

RAPORT ȘTIINȚIFIC proiect complex ReTeRom, etapa I - iunie 2018

Proiectul 1: COBILIRO

Activitatea A1.1

Denumire Etapă: **Studiu *state-of-the-art* asupra realizării corpusurilor bimodale**

Autori: UAIC, ICIA, UTCN, UPB

REZUMATUL ETAPEI

Prima etapă a proiectului 1 COBILIRO pregătește cu date și modele realizarea prototipului, ce va fi realizat pe parcursul etapei a II-a. În esență, activitățile prevăzute în acest an au ca principal obiectiv stabilirea convențiilor de adnotare lingvistică și adnotarea unui eșantion de texte, ce vor fi apoi luate drept model pentru un proces automat, propunerea unui inventar de date geografice și specificații pentru realizarea interfețelor utilizator. Alte activități, urmăresc realizarea comunicării în cadrul consorțiului, diseminarea rezultatelor, coordonare și management și raportare.

Rezultate Etapa: Raport asupra progreselor privind realizarea corpusurilor modale

În această etapă au fost identificate, analizate și evaluate studiile centrate pe realizarea corpusurilor bimodale în vederea selectării acelor care pot constitui modele necesare producerii unei tehnologii integrate pentru procesarea limbajului natural în limba română, procesarea și adnotarea pe diferite niveluri lingvistice a corpusului bimodal generat în cadrul Proiectului 1 COBILIRO.

1. INTRODUCERE

1.1 Mediu, modalitate, corpus

În semiotică, prin modalitate¹ se înțelege modul prin care informațiile trebuie codificate pentru a fi prezentate oamenilor, adică tipul de semn și statutul realității atribuite sau revendicate de un semn, fie el text sau alt gen. Modalitățile senzoriale vor fi: vizuale, auditive, tactile, olfactive, gustative, kinestezice etc.

O listă a tipurilor de semne ar include: scrierea, simbolul, indexul, imaginea, harta, graficul, diagrama etc. Unele combinații de semne pot fi multi-modale, adică diferite tipuri de semne pot apărea grupate împreună pentru un efect amplificat. Trebuie făcută distincția dintre mediu și modalitate. Modalitatea se referă la un anumit tip de informație și / sau la formatul de reprezentare în care sunt stocate informațiile. Mediul este mijlocul prin care această informație este transmisă simțurilor interpretului uman.

Limbajul natural este modalitatea principală, având proprietăți invariabile în mediul auditiv - ca limbă vorbită, în mediul vizual - ca limbă scrisă, în mediul tactil - ca Braille și în mediul cinetic - ca limbaj al semnelor. Din acest punct de vedere, termenul de bi-modal (care înseamnă utilizarea combinată a două modalități) are un înlocuitor preferat, care este

¹ [https://en.wikipedia.org/wiki/Modality_\(semiotics\)](https://en.wikipedia.org/wiki/Modality_(semiotics))

bi-medial (care înseamnă, utilizarea combinată a două medii). Cele două medii de comunicare ale limbajului, în cazul proiectului nostru, sunt textual și auditiv.

În contextul interacțiunii om-calculator și în interesul proiectului nostru vom considera că **modalitatea** este clasificarea unui canal, independent de tipul de intrare / ieșire senzorială între om și calculator. Un sistem este proiectat unimodal dacă are o singură modalitate implementată și multimodal dacă are mai multe (Karray *et al.*, 2015). În cazul în care sunt disponibile mai multe modalități pentru o sarcină, se spune că sistemul are modalități redundante. Modalități multiple pot fi utilizate în combinație pentru a oferi metode complementare care, deși redundante, pot transmite informații mai eficiente (Palanque *et al.*, 2001, p. 43).

Modalitățile pot fi implementate în ambele direcții: dintre om spre calculator și invers. În aceste transferuri informaționale se utilizează o gamă largă de tehnologii: (a) modalități comune (vizuale - grafică computerizată, de obicei, printr-un ecran; auditive - ieșiri audio; tactile - vibrații sau alte mișcări), (b) modalități sofisticate (gustative; olfactive; termoreceptive; nonperceptive; de echilibru). Modalitățile vizuale și auditive sunt cel mai frecvent utilizate deoarece sunt capabile să transmită informații cu o viteză mai mare față de altele: de la 250 la 300 cuvinte pe minut (WPM) (Ziefle, 1998) și, respectiv, de la 150 la 160 cuvinte pe minut (Williams, 1998). Deși implementarea în direcția de modalitate calculator => om este rară, tactilul poate atinge o medie de 125 WPM prin utilizarea unui afișaj Braille actualizabil. Alte forme comune ale tactilului sunt vibrațiile de smartphone și controlerul jocului.

