există și vor exista culturi ale profesiilor, ale domeniilor cunoașterii, ale națiunilor, etniilor, grupurilor, ale comunităților constituite pe Internet, ale instituțiilor și localităților virtuale, ale mașinilor inteligente etc. Lumea devine tot mai pluriculturală. Probabil aceasta este trăsătura cea mai importantă a postmodernității [20].

Momentul actual ar trebui să fie acela al tendinței spre cunoaștere și cultură (cu înțelesul ei total) pentru întreaga populație a omenirii, fiecare zonă locală, geografică sau virtuală, trebuind să fie preocupată activ de realizarea concretă a acestei tendințe.

Referințe bibliografice

[1] Mihai Drăgănescu: *Lucrări despre cultură*:

- I. Mihai Drăgănescu, **Cultura și marile tehnologii**, conferință, Universitatea Populară de Vară "Nicolae Iorga", Vălenii de Munte -30 august 1996.
- II. Mihai Drăgănescu, **Perspectiva informațională a culturii**, Contemporanul, 27 mai 1983.
- III. Mihai Drăgănescu, **Dimensiunile europene ale culturii române**, expunere, Vălenii de Munte, 1992, publicată în *Academica*, 1992.
- IV. Mihai Drăgănescu, **Arta și societatea**, cuvânt, Ploiești, 4 noiembrie 1991, publicat în *Academica*, 1991.
- V. Mihai Drăgănescu, **Criterii transpolitice și transnaționale în cultură**, 18 mai 1997, *Caiete Critice*, 1997, nr.3-4, p.145-147.
- VI. Mihai Drăgănescu, **Spirit enciclopedic și enciclopedism**, conferință, Vălenii de Munte, 22 august 1993 (publicată în *Academica* și în volumul autorului, *Cariatidele gândului*, Ed. Academiei Române, 1996, p. 163-168).

[2] Ernst Cassirer, *Substanzbegriff und Funktionbegriff*, 1910; *Die Philosophie der Symbolischen Formen*, 1923-1929 (3 vol).

[3] Oltea Mișcol, Elena Gheorghe, **Repere istorice în filosofia culturii**, *Revista de filosofie*, XLVII, Nr. 5-6, 2000, p.449-459.

-
- [4] Mihai Drăgănescu, *Categories and functors for the Structural Phenomenological Modeling*, Proceedings of the Romanian Academy, Series A, Vol.1, No.2, 2000, p.111-115.
- [5] Grigore T. Popa, *Reforma spiritului*, volum în editare, conținând lucrări ale acestui autor prezentate și publicate la Academia Română în anii 1940 (vezi și prefața: Mihai Drăgănescu, *O gândire asupra conștiinței, moralității și societății*).
- [6] UNESCO, definiția culturii intangibile, web-site UNESCO.
- [7] Mihai Drăgănescu, *Societatea Informațională și a Cunoașterii. Vectorii Societății Cunoașterii*, Academia Română, București, 9 iulie 2001, publicat în vol. coord. Florin Gh. Filip, *Societatea informațională-Societatea cunoașterii. Concepte, soluții și strategii pentru România*, Academia Română, 2002, p.43 - 112).
- [8] Moravec H., *Rise of the Robots*, Scientific American, December 1999, p. 86-93.
- [9] Moravec H., *Robot. Mere Machines to Transcendent Mind*, Oxford University Press, Oxford, 1999.
- [10] Kurzweil R., *The Age of Spiritual Machines*, Penguin Books, 2000.
- [11] Broderick D., *The Spike*, New York, 2002, paperback.
- [12] Buttazzo G., *Artificial Consciousness. Utopia or Reral Possibility?* Computer (IEEE), July 2001, p.24-30.
- [13] Interviews of Menzel P. and D'Aluisio F., *Robo Sapiens. Evolution of a new species*, MIT Press, Cambridge, Massachusetts, 2002.
- [14] Drăgănescu M., *Din lucrările despre minte și conștiință*:
- A. Mihai Drăgănescu, *The Interdisciplinary Science of Consciousness*, Noetic Journal, Vol.3, No.1, Jan.2000, p. 37-46; republicat în eds. Richard L. Amoroso et al, *Science and the Primacy of Consciousness, Intimation of a 21st Century Revolution*, Chapter 5, pp. 46-59, Orinda: The Noetic Press, 2000.
 - B. Mihai Drăgănescu, *Theories of Brain, Mind and Consciousness: Still Great Divergences*, Noetic Journal, vol.3, No. 2, Apr. 2000, p.125-139.
 - C. Mihai Drăgănescu, *The Brain as an Information Processor*, NOESIS, XXV, 2000, p. 9-20.
 - D. Mihai Drăgănescu, *On the Structural-Phenomenological Theories of Consciousness*, NOETIC JOURNAL, Vol.1, No.1, June 1997.
 - E. Mihai Drăgănescu, *Continuities and Discontinuities in the realms of life and mind*, Revue Roumaine de Philosophie, Tome 41,1997, Nos 1-2, p.3-9.
 - F. Mihai Drăgănescu, *De la filosofia la știința mentalului*, Revista română de filosofie, XLIV, Nr.5, sep-oct 1997, p. 457-464.
 - G. Mihai Drăgănescu, *Procesarea mentală a informației*, Memoriile Sect. St. ale Acad. Române, SERIA IV, Tom. XX, 1997, p.263-284.
- [15] Kafatos M., Draganescu M., *Preliminaries to the Philosophy of Integrative Science*, E-book (Microsoft Reader), ISBN 973-10-02510-X, Editura ICI, Bucharest, 2001.
- [16] Draganescu M., Kafatos M., *Generalized Foundational Principles in the Philosophy of Science*, paper presented at the Conference on "Consciousness in Science and

- Philosophy" in Charleston, Illinois, 6-7 Nov 1998, published in The Noetic Journal, Vol.2, No.4, Oct. 1999, p. 341-350, republished in the vol. *Science and the Primacy of Consciousness, Intimation of a 21st Century Revolution*, Richard L. Amoroso and others (eds), Orinda: The Noetic Press, 2000, Chapter 9, pp. 86-98.
- [17] Mihai Drăgănescu, ***Advancement in Neural Engineering and Neuroelectronics Put Forward Artificial consciousness***, Communication at the INGIMED II Conference, Bucharest, Dec. 13, 2001; E-PREPRINT, MSReader Format, 2002.
- [18] Mihai Drăgănescu, ***Conștiința, frontieră a științei, frontieră a omenirii***, Revista de Filosofie, XLVII, nr. 1-2, 2000, p.15-22.
- [19] Mihai Drăgănescu, ***Societatea conștiinței, o viitoare etapă a erei informației. Vectorii societății conștiinței***, studiu, Academia Română, în pregătire.
- [20] După Alain Fienckielkrant, apud [3], p.458-459.

Între lingvistica matematică și cea computațională

Acad. Solomon MARCUS
Secția de Științe Matematice a Academiei Române
solomon.marcus@imar.ro

Mă simt obligat să reacționez la un anumit mod de prezentare a evoluției ideilor, în cea de a doua jumătate a secolului al XX-lea, în articolul [1] al d-lui Dan Tufiș (de aici mai departe DT), membru corespondent al Academiei Române. Precizez de la început că nu contest interesul și utilitatea direcției de preocupări prezentate în [1]; am în vedere numai modul în care această direcție este pusă în relație cu alte cercetări dedicate limbajului.

Cităm din [1: 133]:

“Desprinzându-se din lingvistica formală, “lingvistica matematică” a încercat dezvoltarea unor modele matematice de reprezentare a limbajelor naturale sau formale (în general al aspectului lor sintactic, gramatical), căutând soluții abstracte de modelare generativă de tip universal a ceea ce se presupunea (la nivelul cunoașterii științifice a anilor 1960) a fi facultatea limbajului”.

Nu știu ce înțelege DT prin “lingvistica formală”, o sintagmă nu prea folosită în perioada de emergență a lingvisticii matematice; există lingvistica structurală (altceva decât ceea ce ar putea fi lingvistica formală, adică bazată pe formalizare în sensul logicii matematice moderne), care desigur a constituit una din sursele lingvisticii matematice (de aici mai departe LM), așa cum i se pot indica și alte surse (biologice, logice, matematice, psihologice etc.), dar factorul determinant în nașterea LM, în a doua jumătate a anilor '50, a fost dezvoltarea calculatoarelor electronice și, împreună cu ea, a primelor preocupări sistematice de LC (prescurtare a lingvisticii computaționale), numite atunci traducere automată, documentare automată, prelucrarea automată a limbajului, cu diverse variante ale lor în engleză (de exemplu, “machine translation”), franceză, rusă, germană, italiană etc. Din aceste preocupări s-au inspirat primele modele care au constituit noua disciplină a LM.

Vorbesc despre lucruri trăite. Punctul meu de plecare s-a aflat în lucrările unor Kulagina și Melciuk, puternic implicați în studiile de traducere automată rusă-franceză, Yves Lecerf, implicat în problemele de documentare automată, D. G. Hays, implicat în traducerea automată din rusă în engleză și reciproc, B. Vauquois, cu preocupări de informatică lingvistică la Grenoble. De la ei, ca și de la alți autori similari, am preluat în bună măsură ștabela pe care am căutat s-o duc mai departe. Ceea ce afirm despre mine este valabil pentru cei mai mulți cercetători din domeniul LM din anii 1950 și 1960, cum ar fi

Maurice Gross, Masami Ito, A. Trybulec și mulți alții. Dubioasă mi se pare sintagma “soluții abstracte”, probabil efectul unui obicei binecunoscut de a diaboliza abstractul.

În ceea ce privește sintagma “lingvistică formală”, ea a căpătat o anumită utilizare în anii târzii 1960 și în anii următori, iar personal am folosit-o în unele lucrări, după cum se va vedea imediat, considerând-o oarecum echivalentă cu LM; dar chiar dacă nu acceptăm această echivalență, nu putem eluda faptul că lingvistica formală se află în imediata vecinătate a LM. DT pretinde ca LM “a încercat”, sugerând astfel că ea a eșuat în tentativă de modelare a limbajului natural. Ceea ce este deocamdată numai o sugestie devine, după cum se va vedea, o certitudine pentru DT.

Într-adevăr, iată ce scrie mai departe DT ([1]: 133):

“Curând metodele lingvisticii matematice și-au atins limitele drept care, în anul 1966, la propunerea lui David Hays, domeniul de cercetare al limbajelor naturale, din perspectiva utilizării acestora în interacțiunea cu calculatoarele electronice, este individualizat sub numele de lingvistica computațională”.

Chestiunea cu atingerea limitelor ține de domeniul umorului involuntar și trecem peste ea, dar nu ne miră, după ce am văzut la ce se reduce LM pentru DT. Nu mi-am imaginat niciodată ca între LM și LC ar putea avea loc o competiție, prima definindu-se prin metoda (căci ce altceva este LM decât studiul limbajului cu ajutorul matematicii ?) iar a doua prin obiectivul pe care și-l propune. LM nu poate ignora problematica LC iar LC nu-și poate realiza proiectele fără LM. Probabil însă că DT lucrează cu o definiție specială a LM, pe care am dori s-o aflăm. Modul simplificator în care DT se referă la generativismul lingvistic, într-o logică binară care eludează faptul că în materie de modelare se lucrează cu grade de adecvare și relevanță, este însă simptomatic pentru viziunea sa limitativă în problema în discuție.

Crede DT că gramaticile lui Joshi, atât de importante în LC, puteau fi concepute fără să fi fost precedate de cele ale lui Chomsky ? Da, Chomsky a fost tot timpul foarte controversat, dar fără stimulentele sau nu știu ce ne-am fi făcut, inclusiv în LC și în LM, în ciuda faptului că el nu s-a prea referit explicit nici la LC, nici la LM. Faptul că gramaticile context free se află din nou, începând cu anii '80, în centrul atenției în LC nu spune ceva ? Iar faptul că aceleași gramatici (cu extensiunile lor) au marcat, încă din anii “60, teoria limbajelor de programare, domeniu în care ținta programării în limbaj natural se află în actualitate, nu este și el semnificativ ? LC are mai multe părți, mai multe orientări, mai multe niveluri de abstracție, care comportă criterii diferite de evaluare. DT îl asociază pe D. Hays la ideea sa privind falimentul LM și lansarea, drept consecință, a LC. Ca unul care a cunoscut bine cercetările lui Hays (a se vedea frecvența citărilor numelui său în lucrările subsemnatului) și l-a cunoscut și personal foarte bine, fiind invitatul său ca “plenary speaker” la Institutul de lingvistică al Americii (SUNY, Buffalo, 1971), pot depune mărturie că acest autor vedea în LM și LC două domenii solidare, două fețe ale aceleiași medalii, așa cum se va vedea din citatul pe care-l vom da mai jos. Desigur, Hays a avut un rol important în anii de pionierat ai LM și LC, dar ideea unei competiții între ele i-a fost străină. Voi evoca aici intervenția sa la cea de a treia Conferință Internațională de LC

(COLING, September 1971): “The field and scope of Computational Linguistics” [2]. Cităm ([2]:p.23):

“Solomon Marcus says that formal linguistics is a pilot science, emphasizing at the same time that the ordinary field of linguistics is not. But that is to say that linguistics as a branch of mathematics will supply methods to many fields of science, whereas linguistics as a descriptive field, a branch of natural history or natural science, does not. [...] A four-way scheme can be arranged, with psychology, computation, formal linguistics, and descriptive linguistics at the poles. Psychology and computation are about performance, formal and descriptive linguistics are about competence, computation and formal linguistics are abstract, and psychology and descriptive linguistics are sciences. But two other fields have to find places in this scheme: psycholinguistics joins psychology with linguistics and seems at this time a most fruitful field, one in which great progress can be made with benefit to both parent fields. Correspondingly, on the abstract side, COMPUTATIONAL LINGUISTICS JOINS COMPUTATION WITH FORMAL LINGUISTICS (subl. mea, S. M.) and also seems a fruitful area, one in which RAPID PROGRESS CAN BE EXPECTED WITH BENEFIT TO BOTH PARENT FIELDS (subl. mea, S. M.) and with beneficial application to psycholinguistics”.

Referirea pe care o face Hays la subsemnatul are în vedere sloganul, pe care l-am folosit de mai multe ori, “formal linguistics as a pilot science”, unde sintagma “formal linguistics” era folosită ca un echivalent al LM. Iată deci că Hays vedea în LC o alianță a LM cu computaționalul, alianță de natura să imprime un progres rapid atât în LM cât și în domeniul computațional. Cei 30 de ani scurși de atunci au confirmat-o pe deplin. Denumirile folosite pentru preocupările la interferența limbajelor, informaticii și matematicii au variat tot timpul și nu cred că acest aspect merită prea multă atenție. Lingvistica matematică? computațională? inginerească? algebrică? cognitivă? aplicată? cantitativă? teoretică? statistică? probleme matematice ale semioticii? tehnologia limbajului? limbajul în inteligența artificială? lingvistica inginerească? procesarea limbajului natural? “information storage retrieval”? lingvistica cibernetică? pe fiecare dintre acestea am întâlnit-o și propriile mele articole au fost publicate aproape sub fiecare dintre etichetele de mai sus. Iată și câteva detalii semnificative ale istoriei.