În cazul comunicării om => calculator, calculatoarele pot fi echipate cu diferite tipuri de dispozitive de intrare și senzori pentru a le permite să primească informații de la oameni. Dispozitivele de intrare comune sunt adesea interschimbabile dacă au o metodă de comunicare standardizată cu computerul și permit ajustări practice utilizatorului. Anumite modalități pot oferi o interacțiune mai bogată în funcție de context, iar opțiunile de implementare permit sisteme mai robuste (Bainbridge, 2004: 483). Vorbim de modalități simple (tastatură, dispozitiv de indicare, ecran tactil) și modalități complexe (recunoașterea imaginilor, a vorbirii, interfețe haptice, orientare etc.).

Odată cu creșterea popularității smartphone-urilor, modalitățile complexe devin tot mai pe plac publicului larg. Recunoașterea vorbirii, după introducerea lui Siri (Epstein, 2015), a reprezentat un punct de vânzare important al produselor Apple. Această tehnologie oferă utilizatorilor un mod alternativ de a comunica cu calculatoarele atunci când tastarea este mai puțin dezirabilă. Cu toate acestea, într-un mediu cu zgomot puternic, modalitatea auditivă nu este destul de eficientă. Acest lucru justifică faptul că anumite puncte forte ale modalităților diferă în funcție de situație (Kurkovsky, 2009: 210-211). Alte modalități complexe, cum ar fi cele legate de vizual (de exemplu, Kinect Microsoft, sau alte tehnologii similare) pot face ca sarcinile sofisticate să fie mai ușor de comunicat unui calculator, mai ales sub forma unei mișcări tridimensionale (Kurosu, 2013: 366).

Desigur, un sistem cu multiple modalități oferă mai multe posibilități utilizatorilor și robustețe sistemului. De asemenea, permite o mai mare accesibilitate utilizatorilor. De reținut rămâne faptul că modalitățile multiple pot fi folosite ca rezervă atunci când anumite forme de comunicare nu sunt posibile. Acest lucru este valabil mai ales în cazul

modalităților redundante în care două sau mai multe modalități sunt utilizate pentru a comunica aceleași informații.

Există șase tipuri de cooperare între modalități și ele ajută la definirea modului în care o combinație sau fuziune a modalităților cooperează pentru a transmite mai eficient informațiile (Grifoni, 2009: 37): (1) echivalența (informațiile sunt prezentate în mai multe moduri și pot fi interpretate ca aceleași informații), (2) specializarea (când un anumit tip de informații este prelucrat întotdeauna prin aceeași modalitate), (3) redundanța (mai multe modalități procesează aceleași informații), (4) complementaritatea (multiplele modalități iau informații separate și le combină), (5) transferul (o modalitate consumă o altă modalitate pentru a produce noi informații), (6) concurența (mai multe modalități iau informații separate care nu sunt fuzionate).

În accepțiunea proiectului ReTeRom, prin corpus bimodal² vom înțelege o colecție de înregistrări orale însoțite de transcrierile lor și de metadatele corespunzătoare. Un corpus bimodal este găzduit pe o platformă specializată, împreună cu serviciile și aplicațiile web de acces, dezvoltare și întreținere ale lui, unde sunt specificați algoritmi de utilizare ai corpusului și, în unele cazuri, de unde pot fi descărcate exemple de aplicații care utilizează corpusul.

Un corpus bimodal este un caz particular de corpus multimodal, care, la rândul lui, reprezintă un tip particular de corpus³. Corpusurile pot conține texte scrise, înregistrări orale sau ambele modalități de redare a unei limbi naturale. În proiectul de față ne interesează ultimul caz.

Se remarcă o evoluție a conceptului de corpus: prima generație de corpusuri conținea doar transcrieri (ex. *British National Corpus*), a doua generație conținea și textul oral și pe cel scris (ex. *The Michigan Corpus of Academic Spoken English*). A treia generație conține în plus alinierea dintre componenta scrisă și cea orală. Există și o ultimă etapă de evoluție, caracterizată prin includerea dimensiunii video în corpusuri (Rasso and Mello, 2014:29). Obiectul interesului nostru este a treia generație de corpusuri.

Uneori⁴ se face distincția dintre corpus de vorbire (engl. *speech corpus*) și corpus oral (engl. *spoken corpus*), deși, adesea, cel puțin în engleză, termenii sunt utilizați neglijent, fără o distincție clară între ei. În continuarea acestei secțiuni vom arăta cum pot fi departajate aceste două categorii de corpusuri, pentru ca mai departe să ne referim la amândouă prin termenul de corpus oral.

Prin **corpus de vorbire**, în general, se înțelege o bază de date de fișiere audio și de transcrieri textuale ale acestora, într-un format care poate fi folosit pentru a crea modele acustice ce devin motorul cercetărilor asupra recunoașterii vorbirii. Un exemplu sunt transcrierile Switchboard revizuite la Institutul pentru Procesarea Semnalelor și Informațiilor (ISIP) începând cu anul 1997⁵ (Godfrey and Hollman, 1997). Fișierele audio și transcrierile lor pot fi aliniate (la nivel de fonem, silabă și cuvânt). În cadrul sistemelor de recunoaștere vocală, modelele prozodice sunt folosite în principal pentru predicția evenimentelor

² V. și <http://www.aclweb.org/anthology/R09-1044>

³ Uneori este utilizat termenul de *corpus lingvistic*, evident - un pleonasm, pentru că prin corpus se înțelege oricum o colecție de date despre limbă, în format scris sau audio.

⁴ <http://www.ilc.cnr.it/EAGLES96/spokentx/spokentx.html>

⁵ <https://catalog.ldc.upenn.edu/ldc97s62>

prozodice (accente sintactice și semantice) asociate unui text. Implementarea prozodiei în procesarea vocală a permis obținerea de semnal vocal sintetizat cu conținut retoric/pragmatic. Corpusul de vorbire interesează cercetările în care accentul se pune pe semnalul sonor, pe proprietățile sale acustice și proprietățile articulatorii ale tractului vocal. Reprezentarea simbolică în acest caz se face utilizând alfabetul fonetic.

În schimb, **corpusul oral** prezintă interes, cu precădere, pentru cercetările în care se studiază modul de utilizare a limbii și caracteristicile diverselor niveluri lingvistice: lexical, morfologic, sintactic, semantic, analiza discursului, a conversației, studii de pragmatica comunicării, sociolingvistică, dialectologie etc., dar și în tehnologia vorbirii, atât pentru sisteme de înțelegere cât și de generare a vorbirii. Transcrierea marchează, într-un sistem ortografic îmbogățit, diverse fenomene auditive care însoțesc pronunția: bâlbâieli, pauze, tuse, răs etc., uneori dublat și de un alfabet fonetic. Și aici înregistrarea verbală este aliniată la interpretarea sa textuală. Ideal ar fi atât cercetătorii lingviști cât și cei din tehnologia limbajului natural să folosească același set de date, cu condiția ca acesta să îndeplinească cerințele ambelor comunități în ceea ce privește colectarea, transcrierea, codificarea și adnotarea datelor.

1.2 Tipuri de corpusuri orale

Așa cum am precizat mai sus, ne referim aici la baze de date de înregistrări vocale dublate de transcrierile lor textuale, utile în crearea de *modele acustice*, care, la rândul lor sunt folosite în aplicații de recunoaștere a vorbirii (*speech recognition engines* sau *speech-to-text systems*, S2T). Pe de o parte, corpusul oral înseamnă un corpus de conversație informală, improvizată, fără implicare mediatică (Chafe, 1995), iar pe de altă parte, el este folosit pentru a înțelege limba a cărei formă de manifestare predominantă este cea orală.

Corpusurile orale sunt clasificate în două tipuri:

1) **voce-în-citire** (*read speech*), care include:

- lecturi din cărți;
- știri;
- liste de cuvinte;
- secvențe de numere.

2) **vorbire spontană** (*spontaneous speech*), care include⁶:

- dialoguri - între două sau mai multe persoane;
- narrative - o persoană redă o întâmplare sau spune o poveste;
- trasee pe hartă - o persoană explică alteia un traseu pe o hartă, persoanele nu se văd, ambele au în față aceeași hartă, dar doar harta celei care explică are marcat un traseu, cealaltă urmând să refacă traseul pe harta proprie;
- stabilirea unei întâlniri - două persoane încearcă să fixeze parametrii unei întâlniri, în funcție de restricțiile individuale din calendare;
- simulări “Vrăjitorul din Oz” - simularea unei situații din viața reală, ca de exemplu, dialogul pentru rezervarea unui zbor între un client și un funcționar al unei agenții de turism;

⁶ A se vedea, ca exemplu, corpusul SRI:

https://web.stanford.edu/dept/linguistics/corpora/material/X_Speech_Corpora.pdf

Limba engleză, este cel mai bine reprezentată de corpusuri orale, dar pe web pot fi găsite și corpusuri ale altor limbi: poloneză, arabă, bulgară, cantoneză, cehă, farsi, franceză, germană, hindi, japoneză, coreeană, mandarină, portugheză, rusă, spaniolă, tamil, vietnameză etc.

Avantajele utilizării unui corpus oral sunt:

1. economie de timp la antrenarea sistemelor, nu este necesară colectarea spontană, înregistrările sunt pre-procesate, pot fi utilizate în mod repetat;
2. corpusurile pot conține cantități mari de date, necesare în procesele de învățare;
3. sunt interogabile;
4. pot fi create sub-corpusuri pe bază de selecții;
5. exemplele de limbaj din corpusurile de vorbire spontană sună natural, pentru că în colectarea datelor se pune accent pe spontaneitate.

Dificultăți/provocări în utilizarea unui corpus oral:

1. uneori, calitatea înregistrării este mai scăzută decât într-un laborator fonetic;
2. sunt greu de procurat, procesul de introducere de metadate și de prelucrare este laborios, consumator de timp;
3. dacă se lucrează cu voluntari, mai ales într-un regim *crowdsourcing*, datele trebuie verificate, validate, curățate, pentru că, uneori, voluntarii sunt neglijenți ori chiar rău intenționați;
4. înregistrările trebuie asigurate prin contracte de drepturi de autor și de permisiune de utilizare (*data privacy*) cu vorbitorii.