În 1962 s-a înființat în USA “Association of Computational Linguistics”.

În 1963 Ferenc Kiefer a demarat la Budapesta revista “Computational Linguistics”, care a trăit peste zece ani. Conferința de la Grenoble de “traitement automatique des langues” din 1967 era a treia de acest fel, fiind precedată de o alta, la New York, în 1965 și de una în Anglia, probabil în 1963, organizată de M. Masterman. Între timp, la ruși, numeroase conferințe au avut loc pe tema “avtomaticeskaja obrabotka tekstov” iar “Sprachkunde und Informationsver- arbeitung” a fost uneori eticheta folosită de germani ș.a.m.d. Nu negăm rolul important pe care l-a avut David G. Hays în dezvoltarea CL, dar acest rol a fost altul decât cel afirmat de DT. Emergența LC s-a produs încă din anii ‘50, sintagma LC a devenit curentă încă de la începutul anilor ‘60. Șirul de conferințe COLING nu a făcut decât să continue această tradiție. Alții au preferat folosirea

sintagmei LM (a se vedea, de exemplu, "Prague Bulletin of Mathematical Linguistics", "Prague Studies of Mathematical Linguistics", revista japoneză "Mathematical Linguistics" (în echivalentul ei japonez) etc. În ceea ce privește însă profilul acestor reviste, nu am constatat o diferență față de cele de CL. Desigur, între timp au început să apară și unele publicații mai specializate, cu referire la părți determinate ale CL (cum ar fi cea relativă la corpusul lingvistic). Etichetele nu au avut importanță și nu știu să se fi desfășurat vreo competiție între ele. Chiar Hays a folosit diverse etichete, de exemplu cea din [3]. Dar DT merge mai departe pe ideea sa și afirmă (în completă discordanță cu viziunea lui Hays, de la care se reclamă) că "metodele LM și-au atins limitele" (încă în urmă cu peste 30 de ani!), pentru ca numai două pagini după aceasta afirmație (deci la pagina 135 din [1]) să afirme că e nevoie de "modele formale ale limbii la toate nivelurile ei (fonetică, morfologie, sintaxă, discurs) gramatici formale [...]". Cum vede DT aceste modele formale altfel decât sub formă logico-matematică? Știe oare că multe modele de acest fel există de câteva decenii? Indicații bibliografice asupra lor sunt date parțial în [4], [5], [6], [7] iar pentru cercetările românești în [8], [9]. Desigur, aceste modele sunt inegale ca valoare, au nevoie de continuări, modificări, ameliorări, dar ele nu pot fi ignorate. Fonetica, fonologia, vocabularul, morfologia, sintaxa, semantica lingvistică și lingvistica istorică au beneficiat din plin de metodele matematice, așa cum se poate vedea din impactul deosebit al lucrărilor respective în literatura de specialitate; DT indică, drept domeniu al LM, numai "aspectul sintactic, gramatical", despre celelalte nu a aflat. Nu a aflat nici că LM a abordat și aspecte analitice, nu numai pe cele generative. DT definește "dimensiunea fundamentală" a LC prin "fezabilitatea instanțierii unei descrieri lingvistice cât mai complete, mentenabilitatea acestei instanțieri și, desigur, conformantă cu realitatea uzului limbii". ([1]: 133). Cu un mic efort înțelegem despre ce este vorba. Desigur că problemele de complexitate, de cost, nu puteau fi încă abordate în anii '50 și '60 cu mijloacele cu care ele au început a fi studiate în a doua jumătate a anilor '70, când instrumentele elaborate în informatica matematică deveniseră mult mai perfecționate. Dar acest fapt nu ține, cum crede DT, de alegerea între LM și LC, ci de progresul general realizat în știință. Pentru a mă referi la propria noastră experiență, atunci când, în 1969, prezentăm la COLING-ul din Suedia gramaticile contextuale nu aveam cum să mă ocup de aspectul complexității acestor gramatici în maniera în care s-a putut face acest lucru ulterior (a se vedea, de exemplu, [10]). Dar acest fapt nu are nici o legătură cu eticheta folosită.

Anii '80 și '90 au confirmat necesitatea unui orizont cât mai larg în domeniul computațional. Nu m-am mirat atunci când "Encyclopedia of Microcomputers" și "Encyclopedia of Computer Science and Technology" mi-au solicitat o contribuție cu tema "Semiotics and Formal Artificial Languages" (a se vedea [11]) și nici când "Handbook of Formal Languages" mi-a solicitat un capitol privind "Contextual Grammars and Natural Languages"[12] iar o lucrare preponderent teoretică a fost inserată în "Computational Linguistics in the Netherlands 2000"[13]. Nu m-am mirat nici când am văzut că o revista cu titlul "Linguistics and Philosophy" publica articole excelente de LC. Interferențele merg în toate direcțiile și ele caracterizează cultura contemporană. În acest orizont trebuie să ne plasăm, cred, atunci când ne referim la disciplinele cognitive care se dezvoltă sub ochii

noștri și își pun amprenta pe modul nostru de gândire și de comportare. Un tratat ca "Mathematical Methods in Linguistics" [14] include multe fapte de LC, deși în titlul său nu figurează epitetul "computational". O revistă ca "Theoretical Linguistics" (1970-2000), publicată de Walter de Gruyter (Berlin-New York) a inclus multe articole vizând aspecte matematice și/sau computaționale, deși numele revistei nu indică acest lucru. Chiar o revistă mai tradițională, ca "Linguistics" a inclus de multe ori articole de LM și nici "Foundations of Language" nu a procedat altfel. Multe fapte de LM și de LC se plasează în mod natural în orizontul semioticii computaționale. Era internetului impune desigur o problematică nouă, față de care abordările anterioare se pot dovedi insuficiente. Salutăm inițiativa noii generații de cercetători de a se dedica noilor probleme. Dar trecerea de la ieri la azi și de la azi la mâine nu poate fi decât una care ține seama în mod critic de experiența acumulată. Din tot ceea ce am prezentat mai sus rezultă clar că LM și LC au fost mereu împreună și că, în general, etichetele nu au contat prea mult. Unii au mers chiar mai departe; astfel, în capitolul 4, "Mathematical and Computational Linguistics", din [15], se afirmă pur și simplu (p.86): "Mathematical linguistics has also been called theoretical linguistics and even computational linguistics". Iar mai departe, în același loc: "Computational Linguistics originated around 1950 with the initiation of research on automatic translation" (se trimite la o carte editată de D.G.Hays [3] și la o alta avându-l ca autor pe acesta [16]).

Ca unul care crede în legătura naturală a lingvisticii cu matematica, am încercat o deosebită satisfacție să trăiesc momentul în care această legătură a fost acceptată de ambii parteneri și că de multe ori nici nu mai e nevoie de accentul retoric al epitetului "matematica"; LM este acceptată pur și simplu ca lingvistică. Suntem convinși că o traiectorie similară o urmează și LC iar unele semne în această privință există de pe acum, așa cum am arătat mai sus.

LC este de mai mulți ani o secțiune la congresele internaționale de lingvistică iar LM și LC au secțiunea lor în reviste internaționale de referate ca "Language and Language Behavior Abstracts". În România, minți luminate ale anilor '60, ca profesorii Al. Rosetti, Grigore Moisil și Tudor Vianu, au înțeles schimbările care se profilau și au sprijinit proiectul înființării unei secțiuni de "lingvistica aplicată" la Facultatea de Limbă și Literatură Română a Universității din București, dar s-au găsit alții care să-i torpileze.

La Academia Română a funcționat mulți ani "Comisia de Lingvistică Matematică" iar revista "Cahiers de Linguistique Theorique et Appliquee", înființată în 1962, a fost multă vreme expresia colaborării lingvisticii cu matematica și cu informatica. În ciuda forțelor adverse, s-a reușit în acei ani atragerea unor studenți străluciți ai unor facultăți umaniste la cercetarea limbii și literaturii cu mijloacele LM și LC. Pentru a da numai două exemple de actuali profesori universitari care au susținut teze de doctorat de acest tip, voi menționa pe Pia Brinzeu, de la Catedra de Engleză a Universității din Timișoara și pe Mihai Dinu, de la Facultatea de Litere a Universității din București. Tot în acea perioadă și-a susținut teza de doctorat Sorin Cristian Niță, pe o tema de critică textuală automată

privind înlănțuirea (filiația) diferitelor variante ale “Istoriei Țării Românești” (Șerban Cantacuzino).

Iată însă că, în pofda realităților puse în evidență mai sus, în ([1]: 134) se scrie: “În România, cercetările în domeniul LC și al prelucrării limbajului natural, precum și primele rezultate practice au apărut la începutul anilor ‘80 [3, 4, 5, 6]”.

La ce trimit numerele indicate în paranteze? La o bibliografie de 24 de titluri în care aproape toate (dar toate cele indicate între paranteze) încep cu DT (ignorându-se regula generală în lumea științifică, a așezării numelor autorilor aceluiași articol în ordine alfabetică; dar nu acest fapt este cel care ne interesează în momentul de față). Să observăm că încă în 1978, în articolul “Mathematical and Computational Linguistics” [9] de prezentare a activității din România în domeniul LM și LC se face referire la peste 400 de articole publicate de 130 de autori români și sunt menționați peste 300 de autori străini (unii dintre ei, nume de vază ale LM și LC din acea perioadă) care au citat și continuat cercetările românești. Să mai adăugăm că numeroși lingviști români dintre cei mai importanți au citat și folosit rezultatele școlii românești de LM și LC. Iată că vine acum DT și face (deliberat sau nu) din tot acest efort un teren viran care-l aștepta pe DT să tragă primele jaloane. Nu e cam mult?

Să fim bine înțeleși. Nu noi avem nevoie de încă o citare pe lângă miile de citări deja acumulate, ci noile generații de studenți și de cercetători au dreptul la o informare corectă asupra dezvoltării LM și LC în general și, în particular, asupra LM și LC în România. DT a mai publicat, în urmă cu câțiva ani, un articol în care se schița o privire istorică asupra LC în România, cu câteva citări la întâmplare, care trădau necunoașterea situației reale.

Mai este un aspect care cere o precizare. În conformitate cu specificul volumului în care apare articolul [1], DT face numeroase referiri la acte și documente ale unor organisme europene și internaționale, cum este și firesc, pentru a nu mai vorbi de aspectul financiar al colaborării cu organismele respective. Această situație a existat de la începutul LM și LC (chiar dacă nu a avut amploarea de azi), datorită faptului că LM și LC au apărut și ca urmare a unor comandamente sociale, privind precaritatea mijloacelor de prelucrare a informației. Îmi amintesc de faimoasele Rapoarte CETIS care veneau de la EURATOM, Bruxelles, pe teme legate de analiză și prelucrarea automată a limbajului, traducere automată și documentare automată. În USA, diferite corporații (cum ar fi RAND Corporation, Santa Monica, Calif.) finanțau cercetări similare. O întâlnire semnificativă a fost aceea din 1962, organizată de “NATO Advanced Summer Institute”, la Veneția, Italia, privind traducerea automată. De numele acestui Institut este legat un document care a marcat evoluția cercetărilor de traducere automată: seria de expuneri prezentate de Y. Bar-Hillel [17]. În legătură cu aceste activități dirijate și finanțate de diferite organisme europene și internaționale, trebuie să observăm că cei implicați au avut înțelepciunea și pricepera necesare pentru a nu reduce proiectele respective la dimensiunea lor exclusiv utilitară, ci de a o subordona pe aceasta unei perspective mai ample, care lua în considerare orizontul științific real al problemelor. Pentru a da un prim exemplu, mă voi referi la faptul

că mai multe rapoarte CETIS au pus în discuție un concept care, născut din experimentele de traducere automată, avea să se dovedească de o deosebită semnificație pentru teoria sintactică în toată generalitatea sa; este vorba de conceptul de proiectivitate sintactică, cu consecințe bogate în studiul structurilor arborescente și al gramaticilor de dependență. Azi putem spune că și sintaxa limbajului natural și teoria matematică a grafurilor au profitat esențial de conceptul respectiv (folosit până și de Rene Thom, în probleme de morfogeneză [17]). Această expansiune a unui concept sau rezultat dincolo de motivația sa inițială este testul cel mai convingător al interesului său. Un al doilea exemplu se referă la titlul provocator folosit de Bar-Hillel pentru expunerile sale: “Patru conferințe despre lingvistica algebrică și traducerea automată”.

Simpla alăturare a celor două sintagme, una foarte teoretică, cealaltă aparent tehnologică, avea menirea să-i avertizeze pe cei care presau să se obțină cât mai repede rezultate practice asupra faptului că proiectele de traducere automată nu se pot finaliza de azi pe mâine, ci au nevoie de un lung itinerar lingvistic, matematic și computațional. Acum știm că acest itinerar continuă și azi, cu tatonări și reveniri, și, chiar dacă nu a dus încă la rezultatele visate, a impulsionat în mod esențial cercetările de AI, cu consecințe benefice pentru aspectele logice și semantice ale limbajului natural.

Întrebarea pe care ne-o punem, dar o lăsăm deocamdată fără răspuns, deoarece nu suntem pregătiți pentru a-l da, este următoarea: Nu cumva aspectele pe care le-am criticat mai sus sunt consecința unui fenomen mai general, acela al unui orizont insuficient de cuprinzător, al unei prea mari dependențe de factori utilitari imediați? Știința a oscilat mereu între cognitiv și utilitar, dar istoria arăta că funcția utilitară s-a manifestat în toată profunzimea ei atunci când ea a fost fructul unei evoluții firești a funcției cognitive, evoluție care poate fi de doi ani, de 20 de ani, de 200 sau de 2000 de ani. Cu un ochi îndreptat spre comisiile europene, suntem obligați totuși să ținem treaz și celălalt ochi, îndreptat spre ceea ce se întâmplă pe scena cercetării științifice vii, așa cum apare ea în revistele de specialitate și la întâlnirile științifice de profil. Istoria generală a științei și, în particular, scurta istorie a LM și LC, sunt pline de învățăminte în această privință.

Referințe bibliografice:

- [1] D. Tufi°. *Promovarea limbii române în SI-SC*. În *Societatea Informațională – Societatea cunoașterii* (coord. F. Gh. Filip). Ed. Expert, București, 2001, 131–142.
- [2] D. G. Hays. *The field and scope of computational linguistics*. Papers in Computational Linguistics (eds. F. Papp, G. Szepe). Proceedings of the Third International Meeting of Computational Linguistics, held in Debrecen, Hungary, 1971. Akademiai Kiado, Budapest, 1976, 21–26.
- [3] D. G. Hays (ed.). *Readings in Automatic Language Processing*, American Elsevier, New York, 1967.
- [4] S. Marcus. *Mathematical Linguistics in Europe. Current Trends in Linguistics* (Th. A. Sebeok, ed.), vol.9, Mouton, The Hague, 1972, 646–687.