1.3 Transcrierile

Transcrierile sunt redările textuale ale înregistrărilor vocale. În realizarea corpusurilor, transcrierile pot urma sau precede înregistrările vocale, dar ele sunt componente obligatorii ale corpusurilor bimodale. Tipurile de transcrieri posibile într-un corpus oral sunt:

- transcrierea ortografică fără aliniere în timp cu semnalul sonor;
- transcrierea ortografică aliniată în timp;
- transcriere fonetică realizată cu ajutorul simbolurilor, la rândul lor clasificate în semne grafice care au drept corespondente sunete „primare” și semne grafice care au drept corespondente sunete marcate de unul sau mai multe fenomene fonetice (Bejinariu *et al.*, 2007).

2. ACHIZIȚIA CORPUSURILOR BIMODALE

Cercetările dedicate recunoașterii automate a vorbirii au cunoscut o creștere spectaculoasă începând cu anii '60 (Halle și Stevens 1962; Denes & Mathews 1960; Denes, 1960), deși utilizarea de corpusuri orale ca depozitare a limbii ce ar trebui interpretată sau generată de mașină este de dată mult mai recentă. În vederea achiziționării de astfel de corpusuri au început să fie create interfețe specializate. Am considerat util să prezentăm în această secțiune caracteristicile unei aplicații dedicate achiziției corpusurilor bimodale, care

concentrează principalele funcționalități necesare construcției de corpusuri aliniate voce-text de tip voce-în-citire, independente de limbă.

Google for Android (Hughes *et al.*, 2010)⁷

Arhitectura sistemului este una *client-server*, cu clienți rulând o aplicație ecran sub Google Android, serverul rulând Java și tratând cereri HTTP simultane de la clienți.

Aplicația server pune la dispoziție clienților rânduri de text, care sunt scurte segmente de text (propoziții, cuvinte). Aceste prompt-uri sunt generate la inițializarea aplicației din cereri de căutare Google Search, filtrate pentru a elimina cuvinte pornografice și cuvinte tipărite greșit; limba cererii se detectează automat, inclusiv utilizând informații de localizare; la primirea unei cereri de date (prompt) de la un client, serverul selectează la întâmplare din această listă de texte un rând și îl trimite clientului (fără protecția de a nu trimite duplicate); el primește apoi înregistrarea sonoră și metadatele aferente ei via HTTP POST request; tripletul text-metadate-înregistrare este apoi stocat într-o bază de date și clientul este informat că înregistrarea poate fi ștearsă din cardul SD; înregistrărilor sonore li se adaugă automat o margine scurtă de sunet de fundal la început și la sfârșit.

În *aplicația client*, la deschiderea unei sesiuni de lucru, clientului i se deschide o interfață în care i se cere mai întâi să completeze metadatele următoare:

- mediul acustic (în casă, afară, în mașină, nivelul de zgomot);
- genul vorbitorului;
- vârsta (pe decade) a vorbitorului;
- accentul vorbitorului;
- un nume utilizator, pentru identificarea sesiunii.

În completare la acestea, aplicația mai include automat:

- data și ora/minutul curent;
- versiunea hard a telefonului și versiunea de SO Android;
- numărul IMEI al telefonului (un ID unic de identificare a dispozitivului mobil);
- locația geografică, determinată de GPS-ul mobilului (atenție, o parte din aceste informații pot ridica probleme în contextul curent al legii intimității informatice!).

Urmează o serie de interacțiuni, în care utilizatorul primește în ecran textul și trebuie să-l pronunțe. Vocea utilizatorului este stocată și trimisă pe server împreună cu metadatele completate de client. Pentru a permite utilizarea în locuri fără conexiune la Internet, clientul necesită doar conexiuni intermitente la server pentru a descărca și încărca înregistrările. Sesiunea este terminată de utilizator la dorință sau după un număr fixat de intrări completate de utilizator.

Sistemul a fost utilizat până acum pentru a colecta peste 3000 de ore de transcriere audio în 17 limbi, inclusiv limba rusă.

3. O TRECERE ÎN REVISTĂ A UNOR CORPUSURI BIMODALE

Oferim în această secțiune o seamă de informații relative la corpusuri bimodale existente în lume. Informațiile sunt grupate în 10 categorii, care au fost identificate în majoritatea descrierilor.

⁷ <https://static.googleusercontent.com/media/research.google.com/ro//pubs/archive/36801.pdf>

Santa Barbara Corpus of Spoken American English⁸ (Du Bois *et al.*, 2000-2005)

1. **limba:** engleză;
2. **tip corpus:** în principal vorbire spontană dar și voce-în-citire, predomină conversația față în față, dar conține și convorbiri telefonice, discuții în timpul jocului de cărți, al gătitului, discuții la serviciu, prelegeri la clasă, predici, spunere de povești, întâlniri profesionale, prezentări din timpul ghidajelor turistice, texte citite etc.;
3. **dimensiune:** înregistrări totalizând 249.000 cuvinte;
4. **vorbitori:** oameni din diverse regiuni geografice de pe întreg teritoriul SUA, de vârste diferite, ocupații diferite, de sexe diferite, etnii și categorii sociale diferite;
5. **informații private relative la vorbitori:** numele participanților la conversații și alte date personale sunt anonimizate;
6. **metadate:** pentru fiecare vorbitor, sunt specificate: ID nume, sex, vârstă, localitatea de origine, statul de origine, statul în care trăiește la momentul înregistrării, educație, ani de educație, ocupație, etnie;
7. **conținut și mod de realizare:** înregistrare sonoră, transcriere, aliniere voce-timp-text (engl. *timestamps*) la nivelul unităților intonaționale;
8. **formatul înregistrărilor:** MP3 și WAV;
9. **suport pentru realizare:** -
10. **distribuție și copyright:** cu licență CC-BY-ND.

The Buckeye Speech Corpus⁹ (Pitt *et al.*, 2007)

1. **limba:** engleză
2. **tip corpus:** vorbire spontană în dialog, temele conversațiilor sunt cotidiene (politică, sport, transport, școală);
3. **dimensiune:** 40 de ore și 300.000 de cuvinte;
4. **vorbitori:** 40, intervievați fiecare câte o oră, provenind din Columbus, Ohio (născuți aici sau care au ajuns aici la o vârstă mai mică de 10 ani, ceea ce face din acest corpus unul dialectal); stratificare pe vârstă și sex;
5. **informații private relative la vorbitori:** -
6. **metadate:** -
7. **conținut și mod de realizare:** înregistrările sunt transcrise (folosind Soundsciber); aliniere automată (folosind Entropics Aligner) la nivel de cuvânt; ulterior, s-a corectat manual această aliniere;
8. **formatul înregistrărilor:** -
9. **suport pentru realizare:** -
10. **distribuție și copyright:** accesibil gratuit pentru cercetare (pe bază de cont); se poate descărca prin semnarea unei licențe.

The Spoken Turkish Corpus¹⁰ (Ruhi, 2011)

1. **limba:** turcă;
2. **tip corpus:** multimodal (conține și înregistrări video, nu doar audio); conversații față în față sau interacțiune mediată (la telefon);

⁸ <http://www.linguistics.ucsb.edu/research/santa-barbara-corpus>

⁹ <https://buckeyecorpus.osu.edu/>

¹⁰ <https://std.metu.edu.tr/en/>

3. **dimensiune:** 10 milioane de cuvinte (în dezvoltare);
4. **vorbitori:** reprezentând limba turcă contemporană, în varianta sa standard, dar și forme dialectale;
5. **informații private relative la vorbitori:** vorbitorii sunt anonimizati;
6. **metadate:** -
7. **conținut și mod de realizare:** textul este adnotat morfologic și la acte de vorbire; se folosește EXMARaLDA pentru transcriere și interogare;
8. **formatul înregistrărilor:** -
9. **suport pentru realizare:** realizat prin finanțare, dar și prin voluntariat;
10. **distribuție și copyright:** parțial accesibil gratuit pentru cercetare (demo), descărcare - prin semnarea unei licențe.

CORP-ORAL (Freitas and Santos, 2008)

1. **limba:** portugheza europeană;
2. **tip corpus:** dialoguri față în față, vorbire spontană în sensul că nu se impune o temă, dar vorbitorii știu că sunt înregistrați;
3. **dimensiune:** 53h, dintre care transcrise ortografic: 32h, dintre care transcrise fonetic: 1h;
4. **vorbitori:** 60, din zona Lisabonei, cu vârste între 17-74 ani, nivel educațional: de la 9 clase la doctorat;
5. **informații private relative la vorbitori:** -
6. **metadate:** -
7. **conținut și mod de realizare:** transcriere ortografică cu ELAN (cu marcarea pauzelor, respirației, râsului etc.), transcrierea fonetică cu Praat, adnotare prozodică;
8. **formatul înregistrărilor:** -
9. **suport pentru realizare:** -
10. **distribuție și copyright:** poate fi interogat la <http://spock.iltec.pt/>

Spoken Portuguese Corpus

1. **limba:** portugheza;
2. **tip corpus:** comunicare orală spontană, pe subiecte legate de viața curentă;
3. **dimensiune:** 86 de înregistrări: Portugalia (30), Brazilia (20), Angola (5), Cap Verde (5), Guinea-Bissau (5), Mozambic (5), Sao Tome and Principe (5), Macao (5), Goa (3), East-Timor (3). Total: 8h 44min, 153.588 tokeni;
4. **vorbitori:** diversitate sociolingvistică, perioada: 1970-2001;
5. **informații private relative la vorbitori:** -
6. **metadate:** -
7. **conținut și mod de realizare:** transcrieri aliniate în format XML Exmaralda și transcrieri în text simplu, adnotate la POS;
8. **formatul înregistrărilor:** fișiere audio WAV;
9. **suport pentru realizare:** -
10. **distribuție și copyright:** ELRA, gratuit pentru cercetare.