-
- [5] S. Marcus. *Mathématique et Linguistique*. In *Mathématique, Informatique et Sciences Humaines*, Paris, 26, 1988, 103, 7–21.
- [6] S. Marcus. *The status of research in the field of analytical algebraic models of language*. In *Current Issues in Mathematical Linguistics* (C. Martin–Vide, ed.). Elsevier – North Holland, Amsterdam, 1994, 3–21.
- [7] S. Marcus. *Lingvistica matematica, azi*. In *Matematica in lumea de azi si de maine* (C. Iacob, coord.), Editura Academiei, Bucuresti, 1985, 182–186.
- [8] S. Marcus. *Recent Romanian investigations in the field of mathematical and computational linguistics*. Avtomatizeskaja Obrabotka Tekstov, Matem. Fyz. Fakulta, KL Praha, 1973, 15–42.
- [9] S. Marcus. *Mathematical and computational linguistics*. In *Current Trends in Romanian Linguistics* (A. Rosetti, S. Golopentia Eretescu, eds.). Revue Roumaine de Linguistique 23, 1978, 1–4, 559–588.
- [10] S. Marcus, C. Martin–Vide, G. Paun. *Contextual grammars as generative models of natural languages*. Computational Linguistics 24, 1998, 2, 245–274.
- [11] S. Marcus. *Semiotics and formal artificial languages*. In *Encyclopedia of Computer Science and Technology* (A. Kent, J.C. Williams, eds.) 29, Ed. Marcel Dekker, New York, 1994, 393–405; also in *Encyclopedia of Microcomputers* (A. Kent, J.C. Williams, eds.) 15, 1995, 299–312.
- [12] S. Marcus. *Contextual grammars and natural languages*. *Handbook of Formal Languages* (G. Rozenberg, A. Salomaa, eds.), 2, Springer, Berlin, New York, 1997, 215–235.
- [13] S. Marcus, C. Martin–Vide, G. Paun. *A new–old class of linguistically motivated regulated grammars*. *Computational Linguistics in the Netherlands 2000* (W. Daelemans et al., eds.), Selected Papers from the Eleventh CLIN Meeting, Ed. Rodopi, Amsterdam, New York, 2001, 111–125.
- [14] B. H. Partee, A. Ter Meulen, R. Wall. *Mathematical Methods in Linguistics*. Kluwer, Dordrecht et al, 1990.
- [15] E. F. Beckenbach, Ch. B. Tompkins (eds.). *Concepts of Communication: Interpersonal, Intrapersonal and Mathematical*. John Wiley and Sons, New York, 1976.
- [16] D. G. Hays. *Introduction to Computational Linguistics*. American Elsevier, New York, 1967.
- [17] R. Thom. *Stabilité Structurale et Morphogenèse*. John Benjamins, New York, 1970.
- [18] Y. Bar–Hillel. *Four Lectures on Algebraic Linguistics and Machine Translation* revised version of a series of lectures given in July 1962, before a NATO Advanced Summer Institute, Venezia, Italy.

Între lingvistica matematică și cea computațională: o altă perspectivă

Dan TUFÎȘ

1. În loc de introducere

Dat fiind că acest articol este un comentariu asupra filipicei de neînțeles "Între lingvistica matematică și cea computațională" a domnului Solomon Marcus, membru titular al Academiei Române, mărturisesc că elaborarea sa fost o întreprindere asupra căreia am avut multe ezitări iscate din incertitudinea receptării sale corecte, constructive. Din păcate majoritatea afirmațiilor și implicațiilor pe care domnia sa le face în articolul amintit, sunt inexacte și umorale. Nu mai insist și asupra decontextualizării citatelor din lucrarea mea [1], procedeu neelegant. Este binecunoscut din logica clasică faptul că dintr-o serie de premise false se poate demonstra orice. În ciuda ezitărilor amintite, violenta polemică lansată de domnul Solomon Marcus prin articolul menționat îmi oferă posibilitatea de a aduce în discuție elemente de istorie a domeniului ce ar putea fi de interes, cu precădere pentru cititorii al căror domeniu de specialitate nu este prelucrarea automată a limbajului natural. Pentru specialiștii în domeniul prelucrării limbajului natural, majoritatea argumentelor pe care le voi aduce sunt bine cunoscute.

Ca modalitate de documentare, am optat pentru includerea integrală a materialului produs de domnul Academician Marcus, indentat și redat cu caractere italice. De asemenea, am păstrat secțiunea domniei sale de referințe bibliografice. Lucrările pe care le-am citat eu sunt documentate în cuprinsul textului, prin includerea referinței complete între paranteze rotunde. Singura excepție este lucrarea mea, sursa nemulțumirii domnului Marcus, care este referită de amândoi ca [1]. Cititorul va putea face astfel mai ușor distincția între cele două categorii de referințe. Înainte de a proceda la analiza afirmațiilor domnului Academician Marcus, aș dori să fac unele precizări:

- contextul discuției în [1], ca și în cele ce urmează, este cel al tehnologiei limbajului, al cercetărilor foarte intense în întreaga lume pentru dezvoltarea de sisteme inteligente capabile să faciliteze comunicarea dintre doi sau mai mulți conlocutori (oameni sau sisteme software), prin intermediul limbajului natural;
- în raport cu lucrarea [1] domnul Academician Marcus se oprește cu îndârjire asupra a doar trei fraze interpretate ca atac la persoana sau activitatea sa

științifică și se referă ironic (și după cum se va vedea în continuare, în mod nejustificat) la alte două, făcând abstracție de restul prezentării care nu are nici o contingentă cu domnul Marcus. Domnul Academician are merite pe care nu i le poate lua nimeni, are contribuții importante în mai multe domenii și este creatorul școlii românești de lingvistică matematică. Interesul domniei sale pentru aspectele legate de implementarea pe calculator a programelor de prelucrare a limbajului natural a fost minim. Îmi reamintesc o discuție pe care am avut-o în anul 1991 la câțva timp după ce mă întorsesem de la Conferința Europeană de Lingvistică Computațională organizată la Berlin de profesorul Jurgen Künze. Cu acea ocazie, domnul Academician Marcus mi-a mărturisit că îl cunoaște de multă vreme pe profesorul Künze și că au și colaborat o perioadă cât amândoi au avut ca domeniu de preocupări lingvistica matematică. La sfârșitul anilor '60, mai spunea domnul Marcus atunci, drumurile celor doi s-au despărțit, profesorul Künze optând pentru noua paradigmă a lingvisticii computaționale.

- Domnul Academician Marcus a scris enorm, în domenii extrem de variate, aici mă refer în special la cele legate de studiul limbii, și prin urmare era inevitabil să nu atingă subiectul foarte actual al prelucrării automate a limbajului natural. A făcut-o însă detașat de nivelul inerent perisabil al tehnologiei informatice. O teorie științifică, un model formal teoretic sau transpus într-o implementare a unui program software sunt inevitabil supuse „eroziunii” timpului, unele mai rapid altele mai lent. Lucrarea [1], despre care discutăm, ia în discuție exact acest cadru al investigației tehnologice și a măsurilor științifice, tehnice, organizatorice și chiar legislative pentru a crea o bază perenă a cercetării și dezvoltării tehnologice privind prelucrarea automată a limbii noastre: resursele computaționale fundamentale ale limbii române. Societatea Informațională-Societatea Cunoșterii este caracterizată de vectori tehnologici și funcționali [M. Drăgănescu: „Societatea informațională-societatea cunoașterii. Vectorii societății cunoașterii”, In *Societatea Informațională – Societatea cunoașterii* (coord. F. Gh. Filip). Ed. Expert, București, 2001, 43–112.] a căror ignorare este nu numai neproductivă dar și periculoasă. „În era electronică, **este esențial pentru supraviețuirea unei limbi ca ea să fie folosită în sistemele de informare electronică**” afirmă fără echivoc Alain Danzin în influentul raport al Comisiei Europene „Towards a European Language Infrastructure” întocmit în 1992 prin consultarea a 182 de specialiști din cercetare și industrie. Promovarea limbii române în contextul informațional al societății cunoașterii este un obiectiv actual și de viitor și nu poate fi subiect de dispută în viața științifică românească;
- deși este un truism, cred că pentru evitarea unor interpretări greșite este necesar să subliniez faptul că în dezvoltarea programelor de inteligență artificială, de prelucrare a limbajului natural sau în general în ingineria software, o mulțime de discipline matematice (teoria algoritmilor, teoria complexității, teoria limbajelor formale, teoria categoriilor, statistica

matematică și multe, multe altele) sunt fundamente indispensabile în avansul științific și tehnologic al acestor discipline (și desigur nu numai al lor). Programarea (ca și matematica elementară) sau utilizarea de produse informatice sunt activități la îndemâna tuturor (de altfel reflectate și în programele școlare de învățământ), dar proiectarea și realizarea de programe software inteligente necesită o pregătire teoretică solidă, talent și multă muncă. Diferența între două programe care realizează aceleași prelucrări dar unul în câteva secunde și altul în câteva ore, apare tocmai din diferența de pregătire teoretică și talent a autorilor lor.

- domeniul științei și tehnologiei informației este poate cel mai dinamic sector al activității creative: Bill Gates spunea că dacă de pildă industria automobilelor ar fi avut aceeași dinamică cu cea a calculatoarelor, acum o mașină ar trebui să coste 1 dolar. Fantasticul ritm de dezvoltare al tehnologiei hardware (bazată pe importante descoperiri științifice obținute în ultimii 50 de ani) nu a fost nici pe departe egalat de ritmul dezvoltării în domeniul software. În ciuda acestui decalaj, știința ingineriei software și-a reinnoit instrumentarul teoretic (modele și/sau formalisme) cu o viteză neîntâlnită în alte domenii științifice. Dinamica fără precedent a cunoașterii în știința și tehnologia informației obligă omul de știință din acest domeniu la o informare continuă, din ce în ce mai specializată și mai selectivă. Se estimează că în acest domeniu se scriu în fiecare zi mai multe articole decât poate citi un om în întreaga sa activitate și că informația mai veche de 15-20 ani este foarte probabil să fie perimată (desigur cu excepțiile ce întotdeauna confirmă regula). Evoluția terminologică în acest domeniu este încă o mărturie vie a dinamicii de care aminteam: în domeniul prelucrării limbajului natural se vorbește acum de ontologii lexicale, de gramatici lexicalizate susținute de ontologii, de analiză (parsing) ontologică, de lingvistica WEB-ului și WEB-ul semantic, de resurse lingvistice standardizate și așa mai departe.
- referitor la antinomia „lingvistică matematică-lingvistică computațională” pe care domnul Academician Marcus mi-o atribuie, vreau să precizez că nicidecum nu am afirmat că cele două domenii se exclud reciproc sau că ar fi în competiție; pur și simplu ele sunt subsecvente din punctul de vedere al relevanței față de problemele pe care le discutăm aici. Există fără îndoială o filiație între ele, în sensul că lingvistica computațională a preluat o mare parte din instrumentarul lingvisticii matematice (nici nu se putea altfel) dar ce a adus nou lingvistica computațională, pe lângă noi modele și formalisme, este în primul rând de natură metodologică și tehnologică: experimentul și evaluarea. Ceace se numește astăzi lingvistică computațională teoretică este în mare măsură asimilată cu lingvistica formală modernă. Acest segment al lingvisticii computaționale a moștenit de la lingvistica matematică cel mai mult și adecvându-și metodele la realitățile tehnologice a produs și este de așteptat să producă noi rezultate validabile și incorporabile în sisteme automate de prelucrare a limbajului

natural. Teoriile și formalismele lingvistice, azi în vogă în lingvistica computațională (TAG, LFG, HPSG, CG, CUG), au fost produse de lingvistica formală și prin validarea instanțierilor pe segmente de limbă netriviale, au devenit instrumente operaționale ale prelucrării limbajului natural. Dezvoltarea de modele de limbă, analiza algoritmilor de prelucrare a limbajului (resursele de calcul necesare unei implementări funcționale, viteza de răspuns), construcția (achiziția) resurselor lingvistice standardizate, gradul de acoperire lingvistică al unei formalizări lingvistice (cunoștințe lingvistice=resurse lingvistice), sunt doar câteva direcții definitorii ale metodologiei lingvisticii computaționale.

- în sfârșit, în raport cu obiectivele finale urmărite de **implementarea** unui model de prelucrare a limbajului se remarcă în ultimii circa 10 ani o departajare și chiar o competiție (fără însă a fi o antinomie) între abordările introspective-principiale și cele inductive, bazate pe date. Prima categorie de abordări este caracterizată de dezvoltarea prin introspecție științifică de teorii și formalisme gramaticale **computaționale** (imensa lor majoritate bazate pe restricții și unificare categorială cu accentuată lexicalizare) și mai apoi instanțiate manual de experți lingviști. Cea de a doua abordare, ce câștigă foarte mult teren în ultima perioadă, este cea bazată pe tehnicile învățării automate ce pornesc de la premiza că, într-un corpus lingvistic reprezentativ și de dimensiuni mari, există suficientă informație privind regularitățile dintr-o limbă (cea în care sunt textele ce alcătuiesc corpusul lingvistic) astfel încât, tehnici adecvate de învățare automată să fie capabile să construiască un model de limbă robust și de mare acoperire lingvistică. Aș mai menționa că, în fapt, de multe ori cele două abordări sunt combinate (cu preponderența uneia dintre ele). Într-un anumit sens, acest dualism în abordările modelelor de prelucrare automată a limbajului natural continuă a celebra confruntare de idei între Chomsky și Piaget susținătorii teoriilor înăscutului (innate) și respectiv al învățării în explicarea facultății umane a limbajului.

Cu aceste lămuriri preliminare, voi analiza în continuare afirmațiile domnului Academician Marcus cu sincera speranță că cititorii acestui text, dar mai ales domnia sa, vor înțelege că preocupările mele și ale distinsului profesor au alte obiective, motivații și desigur modalități foarte diferite de finalizare. Acest lucru nu înseamnă că rezultatele fiecăruia dintre noi le anulează sau le diminuează pe ale celuilalt (cu atât mai mult cu cât recunoașterea internațională există pentru amândoi). După cum la fel de bine diferențele de perspectivă și opinii, naturale în fond, nu înseamnă că nu avem a ne spune lucruri interesante unul altuia.

2. O analiză textuală

„Mă simt obligat să reacționez la un anumit mod de prezentare a evoluției ideilor, în cea de a doua jumătate a secolului al XX-lea, în articolul [1] al

d-lui Dan Tufiş (de aici mai departe DT), membru corespondent al Academiei Române. Precizez de la început ca nu contest interesul și utilitatea direcției de preocupări prezentate în [1]; am în vedere numai modul în care aceasta direcție este pusă în relație cu alte cercetări dedicate limbajului.”