Bavarian Archive for Speech Signals Corpora (*Siemens Synthesis Corpus*)¹¹

1. **limba:** germană

¹¹ <http://www.bas.uni-muenchen.de/forschung/Bas/BasKorporaeng.html>

2. **tip corpusuri:** voce-în-citire (dictări, numere de la 1 la 100, nume de străzi, coduri postale, nume de orașe, numere de telefon etc.) și vorbire spontană (dialoguri ale unui dispecer de taxi vorbitor de germană și un client vorbitor de engleză înregistrat prin conexiune telefonică fixă și GSM, în traducere; monologuri spontane pe linii de telefon publice, dialoguri om-mașină purtate pe motocicletă în mers etc.), unele cu zgomot de fond;
3. **dimensiune:** peste 500.000 fraze
4. **vorbitori:** peste 10.000, din toate categoriile de vârstă și sociale, inclusiv adolescenți (13-20 ani), inclusiv vorbitori profesioniști, inclusiv vorbitori nenativi de germană, din toate regiunile Germaniei, inclusiv înregistrări istorice ale unor vorbitori de germană saxonă din Transilvania din perioada anilor 1970¹² etc.;
5. **informații private relative la vorbitori:** -
6. **metadate:** -
7. **conținut și mod de realizare:** parțial adnotat zgomotul de fond, parțial adnotarea prozodiei, segmentări automate la nivel de cuvânt și fonem, segmentare MAUS, segmentare în concordanță cu Standardul CLIPS - BPF, TextGrid, Emu;
8. **formatul înregistrărilor:** BAS Partitur Format files, compatibilitate Verbmobil, SpeechDat Database Format etc.
9. **suport pentru realizare:** Siemens, PhonDat¹³, proiectul Verbmobil, Erlanger Bahnansage¹⁴, SPINA¹⁵ (Robot Commands) etc.
10. **distribuție și copyright:** unele sub ELRA, altele direct la producător, contact: bas@bas.uni-muenchen.de.

GermaParl¹⁶ (Blätte and Blessing, 2018)

1. **limba:** germană;
2. **tip corpus:** protocoale plenare din Bundestag;
3. **dimensiune:** 100 de mil. de tokeni;
4. **vorbitori:** parlamentari, perioada 1996-2002;
5. **informații private relative la vorbitori:** numele vorbitorilor, afilierea la un grup parlamentar sau la un partid politic;
6. **metadate:** numele vorbitorului, afilierea la un grup parlamentar sau la un partid politic, teme de discuții;
7. **conținut și mod de realizare:** se specifică dacă enunțul este o interjecție sau vorbire;
8. **formatul înregistrărilor:** TEI sau pachet de date R;
9. **suport pentru realizare:** -
10. **distribuție și copyright:** domeniul politic nu ridică probleme de copyright.

C-ORAL-ROM¹⁷ (Nascimento *et al.*, 2002)

1. **limbi:** italiană, franceză, portugheză, spaniolă (patru sub-corpusuri comparabile);
2. **tip corpus:** vorbire spontană (vorbirea formală: conversații, inclusiv telefonice, media, conferințe, dezbateri politice, predare; vorbire informală: conversații, monologuri și dialoguri publice și private, interviuri, întâlniri, experiențe profesionale; amintiri personale;

¹² <http://www.bas.uni-muenchen.de/forschung/Bas/BasASDeng.html>

¹³ <http://www.bas.uni-muenchen.de/forschung/Bas/BasPD1eng.html>

¹⁴ <http://www.bas.uni-muenchen.de/forschung/Bas/BasERBAeng.html>

¹⁵ <http://www.bas.uni-muenchen.de/forschung/Bas/BasSPINAeng.html>

¹⁶ <http://www.lrec-conf.org/proceedings/lrec2018/pdf/1024.pdf>

¹⁷ <http://www.elda.fr/en/proj/coralrom.html>

transcrieri de filme reprezentative; eșantioane de înregistrări provenite din înregistrări de radio și TV; interacțiuni om-mașină; dialoguri extrase dintr-un sistem automat de rezervări de bilete de tren etc.);