Așa își începe domnul Academician Marcus articolul solicitat de mine pentru volumul „Limba Română în Societatea Informațională-Societatea Cunoașterii” rezultat al proiectului INFOSOC „SI-SC: Soluții și strategii în România”. Să urmărim un prim citat incriminat (care în transcrierea dlui Academician este trunchiat și conține niște ghilimele ce nu-mi aparțin; redau mai jos varianta publicată) :

[1: p.133]:

“Din acest punct de vedere (*al folosirii calculatorului în prelucrarea limbajului natural – precizarea mea*), este semnificativ a arăta că însuși numele domeniului de cercetare a prelucrării automate a limbajului natural a suferit modificări reflectând progresele științifice și tehnologice: inițial, desprinzându-se din lingvistica formală, *lingvistica matematică* a încercat dezvoltarea unor modele matematice de reprezentare a limbajelor naturale sau formale (în general al aspectului lor sintactic, gramatical), cautând soluții abstracte de modelare generativă de tip universal a ceea ce se presupunea (la nivelul cunoașterii științifice a anilor 1960) a fi facultatea limbajului. “

Ce l-a supărat aici pe distinsul polemist? Ne spune chiar domnia sa:

„Nu știi ce înțelege DT prin “lingvistica formală”, o sintagmă nu prea folosită în perioada de emergență a lingvisticii matematice; exista lingvistica structurală (altceva decât ceea ce ar putea fi lingvistica formală, adică bazată pe formalizare în sensul logicii matematice moderne), care desigur a constituit una din sursele lingvisticii matematice (de aici mai departe LM), așa cum i se pot indica și alte surse (biologice, logice, matematice, psihologice etc.)

Mă surprinde întrebarea retorică cu care începe „argumentația”, și căreia nu-i văd decât un gratuit rol derogativ. Eu nu-mi închipui că domnia sa nu a auzit de antinomia „gramatică descriptivă – gramatică formală” la limitele extreme ea fiind reprezentată de lucrările lui O. Jespersen (O. Jespersen: *The philosophy of Grammar*, Allen & Unwin, London, 1924 și *Analytical Syntax*. Holt Rinehart & Winston, New York, 1937 (republicată în 1969)) și respectiv lucrările timpurii ale lui Chomsky referitoare la lingvistica generativă. Dacă însă mă înșel, o lectură lămuritoare, este influența carte editată de Keith Brown și Jim Miller în Pergamon Press, 1996 numită „Concise Encyclopedia of Syntactic Theories”, cu precădere articolul „Descriptive Grammar and Formal Grammar” de F. Stuurman, al cărui prim capitol se numește chiar Descriptive and Formal Grammar: The Fundamental Opposition. La fel de utilă este și lucrarea monumentală a lui David Crystal „The Cambridge Encyclopedia of Language”, Cambridge University Press, 1987.

Pe de altă parte, o pagină mai încolo, domnul Academician mărturisește că și domnia sa a folosit termenul de lingvistică formală:

În ceea ce privește sintagma “lingvistică formală”, ea a căpatat o anumită utilizare în anii târzii 1960 și în anii următori, iar personal am folosit-o în unele lucrări, după cum se va vedea imediat, considerând-o oarecum echivalentă cu LM (lingvistica matematică);

Pentru lămurirea elementului istoric, furnizez în continuare un citat din recenzia lui R.B. Lees (Language, nr. 33, vol 3, 1957, pp375-408) la faimoasa carte a lui Chomsky (Syntactic Structures, Mouton, The Hague, 1957): „in a sense, transformational analysis is essentially a **formalization** of a long-accepted, traditional approach...”. Citatul apare la pagina 387. Chomsky se pare că a apreciat termenul și l-a adoptat, cel puțin în raport cu propria filozofie generativistă asupra limbajului.

„dar factorul determinant în nașterea LM, în a doua jumătate a anilor “50, a fost dezvoltarea calculatoarelor electronice și, împreună cu ea, a primelor preocupări sistematice de LC (prescurtare a lingvisticii computaționale), numite atunci traducere automată, documentare automată, prelucrarea automată a limbajului, cu diverse variante ale lor în engleza (de exemplu, “machine translation”), franceză, rusă, germană, italiană etc. Din aceste preocupări s-au inspirat primele modele care au constituit noua disciplină a LM.”

Înainte de a face o serie de precizări istorice mai exacte, vreau să notez că de la începutul istoriei sale, domeniul traducerii automate a fost și în mare a și rămas un domeniu distinct de restul preocupărilor legate de prelucrarea limbajului natural. Aș mai observa că textul de mai sus, încearcă să sugereze că LM s-ar fi constituit ca disciplină ulterior LC. Ambiguitatea afirmației de mai sus provine din punerea în relație de concordanță temporală a primelor preocupări în domeniul LC cu apariția domeniului în sine. Oricine știe că un anumit domeniu științific se cristalizează în timp, pe baza unor rezultate științifice promițătoare, a unor experimente convingătoare (în cazul domeniilor tehnologice). Până la sedimentarea elementelor definitorii ale unui domeniu de cercetare, pot coexista sau se pot succeda mai multe direcții de cercetare. Dintre acestea unele pot dispărea sau își pot diminua foarte mult influența în raport cu motivația inițială. Ele își pot continua însă existența prin noi motivații, prin alegerea de noi obiective.

Ca element istoric, aș preciza că în toate evocările pe care le-am citit eu, cel ce pentru prima dată a sugerat ideea folosirii calculatorului și a tehnicilor de decodificare pentru prelucrarea automată a limbajului natural a fost Warren Weaver în 1946. În 1949 el scrie lucrarea „**Translation**” considerată de toți specialiștii în traducere automată ca primul document programatic al acestei discipline. În 1952 a avut loc la Universitatea Georgetown din SUA prima conferință dedicată exclusiv traducerii automate. În 1954, Peter Toma de la Universitatea Georgetown împreună cu un grup de cercetători de la IBM realiza primul experiment de traducere automată (engleza-rusa) folosind un dicționar de 250 de cuvinte și 6 reguli sintactice de rescriere. Acest sistem avea să constituie nucleul faimosului program de traducere automată Systran pe care Peter Toma îl finalizează în 1973.

Punctul meu de plecare s-a aflat în lucrările unor Kulagina și Melciuk, puternic implicați în studiile de traducere automată rusă-franceză, Yves

Lecerf, implicat în problemele de documentare automată, D. G. Hays, implicat în traducerea automată din rusă în engleză și reciproc, B. Vauquois, cu preocupări de informatică lingvistică la Grenoble. De la ei, ca și de la alți autori similari, am preluat în bună măsură ștafeta pe care am căutat s-o duc mai departe. Ceea ce afirm despre mine este valabil pentru cei mai mulți cercetători din domeniul LM din anii 1950 și 1960, cum ar fi Maurice Gross, Masami Ito, A. Trybulec și mulți alții.

Traducerea automată, dar mai ales eșecul primelor încercări de rezolvare a acestui obiectiv încă nerezolvat sau nerezolvat complet, a constituit fără îndoială o motivație a „emergenței” LM. Așa cum voi arăta pe larg mai departe, eșecul proiectelor de traducere automată au fost puse, prin interpretarea unilaterală și tendențioasă a raportului APLAC, exclusiv pe seama inadecvării teoriilor lingvistice folosite atunci și a cantonării în faptul unor limbi particulare. Teoria „facultății innăscute a limbajului” lansată de Chomsky, opunându-se tradiției tipologice de studiu lingvistic prin diversitatea limbilor, a generat o prodigioasă cercetare în direcția determinării principiilor gramaticii universale, în speranța că identificarea și caracterizarea lor riguroasă le-ar putea operaționaliza atât pentru explicarea comunicării umane prin limbaj cât și (un derivat subsidiar al obiectivului lui Chomsky) pentru realizarea de sisteme de traducere automată apropiate de performanța umană.

Dubioasă mi se pare sintagma “soluții abstracte”, probabil efectul unui obicei binecunoscut de a diaboliza abstractul.

Remarca de mai sus mă surprinde de două ori: mai întâi pentru că nu este nimic reprobabil în expresia „o soluție abstractă” (ba chiar dimpotrivă: ”abstract = Care rezultă din separarea și generalizarea însușirilor caracteristice ale unui grup de obiecte sau de fenomene care este considerat independent, detașat de obiecte, de fenomene sau de relațiile în care există în realitate” DEX’96) și apoi referirea la un obicei binecunoscut (al cui?) de diabolizare a abstractului. Nu neagă nimeni că acele soluții abstracte de care aminteam au generat idei valoroase și cercetări computaționale (mai ales în domeniul traducerii automate bazate pe conceptul „interlingua”) dar rezultatele acestor idei și cercetări nu sunt revendicate nici chiar de Chomsky.

În ceea ce privește sintagma “lingvistică formală”, ea a căpătat o anumită utilizare în anii târzii 1960 și în anii următori, iar personal am folosit-o în unele lucrări, după cum se va vedea imediat, considerand-o oarecum echivalentă cu LM; dar chiar dacă nu acceptăm aceasta echivalență, nu putem eluda faptul că lingvistica formală se află în imediata vecinătate a LM.

Cu amendamentele cronologice pe care le-am comentat mai devreme, apropierea între LM și LF (lingvistica formală) este exact ceea ce am afirmat și eu.

DT pretinde ca LM “a încercat”, sugerând astfel ca ea a eșuat în tentativa de modelare a limbajului natural.

În primul rând este vorba de modelarea **computațională** a limbajului. În al doilea rând nu eu pretind acest lucru, dar sunt perfect de acord cu el. Iată câteva opinii ale unor

mari specialiști, activi, din domeniul prelucrării automate a limbajului natural (sublinierile îmi aparțin):

- Christopher Manning and Hinrich Shutze: Foundations of Statistical Natural Language Processing, The MIT Press, 1998:

„...the availability of large text corpora has changed the scientific approach to language in linguistics and cognitive science. Phenomena that were not detectable or seemed uninteresting in *studying toy domains and individual sentences* have moved into the center field of what is considered important to explain.”

- Susan Armstrong-Warwick (editor): Prefața la „Special Issue on Using Large Corpora”, Computational Linguistics, Volume 19, no 1, 1993 p. 4:

„What is that has brought about this rapid growth of interest in corpus-based NLP?...The technological advances in computer power has certainly favoured the approach, as has the growing availability of large-scale textual resources in machine readable form. *More important, perhaps, is the growing frustration of trying to use standard rule-based methods to account for more than a well-chosen fragment of text, regardless of the application.* The data extracted from large corpora have demonstrated that language is more flexible and complex than that which most rule-based systems have up to present tried to account for. The relative lack of practical results at a time when industrial concerns are looking to the CL community to demonstrate progress toward useful applications has also contributed to the growing interest in new methods. And finally, the success rate demonstrated in the speech community offers hope for similar progress in NLP.”

- Nancy Ide and Jean Veronis (editori) Computational Linguistics –Special Issue on Word Disambiguation, Vol. 24, No. 1 1998 p.15:

„Although quantitative methods were embraced in early MT work, in the mid-1960s interest in statistical treatment of language waned among linguists due to the trend toward the discovery of *formal linguistic rules sparked by the theories of Zellig Harris (1951) and bolstered most notably by the transformational theories of Noam Chomsky (1957).* Instead, attention turned toward *full linguistic analysis and hence to sentences rather than texts, and toward contrived examples and artificially limited domains instead of general language.*”

- Victor Yngve: From Grammar to Science:New Foundations for General Linguistics, John Benjamin Publishing Company, 1996:

„there seems to be no scientific way of deciding among the many contenders...We find positions and methods being promoted like a new movie or defended with withering polemics or taken up like the latest fad...*We should abandon logical-domain theories entirely and move to the physical domain...Because this (notation) can be programmed on a computer it can be used to test large-scale models...Gone will be the babel of arbitrary grammatical notations, each to be discarded in turn*”.

Deși nu împărtășesc în întregime poziția extrem de radicală a lui Yngve, ea este simptomatică pentru insatisfacția generală față de abordările tradiționale ale anilor ‘60-‘80.

- R.F. de Bruine (editor) „Synthesis of Proposal for an RTD Programme by Users, Industry and Research in Language and Technology”, DGXIII, Commission of the European Communities, September 1992:

„There is a broad need to further understanding of linguistic phenomena in the context of computerising the analysis and generation of language. General research should be stimulated within the following three main topics:

- research on the linguistic meaning representation at the various level of description, ranging from the lower (e.g. phonetic, morphological and syntactic) and better understood ones to the higher, scientifically more difficult ones (e.g. semantic, pragmatic, contextual and communicative ones). It is foreseen that the former must yield results in the short to medium term. Even if the latter are long-term enterprises, they must be organised in way that ensures availability of usable intermediate results.
- reasearch on *more adequate and efficient computational schemes for natural language processing (e.g. constraints based computing and quantitative aspects) providing the base for robust processing behaviour vz the applications of advanced computer science and statistical methods in close collaboration and synergy with related actions.*
- research into the human factors related with the future spread of advanced language processing technologies taking into account the ergonomics aspects, economic and socio-cultural dimensions.”

Lista unor astfel de citate poate continua pe zeci de pagini, dar am să mă opresc aici nu înainte de a mai reaminti raportul comisiei prezidate de Alain Danzin „Towards a European Language Infrastructure”. Acest document, o adevărată cartă albă a cercetării în domeniul tehnologiilor limbajului, a restructurat complet programele de cercetare și prioritățile pe termen mediu și lung. A o ignora (ba chiar mai mult a o critica fără a-i cunoaște conținutul și a o eticheta ca pe un document birocratic al celor de la Uniunea Europeană) poate fi desigur o opțiune personală, dar cu efectul izolării științifice și mai accentuate.

Ceea ce este deocamdată numai o sugestie devine, după cum se va vedea, o certitudine pentru DT. Într-adevăr, iată ce scrie mai departe DT ([1]: 133):

“Curând metodele lingvisticii matematice și-au atins limitele drept care, în anul 1966, la propunerea lui David Hays, domeniul de cercetare al limbajelor naturale, din perspectiva utilizării acestora în interacțiunea cu calculatoarele electronice, este individualizat sub numele de lingvistică computațională”.

Chestiunea cu atingerea limitelor ține de domeniul umorului involuntar și trecem peste ea, dar nu ne miră, după ce am văzut la ce se reduce LM pentru DT.

În ciuda repetatelor mele clarificări, și după cum se observă și din citatul de mai sus, referirea mea era la utilizarea metodelor lingvisticii matematice în programele de prelucrare a limbajului și nicidecum la domeniul în sine. Probabil că pentru cine nu a încercat să realizeze un sistem de prelucrare a limbajului natural și nu s-a lovit de

problemele implementării unui dicționar și a unei gramatici computaționale e mai greu de înțeles remarcă mea anterioară. Domnul Academician Marcus nu s-a apropiat niciodată de problemele unei implementări și prin urmare nu mă surprinde lipsa de înțelegere a diferenței între o definiție formală a unei gramatici (de exemplu) care se explicitează în câteva rânduri și implementarea unei gramatici computaționale care nu numai că nu încapă în câteva sute sau mii de pagini dar reclamă o muncă exprimată convențional în mii de oameni/an. Gramatica computațională a limbii engleze, dezvoltată în cadrul proiectului Alvey, a fost rezultatul a 10 ani de muncă intensă a celor mai importante 12 colective de cercetare din Anglia, fiecare dintre acestea fiind conduse de cercetători importanți și fiind suplimentate cu numeroși studenți doctoranzi. Gramatica GPSG dezvoltată este unul din exemplele standard de gramatică introspectivă de mari dimensiuni. Un astfel de efort uman și financiar nu este la îndemâna multor societăți. Și experiența a arătat că nici nu este necesar! Ralph Grisman, de la Universitatea din New York a demonstrat că programul sau de inducție gramaticală, pe baza unui corpus de antrenare a generat o gramatică nucleu, a cărei „finisare” a durat mai puțin de două săptămâni și, confruntată cu gramatica Alvey pe un text arbitrar a reușit să analizeze mai multe fraze, cu alte cuvinte a demonstrat o mai mare acoperire lingvistică.

Nu mi-am imaginat niciodată că între LM și LC ar putea avea loc o competiție, prima definindu-se prin metoda (căci ce altceva este LM decât studiul limbajului cu ajutorul matematicii ?) iar a doua prin obiectivul pe care și-l propune. LM nu poate ignora problematica LC iar LC nu-și poate realiza proiectele fără LM. Probabil însă că DT lucrează cu o definiție specială a LM, pe care am dori s-o aflăm.

Nici nu există această competiție decât în imaginația domnului Academician care sugerează mai sus că LC nu folosește matematica sau că atunci când o face, disciplina se numește LM. Ceea ce, așa cum am arătat mai înainte, este fals. Elementele suplimentare, esențiale și definitorii sunt calculatorul, algoritmi eficienți și cunoștințele cu care acesta trebuie „hrănit”. O formalizare a procesului de înțelegere și/sau producere a limbajului natural, de orice sorginte ar fi ea, nu este decât o ipoteză asupra unui fenomen încă neelucidat. Validarea acestei ipoteze este cheia care a diferențiat LC de LM. În anexa acestei lucrări am furnizat două definiții pentru LM și LC. Prima definiție (LM) aparține lui Geoffrey K. Pullum and Andras Kornai iar cea de a doua (LC) se află în pagina WEB a Asociației de Lingvistică Computațională (al cărui membru sunt din 1985). Aș mai face precizarea că lingvistica teoretică modernă (în sensul precizat mai înainte) studiază limbajul nu numai cu ajutorul matematicii. Alături de matematică, sociologia, psihologia, medicina și științele cognitive constituie domenii ale cunoașterii care sunt fundamental implicate în explicarea acestui miracol pe care îl reprezintă comunicarea inter-umană. Incapacitatea actuală de a realiza un procesor artificial de limbaj la nivelul performanței și competenței umane se datorează nedescifrării (încă) a mecanismelor minții și creierului omului. Dihotomia structural-fenomenologic și noile cercetări în direcția unei științe integrative (reprezentată între alții de lucrările de pionierat ale Academicianului Mihai Drăgănescu) sunt fără îndoială porți deschise spre cunoașterea, în viitor, mai exactă a minții și implicit a facultății limbajului. Până atunci, obiectivele LC (realizarea de sisteme automate capabile

să prelucraze limbajul natural) apelează la modele aproximative, a căror acceptabilitate se probează prin implementarea și evaluarea lor pe date reale. Cum între afirmarea unui obiectiv de LC și realizarea sa operațională este o distanță mare, pe care uneori cercetătorii fără o bază în tehnologia programării fie că o ignoră, fie nu vor (și de multe ori nici nu sunt interesați) să o parcurgă, confuzia ce duce la auto-acreditarea într-un domeniu conexe este explicabilă.

Modul simplificator în care DT se referă la generativismul lingvistic, într-o logică binară care eludează faptul că în materie de modelare se lucrează cu grade de adecvare și relevanță, este însă simptomatic pentru viziunea sa limitativă în problema în discuție.

Crede DT că gramaticile lui Joshi, atât de importante în LC, puteau fi concepute fără să fi fost precedate de cele ale lui Chomsky? Da, Chomsky a fost tot timpul foarte controversat, dar fără stimulentele lui nu știu ce ne-am fi făcut, inclusiv în LC și în LM, în ciuda faptului că el nu s-a prea referit explicit nici la LC, nici la LM.

Modul „simplificator” incriminat mai sus se referă la fraza „soluții abstracte de modelare generativă de tip universal”. Având în vedere că în articolul [1] aceasta este singura referire la generativism, bănuiesc că domnul Academician Marcus a vrut să spună „succint”. Apoi, continuarea ce se referă la logica binară pe care o folosesc în interpretare și simptomele viziunii mele limitative asupra problemei discutate desigur sunt efecte stilistice nereușite, întrucât nu am abordat (și nici nu mă interesează în mod deosebit) subiectul pe care îl invocă domnul Academician. Pentru că tot am ajuns aici, țin să-i reamintesc domnului Academician Marcus că Noam Chomsky și-a revizuit complet punctul de vedere care a dominat aproape 15 ani lingvistica mondială. Într-adevăr Chomsky este un mare om de știință, chiar dacă foarte controversat, dar acest statut îi este conferit și de onestitatea cu care s-a detașat de creațiile sale anterioare ce i-au adus notorietatea, dovedite (unele chiar de el însuși) ca fiind depășite, propunând soluții și teorii noi.

Formalismul TAG al lui Joshi este într-adevăr unul foarte important în LC ca și HPSG, LFG, CG și alte câteva. Dar dintre formalismele de lingvistică computațională, TAG este cel mai departe de influența chomskyană. Dacă se poate face o asociere între TAG și vreo teorie generativistă de tip chomskyan aceasta este doar de natură antinomică. Am colaborat cu profesorul Aravind Joshi în 1991 la Institutul Lingvistic de la Universitatea Santa Cruz din California, am fost apoi invitatul său la Universitatea din Pennsylvania, invitație motivată printre altele și de o deosebită apreciere pe care o demonstrație alternativă a mea, mai scurtă și, considerată de profesorul Joshi, mai elegantă a unei teoreme a domniei sale referitoare la categoria de limbaje acoperite de LTAG. Cu acea ocazie, profesorul Joshi mi-a pus la dispoziție trei volume consistente de lucrări asupra TAG tratând foarte amănunțit motivațiile lingvistice, proprietățile computaționale și caracterizarea matematică. Aceste volume i le-am pus la dispoziție și domnului Academician Marcus. Profesorul Joshi a fost în 1997 invitatul profesorului Dan Cristea și al meu la Școala de Vară EUROLAN unde a susținut o serie de prelegeri de înaltă ținută științifică. Am evocat aceste lucruri pentru a-l lămurii pe domnul Academician Marcus că formalismul TAG și varianta sa mai nouă LTAG îmi sunt familiare și

prin urmare mă surprinde afirmația dânsului implicând o filiație între teoriile lui Joshi și Chomsky.

Faptul că gramaticile context free se află din nou, începând cu anii '80, în centrul atenției în LC nu spune ceva ?

Acest lucru este exact și ilustrează foarte bine ceea ce spuneam înainte: contextul computațional în care complexitatea algoritmică este primul mare judecător al adecvării unui model (inerent limitat, după cum arătam mai devreme) bazat pe o anumită teorie lingvistică. În anii de vârf ai lingvisticii matematice, și în cei de început ai lingvisticii computaționale, pornindu-se de la o conjectură a lui Chomsky (limbajele naturale nu sunt limbaje independente de context) demontată în anii '80 de Gerald Gazdar (autorul teoriei GPSG), cercetarea a fost orientată pe identificarea de formalisme lingvistice cât mai puternice, cu puterea generativă cât mai apropiată de cea a gramaticilor universale (echivalente deci cu mașina Turing). Formalismul ATN (Augmented Transition Networks) al lui William Woods de la BBN a fost timp de peste 10 ani suportul standard al majorității sistemelor de prelucrare a limbajului natural. Eu însumi am dezvoltat în anii 1984 și 1985 un mediu de programare lingvistică conținând un editor de gramatici ATN și un compilator ATN. Din punct de vedere formal ATN-ul este echivalent cu o mașină Turing și tocmai această putere formală prea mare l-a scos din competiția soluțiilor utile în lingvistica computațională. La sfârșitul anilor '80 obiectivul major al LC (valabil și astăzi) a devenit identificarea unui formalism de putere generativă cât mai mică dar care să acopere cât mai multe din problemele practice puse de prelucrarea automată a limbajului natural. Așa au revenit în actualitate gramaticile independente de context și s-au dezvoltat abordările lexicalizate. Cele din urmă au fost propuse tocmai pentru a rezolva, în cadrul scheletelor de gramatici independente de context, idiosincraziile limbajului natural cel mai adesea localizate la nivelul lexical. Mai mult, după anii '90, odată cu resurecția interesului față de abordările statistice, gramaticile regulate și automatele finite au căpătat o utilizare foarte largă.

LC are mai multe părți, mai multe orientări, mai multe niveluri de abstracție, care comportă criterii diferite de evaluare.

Este adevărat că actualmente în LC se regăsesc orientări, abordări sau motivații diferite. Dar indiferent de sorginte, ele se plasează (cel puțin declarativ) în contextul computațional prin raportarea la un mediu software de prelucrare. Considerând exemplul HPSG, probabil cea mai în vogă teorie lingvistică computațională actuală, atunci când Ivan Sag analizează sau argumentează adecvarea teoriei sale în descrierea formală a unei limbii naturale (așa cum a procedat în recente sale conferințe la Facultatea de Litere a Universității București și în Aula Academiei Române) el se plasează în sfera lingvisticii teoretice. Atunci când prezintă soluțiile de implementare a unui fragment major al limbii engleze și discută rezultatele generate de analizorul HPSG dezvoltat de grupul sau de la Universitatea Stanford și modalitățile algoritmice de rezolvare a ambiguităților (așa cum a făcut în prelegerea susținută la sediul RACAI, el se plasează în sfera LC.

DT îl asociază pe D. Hays la ideea sa privind falimentul LM și lansarea, drept consecință, a LC.

Afirmația de mai sus conține două lucruri false:

- a) nu am vorbit de falimentul LM ci de insuficiența **metodelor** sale la momentul invocat (cred că citatele pe care le-am prezentat și argumentele aduse până acum sunt lămuritoare).
- b) Eu nu-l pot asocia pe David Hays la o idee pe care nu am exprimat-o.

În textul meu original scriam: „la propunerea lui David Hays, domeniul de cercetare al limbajelor naturale, din perspectiva utilizării acestora în interacțiunea cu calculatoarele electronice, este individualizat sub numele de lingvistică computațională”.

Propunerea lui Hays venea în sprijinul identificării unui nume comun pentru diversele preocupări asupra limbajului din perspectiva implementării de sisteme automate de prelucrare. Traducerea automată, un domeniu care se dezvoltase distinct de celelalte preocupări în domeniul prelucrării automate a limbajului natural, căzuse în disgrație în urma raportului ALPAC (*Languages and machines: computers in translation and linguistics. A report by the Automatic Language Processing Advisory Committee, Division of Behavioral Sciences, National Academy of Sciences, National Research Council. Washington, D.C.: National Academy of Sciences, National Research Council, 1966. (Publication 1416.) 124pp.*). În raportul ALPAC, comandat în 1964 de Academiei Naționale de Științe, în afara criticilor deosebit de dure la adresa realizărilor și abordărilor de până atunci în domeniul traducerii automate existau și o mulțime de recomandări care se refereau *la noi metode de investigație științifică și la abordarea unor obiective mai realiste*. Istoria domeniului a reținut (pe nedrept) doar apriga critică a lui Bar-Hillel care, considerată unilateral, a dus la stoparea pentru circa 15 ani a cercetării oficiale în domeniul traducerii automate în SUA și mai apoi în majoritatea țărilor dezvoltate (o incitantă prezentare a a ceea ce a însemnat proiectul ALPAC este „ALPAC: the (in)famous report”, <http://ourworld.compuserve.com/homepages/WJHutchins/Alpac.htm>, și îi aparține lui John Hutchins). Ceva trebuia făcut pentru a conserva câștigurile științifice obținute până atunci și a permite în noul context continuarea cercetărilor anterioare cu scopul declarat al realizării de programe cu obiective realiste. O serie de minți luminate (John Pierce, David Hays, John Carroll) au văzut pericolul ca, asociate cu domeniul traducerii automate, toate celelalte preocupări privind prelucrarea automată a limbajului puteau fi periclitare, și în acest sens în raport s-a inserat un capitol distinct numit „Automatic language processing and computational linguistics” ce arăta beneficiile aduse de cercetarea în domeniul traducerii automate în domeniile prelucrării automate a limbajului și al lingvisticii computaționale. Printre altele în capitolul respectiv se arată că „... (what is required is) **basic developmental research in computer methods for handling language**, as tools for the linguistic scientist to use as a help to discover and state his generalizations, and ... to state in detail the complex kinds of theories..., **so that the theories can be checked in detail.**” (sublinierea mea, DT). Mai mult președintele comitetului de elaborare a raportului ALPAC, John Pierce, conștient de pericolul interpretării greșite sau al ignorării recomandărilor prezente în anexele raportului (așa cum s-a și întâmplat), a ținut să insereze în raportul final adresat președintelui Academiei Naționale de Științe o secțiune nouă care

sublinia idea de a susține lingvistica computațională în mod distinct de traducerea automată („supporting computational linguistics, as distinct from automatic language translation”). Dezvoltând ideile din capitolul raportului ALPAC referitor la prelucrarea limbajului natural (concept care și atunci și acum este diferit de cel al traducerii automate) Pierce considera că NSF (National Science Foundation) trebuia să asigure fonduri de cercetare pentru dezvoltarea de modele de limbă de dimensiuni mari „**since small-scale experiments and work with miniature models of language have proved seriously deceptive in the past, and one can come to grips with real problems only above a certain scale of grammar size, dictionary size, and available corpus**”.

Acesta este contextul în care David Hays, activ cercetător la începutul anilor '60 în domeniul traducerii automate (de altfel unul din membrii comitetului ce au elaborat raportul ALPAC) a propus individualizarea preocupărilor legate de prelucrarea limbajului natural cu ajutorul calculatorului, dezvoltarea de modele de limbă realiste (nu miniaturi la îndemâna cercetării individuale) și a aplicațiilor „serioase” (în opoziție cu experimentele la scară mică) sub numele de *lingvistică computațională*.

Denumirile folosite pentru preocupările la interferența limbajelor, informaticii și matematicii au variat tot timpul și nu cred ca acest aspect merită prea multă atenție. Lingvistică matematică? computațională? inginerească? algebrică? cognitivă? aplicată? cantitativă? teoretică? statistică? probleme matematice ale semioticii? tehnologia limbajului? limbajul în inteligența artificială? lingvistica inginerească? procesarea limbajului natural? “information storage retrieval”? lingvistica cibernetică? pe fiecare dintre acestea am întâlnit-o și propriile mele articole au fost publicate aproape sub fiecare dintre etichetele de mai sus.

Citatul de mai sus mi se pare extrem de relevant pentru discuția de față și definește clar diferența de opinii. Dacă de pildă distincția dintre *medicină umană și medicină veterinară* sau (coborând în taxonomie) între cardiologie și stomatologie „nu merită prea multă atenție” atunci domnul Academician are dreptate.