3. **dimensiune:** 772 de texte din 121:43:07 de ore de înregistrări audio, aproximativ 300.000 de cuvinte pentru fiecare limbă => total: 1.200.000 de cuvinte;
4. **vorbitori:** 1,427, având caracteristici diferite de vârstă, nivel de educație, origini sociale și geografice;
5. **informații private relative la vorbitori:** -
6. **metadate:** fiecare sesiune conține metadate cu informații despre vorbitori, situația înregistrării, calitatea acustică a înregistrării, sursa, informații despre subiectul sesiunii;
7. **conținut și mod de realizare:** sunt evidențiate principalele evenimente paralingvistice și nonlingvistice, pauzele prozodice și evenimentele verbale, identificarea acestora din urmă fiind o caracteristică a acestui corpus; textele sunt lematizate și adnotate cu părți de vorbire; fiecare text transcris a trecut prin cinci stadii de validare: transcriere, revizuire, adnotare prozodică, aliniere; fiecare etapă a fost făcută de către alt lingvist;
8. **formatul înregistrărilor:** pentru sincronizarea text-sunet s-a folosit formatul WinPitch, cele patru corpusuri au fost transcrise conform standardelor ortografice de transcriere CHAT format (MacWhinney, 1994); sunetul este aliniat la text; marcaje XML;
9. **suport pentru realizare:** -
10. **distribuție și copyright:** corpusul este disponibil în două formate: 8 dvd-uri care conțin toate materialele sau 1 dvd care conține o versiune criptată a corpusului; are acordul fiecărui vorbitor.

Corpus de Référence du Français Parlé (CRFP) (Teston-Bonnard and Véronis, 2004)

1. **limba:** franceză;
2. **tip corpus:** interviuri grupate pe vârstă, educație, tipul conversației (privată, profesională sau publică);
3. **dimensiune:** 137 de înregistrări (37 de ore), 440.000 de cuvinte;
4. **vorbitori:** diferite categorii de vârstă și nivel de educație;
5. **informații private relative la vorbitori:** -
6. **metadate:** -
7. **conținut și mod de realizare:** -
8. **formatul înregistrărilor:** -
9. **suport pentru realizare:** -
10. **distribuție și copyright:** -.

Nagoya University conversation corpus (Fujimura *et al.* 2012)

1. **limba:** japoneză
2. **tip corpus:** conversații spontane între prieteni, membri ai unei familii sau colegi;
3. **dimensiune:** 129 de fișiere, durata: 30-60 min; 1,5 milioane de morfeme;
4. **vorbitori:** 198 de persoane de vârste diverse, nivel de educație variat, la fiecare conversație participă 2-4 persoane; distribuție neuniformă pe sexe (predomină cel feminin) și pregătire profesională înaltă (mulți absolvenți filologi);
5. **informații private relative la vorbitori:** -
6. **metadate:** -
7. **conținut și mod de realizare:** o parte din corpus (97 de dialoguri) a fost adnotată manual cu acte de vorbire și *sympathy* (vorbitorul arată interes pentru subiectul

conversației); transcriere ortografică (nu există accent, intonație, proeminență; există doar intonația urcătoare din întrebări);

8. **formatul înregistrărilor:** -
9. **suport pentru realizare:** -
10. **distribuție și copyright:** disponibil gratuit.

JAIST Annotated Corpus (Shirai and Fukuoka, 2018¹⁸)

1. **limba:** japoneză
2. **tip corpus:** conversații spontane între două persoane despre diverse subiecte;
3. **dimensiune:** 97 de dialoguri; durata: 100h;
4. **vorbitori:** la fiecare conversație participă 2 persoane; distribuție neuniformă;
5. **informații private relative la vorbitori:** -
6. **metadate:** -
7. **conținut și mod de realizare:** o parte din corpus (97 de dialoguri) a fost adnotată manual cu acte de vorbire (act de dialog) și *sympathy* (vorbitorul arată interes pentru subiect);
8. **formatul înregistrărilor:** -
9. **suport pentru realizare:** -
10. **distribuție și copyright:** disponibil gratuit.

4. STANDARDIZAREA CORPUSURILOR BIMODALE

Descriem în această secțiune inițiative de standardizare care pot fi aplicate, cel puțin parțial, în proiectul nostru.

Li and Yin (2006; 2007) descriu o inițiativă chineză de standardizare, care a fost aplicată la constituirea corpusului RASC863. Corpusul cuprinde un număr de 2200 de fraze echilibrate fonetic și 460 de fraze frecvent folosite în viața zilnică, pronunțate de 80 de vorbitori, aleși astfel încât reprezentativitatea pe categorii de vârstă, gen, studii să fie echilibrată. Vorbitorii reprezintă 10 pronunții dialectale diferite, culese de pe toată întinderea Chinei. Înregistrările vocale sunt dublate de transcrieri. Corpusul e utilizat la antrenarea sistemelor text-vorbire. Lucrarea prezintă și proceduri standard de producere a corpusurilor orale, pe care nu le vom detalia aici.