Din punctul meu de vedere însă, este o mare diferență între unele denumiri ale studiului limbii amintite mai sus (la care se mai poate adăuga o listă la fel de numeroasă), ele definind câteva domenii distincte definite prin propriile obiective, competențe, metode și modele.

În 1962 s-a înființat în USA “Association of Computational Linguistics”.

De fapt în 1962 s-a înființat AMTCL, acronim pentru „Association for Machine Translation and Computational Linguistics”, primul președinte al AMTCL fiind Victor Ingve (cel pe care l-am citat mai devreme), iar al doilea fiind David Hays. ACL (Association of Computational Linguistics) a apărut abia în 1968.

În 1963 Ferenc Kiefer a demarat la Budapesta revista “Computational Linguistics”, care a trăit peste zece ani.

Este adevărat, dar conținutul ei era foarte diferit de al revistei „Mechanical Translation and Computational Linguistics' apărută în 1965 ca revistă oficială a AMTCL.

Și tot ca un rezultat al diferențierilor tot mai mari care apăruseră în domeniu, AMTCL își încetează activitatea la începutul anilor '70 fiind înlocuită de „American Journal of Computational Linguistics” care în 1984 devine „Computational Linguistics” (actuala denumire).

Conferința de la Grenoble de “traitement automatique des langues” din 1967 era a treia de acest fel, fiind precedată de o alta, la New York, în 1965 și de una în Anglia, probabil în 1963, organizată de M. Masterman. Între timp, la ruși, numeroase conferințe au avut loc pe tema “avtomaticeskaja obrabotka tekstov” iar “Sprachkunde und Informationsverarbeitung” a fost uneori eticheta folosită de germani s.a.m.d. Nu negăm rolul important pe care l-a avut David G. Hays în dezvoltarea CL, dar acest rol a fost altul decât cel afirmat de DT.

Nu am să reiau explicația faptului că nu i-am atribuit lui Hays nici un rol demolator, dar trebuie să subliniez faptul că inițiativa lui David Hays, de care am discutat mai devreme, a avut un rol **fundamental** în evoluția CL. Așa cum am arătat mai sus, inițiativa disocierii de traducerea automată, pentru a nu periclita restul preocupărilor privind prelucrarea automată a limbajului a fost o necesitate conjuncturală. În 1965, când la New York a avut loc prima conferință COLING, Hays anticipa desigur efectul de bumerang al raportului la elaborarea căruia participa, și a propus chiar atunci, detașarea oficială prin sintagma „computational linguistics” de domeniul traducerii automate (pe care îl părăsise de altfel și Hays cel ce fusese unul dintre principalii specialiști în traducere automată ai RAND Corporation). Deci nu Hays a creat domeniul lingvisticii computaționale, el este cel ce a „oficiat” botezul. Și nu a făcut-o de pe orice poziție ci de pe cea de fost membru al Comisiei Alpac și de președinte al AMTCL.

Emergența LC s-a produs încă din anii “50, sintagma LC a devenit curentă încă de la începutul anilor “60. Șirul de conferințe COLING nu a făcut decât să continue aceasta tradiție. Alții au preferat folosirea sintagmei LM (a se vedea, de exemplu, “Prague Bulletin of Mathematical Linguistics”, “Prague Studies of Mathematical Linguistics”, revista japoneza “Mathematical Linguistics” (in echivalentul ei japonez) etc. În ceea ce privește însă profilul acestor reviste, nu am constatat o diferență față de cele de CL. Desigur, între timp au început să apară și unele publicații mai specializate, cu referire la părți determinate ale CL (cum ar fi cea relativă la corpusul lingvistic). Etichetele nu au avut importanța și nu știu sa se fi desfășurat vreo competiție între ele. Chiar Hays a folosit diverse etichete, de exemplu cea din [3].

Persistența cu care domnul Academician pune semnul egalității între domeniul lingvisticii matematice, în care fără discuție nu a avut sau nu are rival în România, și cel al lingvisticii computaționale sau tehnologia limbajului este aparent foarte curioasă. Nu și dacă observăm următoarele fapte:

- sintagma „lingvistică matematică” este din ce în ce mai puțin utilizată (o căutare pe internet a termenilor „mathematical linguistics”, „computational

linguistics”, „natural language processing” și „language technology” este foarte instructivă: numărul de documente ce îi referă este 4.630, 87.900, 169.000 și respectiv 2.840.000);

- în domeniul strict computațional, la care se referea [1], în România activează de câțva timp o serie de cercetători importanți (majoritatea dintre ei membrii ai Comisiei de Informatizare pentru Limba Română pe care am onoarea să o conduc, și din care de altfel face parte și domnul Academician Marcus);
- domnul Academician Marcus fie nu cunoaște, fie dezavuează rezultatele românești obținute în domeniul **prelucrării cu calculatorul** a limbii române (cel puțin așa poate fi considerată ignorarea completă a acestora în lucrările domniei sale); ori poate consideră că nu reprezintă domeniul său de interes.

Dar DT merge mai departe pe ideea sa și afirmă (în completă discordanță cu viziunea lui Hays, de la care se reclamă) că “metodele LM și-au atins limitele” (încă în urmă cu peste 30 de ani!), pentru ca numai două pagini după această afirmație (deci la pagina 135 din [1]) să afirme că e nevoie de “modele formale ale limbii la toate nivelurile ei (fonetică, morfologie, sintaxă, discurs) gramaticale formale [...]”. Cum vede DT aceste modele formale altfel decât sub forma logico-matematică?

Asupra primei părți a acestei fraze cred că am discutat suficient. Referitor la „contradicția” pe care o semnalează în partea a doua a frazei de mai sus, nu pot să-i recomand domnului Marcus decât să citească încă de câteva ori articolul respectiv (sau să-l citească integral). Este vorba de **NOI modele formale de limbă** (în opoziție cu cele vechi), resurse lingvistice computaționale adecvate momentului actual. Dintre noile teorii care au apărut și s-au și impus așa putea să amintesc teoria optimalității în comunicare dezvoltată de Prince and Smolensky în 1993 (cu implementări în domeniul fonologiei și morfologiei computaționale și cu promițătoare rezultate chiar în sintaxă), teoriile sintactice bazate pe unificare și satisfacerea de restricții, precum și o întreagă pleiadă de teorii ale discursului. În domeniul prelucrării automate a limbajului natural există standarde, există tehnologii specifice, există organizații mondiale specializate, mai toate apărute în ultimii 10-15 ani. Dacă domnul Academician Marcus poate afirma că pentru limba română în domeniul resurselor lingvistice computaționale s-a făcut (sau a făcut) ceva înainte de anii '90 înseamnă că domnia sa are o imagine complet diferită de a tuturor specialiștilor din lume.

Știe oare că multe modele de acest fel există de câteva decenii? Indicații bibliografice asupra lor sunt date parțial în [4], [5], [6], [7] iar pentru cercetările românești în [8], [9]. Desigur, aceste modele sunt inegale ca valoare, au nevoie de continuări, modificări, ameliorări, dar ele nu pot fi ignorate. Fonetica, fonologia, vocabularul, morfologia, sintaxa, semantica lingvistică și lingvistica istorică au beneficiat din plin de metodele matematice, așa cum se poate vedea din impactul deosebit al lucrărilor respective în literatura de specialitate;

Recursul la modelele anilor '60-70 descrise în lucrările menționate ca argument pentru concepte ce au apărut la începutul anilor '90 mă scutește de comentarii. Pe de altă parte, avansul științific în orice domeniu se clădește pe cunoașterea anterioară iar cazurile de „frângere cognitivă”, când salturile științifice neagă cunoașterea anterioară sunt rare și ele de regulă definesc revoluțiile în știință. Filiația sau influențele în dezvoltarea unui domeniu științific (atunci când ele pot fi depistate cu obiectivitate) constituie preocuparea istoricilor științei. Lucrările tehnice, de regulă se raportează la contemporaneitate, ceea ce în termeni temporali poate însemna, în funcție de dinamica domeniului, câțiva ani, un deceniu, mai multe decenii sau perioade chiar mai mari. De pildă, puține lucrări tehnice în domeniul lingvisticii teoretice, al fonologiei se referă la marele gânditor Panini, considerat de mulți oameni de știință creatorul științei limbii. Lucrarea sa fundamentală *Astaka*, cunoscută și sub numele de „gramatica lui Panini” conține descrieri formale ale regulilor de producție ale limbii sanscrite și o clasificare cu peste 1700 de elemente constitutive ale limbajului. Aceste elemente sunt organizate în clase a căror agregare este descrisă prin intermediul unor reguli ordonate, într-o manieră apropiată de teoriile actuale. El poate fi considerat un precursor al teoriei limbajelor formale și al lingvisticii matematice, dar puține cărți sau lucrări de referință în aceste domenii menționează numele genialului savant ce a trăit cu mai bine de peste 2500 de ani în urmă. În schimb, numele său se regăsește în orice lucrare serioasă de istorie a lingvisticii formale.

Obstinația cu care domnul Academician Marcus încearcă să sugereze că eu aș dezavua metodele matematice, sau rezultatele importante ale lingvisticii românești dovedește că domnia sa complet neinformată în ceea ce mă privește.

DT indică, drept domeniu al LM, numai “aspectul sintactic, gramatical”, despre celelalte nu a aflat. Nu a aflat nici ca LM a abordat și aspecte analitice, nu numai pe cele generative.

Fals: „**numai**” este imaginația domnului Academician. Citatul corect este: „în **general** al aspectului lor sintactic, gramatical”.

DT definește “dimensiunea fundamentală” a LC prin “fezabilitatea instanțierii unei descrieri lingvistice cât mai complete, mentenabilitatea acestei instanțieri și, desigur, conformanța cu realitatea uzului limbii”. ([1]: 133). Cu un mic efort intelegem despre ce este vorba. Desigur că problemele de complexitate, de cost, nu puteau fi încă abordate în anii '50 și '60 cu mijloacele cu care ele au început a fi studiate în a doua jumătate a anilor '70, când instrumentele elaborate în informatica matematică deveniseră mult mai perfecționate. Dar acest fapt nu ține, cum crede DT, de alegerea între LM și LC, ci de progresul general realizat în știință. Pentru a mă referi la propria noastră experiență, atunci când, în 1969, prezentam la COLING-ul din Suedia gramaticile contextuale nu aveam cum să mă ocup de aspectul complexității acestor gramatici în maniera în care s-a putut face acest lucru ulterior (a se vedea, de exemplu, [10]). Dar acest fapt nu are nici o legătură cu eticheta folosită.

Efortul (chiar mic) este probabil generat de unii termeni de specialitate nefamiliarii domnului Academician. Voi furniza lămuririle necesare mai jos.

Eu mă refer la perioada actuală când invoc ca dimensiune fundamentală *fezabilitatea instanțierii* unei descrieri lingvistice cât mai complete. Instanțierea unei descrieri lingvistice înseamnă altceva decât complexitatea formală, de care de altfel și amintesc în secțiunea trunchiată a citatului folosit de domnul Academician Marcus mai sus. Este un termen tehnic care se referă la construcția propriu-zisă, în baza unui formalism sau teorii lingvistice, a unei gramatici și a dicționarului aferent, care furnizate ca resurse unui program de prelucrare a limbajului natural, permit acestuia să analizeze sau să genereze un text arbitrar. O astfel de instanțiere este fezabilă dacă ea se poate realiza în condiții de timp și resurse umane rezonabile.

Nu m-am mirat atunci când "Encyclopedia of Microcomputers" și "Encyclopedia of Computer Science and Technology" mi-au solicitat o contribuție cu tema "Semiotics and Formal Artificial Languages" (a se vedea [11]) și nici când "Handbook of Formal Languages" mi-a solicitat un capitol privind "Contextual Grammars and Natural Languages"[12] iar o lucrare preponderent teoretică a fost inserată în "Computational Linguistics in the Netherlands 2000"[13].

Nu văd rostul acestor lămuriri. Toată lumea îl știe, îl recunoaște și nimeni dintre cercetătorii adevărați nu-l contestă pe omul de știință Marcus, important reprezentant român al lingvisticii matematice, creatorul acestei școli în România. În articolul [1] nu m-am referit nici direct nici indirect la domnia sa. Faptul că am evocat criticile pe care le-am comentat anterior la adresa **metodelor** lingvisticii matematice ale începutului deceniului șapte nu are nici o legătură cu realizările (încă o dată, excepționale) ale domnului profesor. Însă probabil că identificându-se cu LM mondială, domnia sa a considerat critica asupra **metodelor LM** din anii '60 un atac la persoana sa, adevărat act de blasfemie.

În anii din urmă, domnul Academician încearcă să transfere în contextul noilor tendințe și tehnologii ale limbajului, ignorând o realitate existentă, tot portofoliul de rezultate pe care le-a obținut anterior creditându-le ca surse primare a tot ceea ce se întâmplă azi în tehnologia limbajului în România (și nu numai). Și cine nu este de acord cu acest lucru (parafrazându-l pe domnul Marcus) trebuie demonizat. Textul pe care îl comentez ca și acțiunile recente declanșate de domnul Academician Marcus, pretinse a fi iscate de conținutul articolului [1], nu fac decât să-mi întărească această impresie. Eu nu am nimic de împărțit cu domnul Academician.

Nu m-am mirat nici când am văzut că o revistă cu titlul "Linguistics and Philosophy" publică articole excelente de LC. Interferențele merg în toate direcțiile și ele caracterizează cultura contemporană. În acest orizont trebuie să ne plasăm, cred, atunci când ne referim la disciplinele cognitive care se dezvoltă sub ochii noștri și își pun amprenta pe modul nostru de gândire și de comportare. Un tratat ca "Mathematical Methods in Linguistics" [14] include multe fapte de LC, deși în titlul sau nu figurează epitetul "computational". O revistă ca "Theoretical Linguistics" (1970

2000), publicata de Walter de Gruyter (Berlin–New York) a inclus multe articole vizând aspecte matematice și/sau computaționale, deși numele revistei nu indică acest lucru. Chiar o revista mai tradițională, ca “Linguistics” a inclus de multe ori articole de LM și nici “Foundations of Language” nu a procedat altfel. Multe fapte de LM și de LC se plasează în mod natural în orizontul semioticii computaționale.

Faptul că tratatul amintit nu incorporează în titlu atributul computațional nu mă surprinde, pentru că ar fi creat o confuzie pe care autorii au evitat-o deliberat. Cartea respectivă nu este o carte de lingvistică computațională, conținutul ei tratează exact ce anunță în titlu: metode matematice folosite în studiul lingvistic. Lingvistica teoretică, puternic formalizată în ultimele decenii apelează inevitabil (ca de altfel marea majoritate a domeniilor științifice) la metode și modele matematice.

Era internetului impune desigur o problematică nouă, față de care abordările anterioare se pot dovedi insuficiente.

Exact aceasta este esența celor 3 paragrafe din [1] incriminate și combătute pe larg de domnul Academician Marcus: insuficiența abordărilor anterioare. Conștientizarea acestei insuficiențe însă a precedat cu câțiva ani apariția internetului.

Salutăm inițiativa noii generații de cercetători de a se dedica noilor probleme.

Nu putem ignora tonul paternalist privind noua generație de cercetători care se dedică problemelor ridicate de internet în prelucrarea automată a limbajului natural. INTERNET-ul este o revoluție! Și implicațiile sale sunt atât de mari încât asigurarea accesului universal la Internet a devenit o problemă fundamentală chiar și pentru o organizație de calibrul UNESCO. Am avut onoarea să fac parte din Comisia de Experți creată de Secretarul General al UNESCO (comisie de cel mai înalt nivel) pentru elaborarea documentului Recommendation on Multilingualism and Universal Access to Cyberspace. Sunt al doilea expert român (după dl. Ambasador Dan Hăulică, Membru Corespondent al Academiei) care a făcut parte dintr-o comisie de experți UNESCO de acest nivel.

Ignorarea în cercetarea privind prelucrarea automată a limbajului natural a fenomenului INTERNET este de neconceput. Societatea cunoașterii are ca una din premisele sale fundamentale accesul universal, neîngrădit de bariere lingvistice la cunoșterea stocată în internet. Alte comentarii sunt de prisos.

Dar trecerea de la ieri la azi și de la azi la mâine nu poate fi decât una care ține seama în mod critic de experiența acumulată.

Nimeni nu neagă acest lucru, și faptul că l-am rugat insistent pe domnul Academician să facă parte din Comisia de Informatizare pentru Limba Română cred că arată buna mea credință și speranța pe care o nutream (și care mai supraviețuiește încă) că experiența domniei sale va fi pusă în slujba obiectivelor pe care nici eu nici domnul Marcus nu le putem atinge singuri. În același spirit, i-am propus domnului Academician Marcus să scriem împreună o antologie a cercetărilor românești în domeniul lingvisticii formale și computaționale, de la începuturile pe care le evocă domnia sa și pînă în zilele noastre. Din păcate propunerea a rămas fără răspuns.

Din tot ceea ce am prezentat mai sus rezulta clar ca LM si LC au fost mereu împreună și că, în general, etichetele nu au contat prea mult. Unii au mers chiar mai departe; astfel, în capitolul 4, "Mathematical and Computational Linguistics", din [15], se afirma pur și simplu (p.86): "Mathematical linguistics has also been called theoretical linguistics and even computational linguistics". Iar mai departe, în același loc: "Computational Linguistics originated around 1950 with the initiation of research on automatic translation" (se trimite la o carte editată de D.G.Hays [3] și la o alta avându-l ca autor pe acesta [16]).

Nu văd în pasajul pe care l-am citat mai sus nici un argument împotriva a ceea ce am susținut în [1] și în cele prezentate aici. Notez în treacăt adverbul „even” cu o valoare discursivă în completă consonanță cu considerentele istorice pe care le-am invocat ale evoluției științifice și tehnologice în domeniul prelucrării limbajului natural.

În România, minți luminate ale anilor "60, ca profesorii Al. Rosetti, Grigore Moisil și Tudor Vianu, au înțeles schimbările care se profilau și au sprijinit proiectul înființării unei secțiuni de "lingvistica aplicată" la Facultatea de Limba și Literatură Română a Universității din București, dar s-au găsit alții care să-i torpileze.

Așa este, și mă bucură elogiul adus acestor corifei ai științei românești. Poate și pentru că alături de câțiva reprezentanți importanți ai lingvisticii românești actuale care au înțeles tendințele și imperatiivele momentului (Prof. Dan Mazilu-decanul Facultății de Litere, Prof. Alexandra Cornilescu, Conf. Emil Ionescu) am participat la reluarea acestei lucrări. Programul de Masterat în Lingvistică Formală și Computațională de la Facultatea de Litere a Universității din București, funcționează de mai bine de 2 ani și nutresc speranța că Ministerul Educației și Cercetării va aproba demersurile noastre privind chiar înființarea unui departament cu acest profil.

În același sens, am participat alături de profesorul Cristea (având fără discuție și sprijinul altor minți luminate ale Universității A.I.Cuza din Iași) la lansarea în 2001 a Masterat-ului în Lingvistică Computațională al Facultății de Informatică. Nu este ușor să pendulezi între Iași și București, dar și domnul profesor Cristea, și doamna profesor Cornilescu și eu o facem pentru ca cele două programe „surori” de master să-și împlinească menirea de a pregăti câți mai mulți specialiști în folosul programelor de informatizare pentru limba română.

La Academia Română a funcționat mulți ani "Comisia de Lingvistică Matematică" iar revista "Cahiers de Linguistique Theorique et Appliquee", înființată în 1962, a fost multă vreme expresia colaborării lingvisticii cu matematica și cu informatica. In ciuda forțelor adverse, s-a reușit în acei ani atragerea unor studenți străluciți ai unor facultăți umaniste la cercetarea limbii și literaturii cu mijloacele LM și LC.

Comisia de Informatizare pentru Limba Română de la Academia Română, înființată în anul 2001, încearcă, ținând cont de realitățile și prioritățile actuale, să armonizeze eforturile celor ce lucrează în domeniul limbii române și care cred în

perspectiva înrolării ei în cadrul limbilor importante ale societății cunoșterii. Eu am convingerea că voi putea spune peste timp același lucru: „*In ciuda forțelor adverse, s-a reușit în acei ani atragerea unor studenți străluciți ai unor facultăți umaniste la cercetarea limbii și literaturii cu mijloacele*” tehnologiei limbajului.

Pentru a da numai două exemple de actuali profesori universitari care au susținut teze de doctorat de acest tip, voi menționa pe Pia Brinzeu, de la Catedra de Engleză a Universității din Timișoara și pe Mihai Dinu, de la Facultatea de Litere a Universității din București. Tot în acea perioadă și-a susținut teza de doctorat Sorin Cristian Niță, pe o temă de critică textuală automată privind înlănțuirea (filiația) diferitelor variante ale “Istoriei Țării Românești” (Șerban Cantacuzino).

Exemple de profesori și cercetători români valoroși, cu contribuții substanțiale în domeniul limbii române se pot da foarte multe. Mulți dintre ei sunt în străinătate și fac o bună propagandă științei românești. Mi-e cunoscută cartea cu adevărat remarcabilă a domnului profesor Mihai Dinu „Personalitatea limbii române”, de altfel premiată de Academia Română. Această lucrare este o solidă cercetare de lingvistică computațională în spiritul actual tocmai pentru că a parcurs acea cale dificilă a instanțierii lingvistice (în cazul său la nivelul componentului lexical).

Iată însă că, în pofida realităților puse în evidență mai sus, în ([1]: 134) se scrie: “În Romania, cercetările în domeniul LC și al prelucrării limbajului natural, precum și primele rezultate practice au apărut la începutul anilor “80 [3, 4, 5, 6]”.

La ce trimit numerele indicate în paranteze ? La o bibliografie de 24 de titluri în care aproape toate (dar toate cele indicate între paranteze) încep cu DT (ignorându-se regula generală în lumea științifică, a așezării numelor autorilor aceluiași articol în ordine alfabetică; dar nu acest fapt este cel care ne interesează în momentul de față).

Înainte de a comenta acest pasaj și pe cel următor, nu pot să trec peste observația absurdă și falsă pusă între parantezele ce trădează totuși o ezitare a probității omului de știință în fața unei răutăți gratuite. Nu există nici o regulă generală de genul celei afirmate. Ordonarea alfabetică este o convenție între autorii cu contribuții egale în redactarea unei lucrări. Am deschis la întâmplare două volume de specialitate, conținând contribuții (S. Armstrong et al. (eds) „Natural Language Processing Using Very Large Corpora, Kluwer, 1999 și T. Strzalkovski (ed) „Natural Language Information Retrieval”, Kluwer, 1999). Din cele 19 lucrări cu mai mulți autori, doar trei urmăresc (probabil din întâmplare) regula generală în lumea științifică pe care o invocă domnul Academician și pe care probabil a impus-o și o impune tuturor celor alături de care publică, indiferent de contribuția fiecăruia.

Să observăm că încă în 1978, în articolul “Mathematical and Computational Linguistics” [9] de prezentare a activității din România în domeniul LM și LC se face referire la peste 400 de articole publicate de 130 de autori români și sunt menționați peste 300 de autori străini (unii

dintre ei, nume de vază ale LM și LC din acea perioada) care au citat și continuat cercetările românești. Să mai adăugăm că numeroși lingviști români dintre cei mai importanți au citat și folosit rezultatele școlii românești de LM și LC. Iată ca vine acum DT și face (deliberat sau nu) din tot acest efort un teren viran care-l aștepta pe DT să tragă primele jaloane. Nu e cam mult?

Deși am repetat de nenumărate ori până în acest moment, o mai fac o dată, precizând că discuția din [1] se referea la **resurse lingvistice computaționale și programe software de dialog în limbaj natural** (în limba română). Acestea erau rezultatele practice pe care le menționeam în citatul comentat cu gratuită aciditate. Poate să-mi menționeze domnul Academician vreun sistem de dialog în limba română implementat înaintea sistemelor pe care le-am realizat eu și colaboratorii mei? Iată câteva repere:

- Sistemul QA (1980) un sistem inferențial de întrebare răspuns în limba română, susținut de un demonstrator original de teoreme în calculul predicatelor de ordin 1;
- SDLR (1981) un sistem de dialog în limba română ce a extins capacitățile lui QA cu operatorii lingvistici ai logicii fuzzy;
- IURES (1983) sistem de generare automată a sistemelor de întrebare răspuns, independent de limbă, pe care l-am realizat împreună cu Dan Cristea, acum decanul facultății de informatică a Universității Cuza. Sistemul IURES a fost omologat internațional în 1988 și a constituit primul produs de inteligență artificială exportat (în același an). Sistemele IURES și SDLR sunt referite printre altele în enciclopedia de lingvistică computațională. Mai important este faptul că sistemele IURES și SDLR sunt amplu descrise în prestigioasa antologie “The Survey of the Current Status Research and Future Trends in Machine Translation and Natural Language Processing” realizat în 1992 de JEIDA (Japan Electronic Industry Development Association), fiind de altfel singurele sisteme de dialog în limbaj natural din întreaga zonă fost comunistă incluse în această carte.

Acestea erau referințele incriminate de domnul Academician și dacă domnia sa poate să-mi indice un singur sistem de prelucrare a limbajului natural realizat în România înaintea celor pe care le-am citat, eu am greșit. Dar mă îndoiesc. Nu cunosc conținutul articolului menționat (pe care i l-am solicitat de altfel domnului Academician, fără a-l primi însă), astfel încât nu pot afirma nimic despre cei 130 de autori români ce au realizat (conform afirmației domnului Marcus) lucrări de lingvistică computațională. Ce pot însă să afirm este că am citit multe din lucrările de lingvistică teoretică contemporană ale marilor noștri lingviști și ele au fost extrem de relevante ca material factual în cercetările mele. Dar lucrările pe care le-am citit (și citat) eu, nu erau din domeniul lingvisticii computaționale. Lucrările domnului Marcus (în special cele din domeniul *limbajelor formale*) apăreau destul de frecvent între referințele bibliografice ale lucrărilor mele de la începutul anilor '80. Eram la început de drum, sursele documentare erau puține și demersul era natural. Pe

atunci, Chomsky era din nou foarte în vogă, noua sa teorie *Government and Binding* impulsionând o serie de cercetări în domeniul formalizării gramaticii universale. Tentația computațională față de această teorie a fost enormă, și chiar dacă actualmente nu există nici o gramatică computațională efectivă a GB, idei fundamentale din GB se regăsesc în formalisme lingvistice computaționale moderne (cum ar fi HPSG).

Să fim bine înțeleși. Nu noi avem nevoie de încă o citare pe lângă miile de citări deja acumulate, ci noile generații de studenți și de cercetători au dreptul la o informare corectă asupra dezvoltării LM și LC în general și, în particular, asupra LM și LC în România. DT a mai publicat, în urma cu câțiva ani, un articol în care se schița o privire istorică asupra LC în România, cu câteva citări la întâmplare, care tradau necunoașterea situației reale.

Cu rezerve față de prima parte a paragrafului, mă opresc la grija domnului Academician pentru dreptul noilor generații de studenți și de cercetători asupra „informării corecte” asupra istoriei LM și LC. Personal, cred că mult mai important pentru ei este să știe prezentul și tendințele viitoare ale domeniului. Astfel de cunoștințe le pot asigura un loc de muncă, o direcție de specializare, o carieră viitoare. Noile generații de studenți și de cercetători sunt utilizatori pasionați ai Internetului. Acest uriaș ocean informațional le asigură un imens volum de cunoștințe, începând cu cursuri on-line (obligatorii pentru profesori la mai toate universitățile importante ale lumii), valome ale conferințelor sau articole extrem de utile, recente și mai puțin recente, cărți electronice. Chiar și relevante lucrări de istorie asupra diverselor domenii științifice. Sistemele moderne de regăsire documentară le asigură și o ierarhizare a acestor surse de informare în raport cu relevanța și cu interesul manifestat de alți cititori. Listele de discuții sau arhivele de întrebări frecvente (FAQ) le pot oferi răspunsuri avizate și obiective la întrebările ce-i preocupă. În anexă este furnizat un exemplu.

În ultima parte a citatului de mai sus, domnul Academician Marcus aduce în discuție o lucrare a mea din 1996 și care arată că frustrările domniei sale sunt mai vechi. Articolul de care amintește domnul Academician mai sus, are titlul **„Resurse lingvistice computaționale: trecut, prezent și viitor”** și a apărut în volumul **„Limbaj și Tehnologie”**, Ed. Academiei, 1996. Cei interesați, pot găsi articolul respectiv în pagina oficială a RACAI (<http://www.racai.ro> secțiunea publicații). Iar cele „câteva citări la întâmplare, care trădau necunoașterea situației reale” apar în capitolul 2. **„Cercetări și realizări românești în domeniul prelucrării automate a limbajului natural”**. Cred că titlul volumului, al articolului și al capitolului sunt lămuritoare pentru ceea ce discutăm acolo, dar probabil fraza, care trimitea la un volum editat de domnul Marcus, „abordările statistice, revenite acum în actualitate, au avut o tradiție strălucită (în România, adăugarea mea DT)” a fost prea scurtă și insuficient de laudativă.

Mai este un aspect care cere o precizare. În conformitate cu specificul volumului în care apare articolul [1], DT face numeroase referiri la acte și documente ale unor organisme europene și internaționale, cum este și firesc, pentru a nu mai vorbi de aspectul financiar al colaborării cu

organismele respective.

Aceasta situație a existat de la începutul LM și LC (chiar dacă nu a avut amploarea de azi), datorită faptului că LM și LC au apărut și ca urmare a unor comandamente sociale, privind precaritatea mijloacelor de prelucrare a informației. Imi amintesc de faimoasele Rapoarte CETIS care veneau de la EURATOM, Bruxelles, pe teme legate de analiza și prelucrarea automată a limbajului, traducere automată și documentare automată. În USA, diferite corporații (cum ar fi RAND Corporation, Santa Monica, Calif.) finanțau cercetări similare. O întâlnire semnificativă a fost aceea din 1962, organizată de "NATO Advanced Summer Institute", la Veneția, Italia, privind traducerea automată. De numele acestui Institut este legat un document care a marcat evoluția cercetărilor de traducere automată: seria de expuneri prezentate de Y. Bar-Hillel [17]. În legătură cu aceste activități dirijate și finanțate de diferite organisme europene și internaționale, trebuie să observăm că cei implicați au avut înțelepciunea și priceperea necesare pentru a nu reduce proiectele respective la dimensiunea lor exclusiv utilitară, ci de a o subordona pe aceasta unei perspective mai ample, care lua în considerare orizontul științific real al problemelor. Pentru a da un prim exemplu, mă voi referi la faptul că mai multe rapoarte CETIS au pus în discuție un concept care, născut din experimentele de traducere automată, avea să se dovedească de o deosebită semnificație pentru teoria sintactică în toată generalitatea sa; este vorba de conceptul de proiectivitate sintactică, cu consecințe bogate în studiul structurilor arborescente și al gramaticilor de dependență. Azi putem spune că și sintaxa limbajului natural și teoria matematică a grafurilor au profitat esențial de conceptul respectiv (folosit până și de Rene Thom, în probleme de morfogeneză [17]). Această expansiune a unui concept sau rezultat dincolo de motivația sa inițială este testul cel mai convingător al interesului său. Un al doilea exemplu se referă la titlul provocator folosit de Bar-Hillel pentru expunerile sale: "Patru conferințe despre lingvistica algebrică și traducerea automată".

Simpla alăturare a celor două sintagme, una foarte teoretică, cealaltă aparent tehnologică, avea menirea să-i avertizeze pe cei care presau să se obțină cât mai repede rezultate practice asupra faptului că proiectele de traducere automată nu se pot finaliza de azi pe mâine, ci au nevoie de un lung itinerar lingvistic, matematic și computațional. Acum știm că acest itinerar continuă și azi, cu tatonări și reveniri, și, chiar dacă nu a dus încă la rezultatele visate, a impulsionat în mod esențial cercetările de AI, cu consecințe benefice pentru aspectele logice și semantice ale limbajului natural.

Întrebarea pe care ne-o punem, dar o lăsăm deocamdată fără răspuns, deoarece nu suntem pregătiți pentru a-l da, este următoarea: Nu cumva aspectele pe care le-am criticat mai sus sunt consecința unui fenomen mai general, acela al unui orizont insuficient de cuprinzător, al unei prea mari

dependențe de factori utilitari imediați? Știința a oscilat mereu între cognitiv și utilitar, dar istoria arată că funcția utilitară s-a manifestat în toată profunzimea ei atunci când ea a fost fructul unei evoluții firești a funcției cognitive, evoluție care poate fi de doi ani, de 20 de ani, de 200 sau de 2000 de ani. Cu un ochi îndreptat spre comisiile europene, suntem obligați totuși să ținem treaz și celălalt ochi, îndreptat spre ceea ce se întâmplă pe scena cercetării științifice vii, așa cum apare ea în revistele de specialitate și la întâlnirile științifice de profil.

Remarcile de mai sus îmi sugerează celebra fabulă cu strugurii cei acri. Cercetarea instituționalizată (în opoziție cu cea „de dragul artei”) are motivații întotdeauna justificabile. Organismele de finanțare a cercetării, naționale sau internaționale, nu fac desigur acte de caritate. Obținerea unei finanțări pentru un proiect de cercetare nu este la îndemâna oricui și el implică nu numai abordarea unei probleme importante, dar și credibilitatea grupului de cercetare. Evaluarea propunerilor de proiecte se face de către experți recunoscuți în domeniul respectiv, angajați și plătiți de agențiile de finanțare a cercetării. În condițiile unei concurențe internaționale acerbe pentru fondurile (din păcate prea mici) destinate cercetării, a lua în derâdere, invocând caracterul utilitar, cercetările ce obțin concurențial finanțarea arată o desprindere de realitate. În luna martie a.c. am participat la evaluarea propunerilor de proiecte europene din cadrul Programului Cadru 5 (apelul 8), și în calitate de raportor al direcției „II.1.1 - **Exploratory High Risk/Long Term Research**”, pot să afirm că propunerile de proiecte pe care le-am văzut erau foarte departe de a avea caracter utilitar. Domnul Academician Marcus lasă fără răspuns o întrebare cu răspuns sugerat, ridicând o problemă discutată cu ceva timp în urmă, anume a tipului de cunoaștere contemporană: enciclopedică (și inerent generalistă) sau specializată. Cel puțin în domeniile tehnologice, viteza fără precedent a apariției de cunoștințe noi face imposibilă cunoașterea enciclopedică și în același timp expertă pe toată lărgimea spectrului cunoașterii actuale chiar și într-un domeniu aparent îngust. Tehnologia limbajului este actualmente termenul ce subsumă toate preocupările legate de **prelucrarea automată a limbajului natural**. Cred că acest lucru spune totul!

3. In loc de concluzii

Ajungând în acest punct al răspunsului meu la atacul domnului Academician Marcus mărturisesc că mă încearcă un apăsător sentiment al deșertăciunii. Nu am dorit această polemică și în nici un caz în acest context. Considerând că ea este nepotrivită față de obiectivele urmărite de proiectul „SI-SC: Soluții și strategii în România”, în calitatea mea de director de proiect și coeditor al volumului de față, am discutat cu membrii comitetului director al proiectului oportunitatea publicării polemicii domnului Academician Marcus (și implicit a răspunsului meu) în volumul destinat unor probleme tehnice. Părerea a fost unanimă că nu este cazul să amestecăm obiectivele proiectului cu discuția de față. Dar transmițând domnului Academician această opinie și făcându-i propunerea de a găzdui această polemică pe internet (în pagina oficială a RACAI) domnia sa s-a simțit cenzurat,

insultat și îndreptățit să facă o serie de afirmații pe care mă abțin să le comentez. Decizia de includere a acestei secțiuni în volumul de față am luat-o fără plăcere pentru că pe de o parte, în ciuda părerii domnului Academician Marcus (*Articolul meu se încadrează perfect în obiectivul pe care pretindeți că-l urmăriți și în acest spirit a fost conceput. Realizați gravitatea deciziei Dv?* - de a nu-l include în volum, precizarea mea, D.T.) continui să cred că nici articolul domniei sale nici al meu nu își aveau rostul aici. Pe de altă parte, nu pot decât să deplâng supărarea pe care i-am provocat-o fără voie domnului Marcus și risipa de energie pe care o depune într-o problemă care din punctul meu de vedere nu există. Drept care sperând că includerea articolului ce *se încadrează perfect în obiectivul...* îi va da domnului Academician satisfacția pe care și-a dorit-o, las cititorii să aprecieze cât de grav ar fi fost pentru obiectivul tehnologiei limbii române în contextul „Societatea Informațională – Societatea Cunoașterii: Soluții și strategii în România” ca cele două articole să nu fi apărut aici.

Referințe bibliografice (secțiune din lucrarea domnului Academician Marcus):

- [1] D. Tufis. *Promovarea limbii române în SI-SC*. In *Societatea Informațională – Societatea cunoașterii* (coord. F. Gh. Filip). Ed. Expert, București, 2001, 131–142.
- [2] D. G. Hays. *The field and scope of computational linguistics*. Papers in Computational Linguistics (eds. F. Papp, G. Szepe). Proceedings of the Third International Meeting of Computational Linguistics, held in Debrecen, Hungary, 1971. Akademiai Kiado, Budapest, 1976, 21–26.
- [3] D. G. Hays (ed.). *Readings in Automatic Language Processing*, American Elsevier, New York, 1967.
- [4] S. Marcus. *Mathematical Linguistics in Europe. Current Trends in Linguistics* (Th. A. Sebeok, ed.), vol.9, Mouton, The Hague, 1972, 646–687.
- [5] S. Marcus. *Mathématique et Linguistique*. In *Mathématique, Informatique et Sciences Humaines*, Paris, 26, 1988, 103, 7–21.
- [6] S. Marcus. *The status of research in the field of analytical algebraic models of language*. In *Current Issues in Mathematical Linguistics* (C. Martin-Vide, ed.). Elsevier–North Holland, Amsterdam, 1994, 3–21.
- [7] S. Marcus. *Lingvistica matematica, azi*. In *Matematica în lumea de azi și de maine* (C. Iacob, coord.), Editura Academiei, București, 1985, 182–186.
- [8] S. Marcus. *Recent Romanian investigations in the field of mathematical and computational linguistics*. Avtomatizeskaja Obrabotka Tekstov, Matem. Fyz. Fakulta, KL Praha, 1973, 15–42.
- [9] S. Marcus. *Mathematical and computational linguistics*. In *Current Trends in Romanian Linguistics* (A. Rosetti, S. Golopenția Eretescu, eds.). Revue Roumaine de Linguistique 23, 1978, 1–4, 559–588.
- [10] S. Marcus, C. Martin-Vide, G. Paun. *Contextual grammars as generative models of*

- natural languages*. Computational Linguistics 24, 1998, 2, 245–274.
- [11] S. Marcus. *Semiotics and formal artificial languages*. In *Encyclopedia of Computer Science and Technology* (A. Kent, J.C. Williams, eds.) 29, Ed. Marcel Dekker, New York, 1994, 393–405; also in *Encyclopedia of Microcomputers* (A. Kent, J.C. Williams, eds.) 15, 1995, 299–312.
- [12] S. Marcus. *Contextual grammars and natural languages*. *Handbook of Formal Languages* (G. Rozenberg, A. Salomaa, eds.), 2, Springer, Berlin, New York, 1997, 215–235.
- [13] S. Marcus, C. Martin–Vide, G. Paun. *A new–old class of linguistically motivated regulated grammars*. *Computational Linguistics in the Netherlands 2000* (W. Daelemans et al., eds.), Selected Papers from the Eleventh CLIN Meeting, Ed. Rodopi, Amsterdam, New York, 2001, 111–125.
- [14] B. H. Partee, A. Ter Meulen, R. Wall. *Mathematical Methods in Linguistics*. Kluwer, Dordrecht, 1990.
- [15] E. F. Beckenbach, Ch. B. Tompkins (eds.). *Concepts of Communication: Interpersonal, Intrapersonal and Mathematical*. John Wiley and Sons, New York, 1976.
- [16] D. G. Hays. *Introduction to Computational Linguistics*. American Elsevier, New York, 1967.
- [17] R. Thom. *Stabilité Structurelle et Morphogenèse*. John Benjamins, New York, 1970.
- [18] Y. Bar–Hillel. *Four Lectures on Algebraic Linguistics and Machine Translation* revised version of a series of lectures given in July 1962, before a NATO Advanced Summer Institute, Venezia, Italy.

ANEXA1: Exemple de căutare într-o arhivă de întrebări frecvente (Usenet FAQ)

The image displays two screenshots of a Microsoft Internet Explorer browser window, showing search results for queries in the Usenet FAQ Archives. The browser's address bar shows the URL <http://www.faqs.org/cgi-bin/faq/faqsearch>.

Top Screenshot: Search for "mathematical linguistics"

The search results page is titled "The Usenet FAQ Archives" and "Results for query 'mathematical linguistics'". It displays a message: "Sorry... NO Matches were found for query 'mathematical linguistics'". Below this, a note states: "Since **nothing** was found... you might want to be a bit less specific in your search please. Just a thought...". A link is provided: "Back To FAQs - Usenet References - FAQ Search Engine". The copyright notice reads: "© Copyright The Internet FAQ Consortium, 1998. All rights reserved."

Bottom Screenshot: Search for "computational linguistics"

The search results page is titled "The Usenet FAQ Archives" and "Results for query 'computational linguistics'". It displays a message: "File name (modification date), and list of matched lines". Below this, a list of results is shown:

1. [mail/college_email/part1](#). (Apr 29 2002)
 - Computational Linguistics: lcl.cmu.edu
2. [ai.faq/general/part6](#). (Apr 5 2002)
 - The Association for Computational Linguistics homepage:
3. [ai.faq/general/part5](#). (Apr 5 2002)
 - The Association for Computational Linguistics (ACL) has a Data

Both screenshots show a search bar on the left with the query "mathematical linguistics" and a list of web pages found, including links to "Buy Current Issues in Mathematical Linguistics by Carlos Martin-Vide at Barnes & Noble.com", "Click here to buy Introduction to mathematical linguistics at Amazon.com", "GSI MC - Research Group on Mathematical Linguistics", "Linguistics 640 Home Page", "Institute of Formal and Applied Linguistics", "University of Sussex - Natural Language Processing", "Oxford Encyclopedia of Linguistics", "U. Penn. Linguistics and Associated Faculty", "LINGUIST List 11.326: Historical Linguistics, Mathematical Linguistics", and "CSE477 Lingu519".

ANEXA 2: Definiții

What is Mathematical Linguistics?

MATHEMATICAL LINGUISTICS is the study of mathematical structures and methods that are of importance to linguistics. As in other branches of applied mathematics, the influence of the empirical subject matter is somewhat indirect: theorems are often proved more for their inherent mathematical value than for their applicability.

Both in phonology/morphology and in syntax/semantics the choice of linguistic formalism is to some extent influenced by considerations that go beyond the primary issue of descriptive adequacy. One important issue is Recognition Complexity. This concerns the complexity of the decision problem for membership in a language: it is assumed that a grammatical theory should have the property of guaranteeing that there is some reasonably rapid (polynomial in the length of the input) computation that will answer the question of whether a given sequence of words is a grammatical expression according to a given grammar. Human beings certainly do much more than this when they listen to an utterance and figure out the meaning of what was said, so a grammatical theory that cannot even guarantee reasonably rapid confirmation of well-formedness is probably not psycholinguistically realistic. Another one is Learnability, which concerns what sorts of mathematically definable procedures could in principle correctly guess the grammars for languages.

(Geoffrey K. Pullum and Andras Kornai)

What is Computational Linguistics?

Simply put, COMPUTATIONAL LINGUISTICS is the scientific study of language from a computational perspective. Computational linguists are interested in providing computational models of various kinds of linguistic phenomena. These models may be "knowledge-based" ("hand-crafted") or "data-driven" ("statistical" or "empirical"). Work in computational linguistics is in some cases motivated from a scientific perspective in that one is trying to provide a computational explanation for a particular linguistic or psycholinguistic phenomenon; and in other cases the motivation may be more purely technological in that one wants to provide a working component of a speech or natural language system. Indeed, the work of computational linguists is incorporated into many working systems today, including speech recognition systems, text-to-speech synthesizers, automated voice response systems, web search engines, text editors, language instruction materials, to name just a few.

(Copyright © 2000, The Association for Computational Linguistics)



București, România
Licența Ministerului Culturii nr. 1442/1992
Tel.: 411.60.75; Fax: 411.54.86
Consilier editorial: **Valeriu IOAN-FRANC**

ISBN 973-8177 - -