În contextul proiectului ReTeRom, ne interesează informațiile codificate în metadate și transcrieri aliniate. Ele sunt următoarele (Li and Yin, 2007):

- specificații relative la vorbitori: vârstă, gen, educație, calitate a vocii, dialect, accent, dar și un ID ales de vorbitor sau generat automat;
- specificații relative la contextul producerii înregistrării: dacă înregistrarea reprezintă vorbire spontană sau voce-în-citire; în cazul vorbirii spontane - dacă ea este solicitată (răspuns la întrebări etc.), în cazul vocilor-în-citire - metadatele textului reprodus, dacă înregistrarea reprezintă un dialog sau un monolog, dacă vorbirea e expresivă;
- specificații relative la înregistrare: nivelul tehnic al echipamentului cu care s-a produs înregistrarea (microfon, platforma soft-hard de înregistrare), contextul înregistrării (în laborator, în casă, pe stradă, în mașină etc.), nivelul zgomotului de fond (în db), detalii

¹⁸ <http://www.lrec-conf.org/proceedings/lrec2018/pdf/179.pdf>

asupra digitalizării semnalului (numărul de biți al eșantioanelor, rata de eșantionare, formatul de înregistrare al formei de undă), data, momentul și locul înregistrării, un ID al înregistrării (ales sau generat) pus în legătură cu numele fișierului¹⁹;

- specificații de adnotare: formatul convențiilor de adnotare ale transcrierii sunet-text (alinieră semnal-timp și text-timp sau direct semnal-offset caracter), adnotări ortografice, fonetice, prozodice, sintactice etc., instrumentele de segmentare și aliniere utilizate;
- criterii de validare: dacă au fost utilizate anumite criterii de validare a corpusului, la nivel de înregistrare dar și global, la nivelul corpusului ca întreg;
- specificații de distribuție și stocare: planul de distribuție, mediul de stocare, condiții de backup etc.

TEI (*Text Encoding Initiative*)²⁰ (Sperberg-McQueen and Burnard, 2018) dezvoltă și menține (prin consorțiul W3C²¹) un standard pentru reprezentarea textelor în format digital. Produsul principal este un set de linii directoare care specifică metode de codificare a textelor în formate interpretabile de către mașină (*machine-readable texts*) cu aplicații în domeniul umanioarelor digitale (*digital humanities*). Începând cu anul 1994, ghidul TEI a fost utilizat pe scară largă de biblioteci, muzee, editori și cercetători pentru a pregăti texte în vederea conservării lor de lungă durată și cercetării. În plus față de liniile directoare, consorțiul oferă o varietate de resurse și evenimente de instruire pentru învățarea TEI, informații despre proiectele care utilizează TEI, o bibliografie a publicațiilor legate de TEI și software dezvoltat pentru TEI sau adaptat acestuia.

Normele TEI P5 includ un capitol numit *Transcriptions of Speech*, în care sunt descrise problemele cele mai comune întâlnite în transcrierile de vorbire, informațiile general utile a fi plasate în metadatele care însoțesc transcrierile vocale, precum și probleme cu caracter specific care pot interveni în transcrierile vocale. Pe baza acestor informații pot fi create standarde ori convenții de reprezentare care să facă față oricăror particularități ale corpusurilor dedicate transcrierilor vocale. Proiectul nostru exploatează nu numai transcrierile în sine, ci pe acestea în pereche cu semnalul vocal original din care au fost generate transcrierile, în vederea dezvoltării de tehnologii de recunoaștere vocală și generare a vocii.

5. ASPECTE LEGALE

Problemele etice în înregistrarea și publicarea în spațiul public a ceea ce a fost produs într-un cadru privat și destinat unei audiențe limitate sunt întâlnite mai frecvent în tratarea textelor vorbite decât în cele scrise (TEI²²). Ca urmare, aspectele legale la constituirea unui corpus bimodal trebuie să aibă în vedere următoarele:

- drepturile de autor ale proprietarilor textelor pronunțate (în cazul corpusurilor voce-în-citire),
- acceptul vorbitorilor de a utiliza înregistrările făcute cu vocea lor, pentru scopuri didactice, de cercetare sau comerciale (depinzând de caz),
- acceptul vorbitorilor de a include în metadate informații de natură personală,
- convenții legale semnate între producătorii, distribuitorii și utilizatorii corpusului.

¹⁹ Lungimea numelui fișierului trebuie să respecte specificațiile atributelor de fișiere descrise în ISO-9960.

²⁰ <http://www.tei-c.org/release/doc/tei-p5-doc/en/html/ref-sp.html>

²¹ <https://www.w3.org/>

²² <https://quod.lib.umich.edu/cgi/t/tei/tei-idx?type=pointer&value=TSOV>

Convenția de colaborare trebuie semnată între producător și vorbitor înainte de începerea înregistrărilor.

6. CONCLUZII

Colecția de resurse ce urmează a fi organizată (dezvoltată, standardizată) în cadrul proiectului ReTeRom nu se adresează strict membrilor consorțiului proiectului. Această activitate are un orizont științific de durată, cu un pronunțat caracter de generalitate. Ea deschide calea pentru formarea unei viziuni științifice care să ghideze colectarea, adnotarea și distribuția de corpusuri de acest gen în România.

Raportul de față lămurește câteva chestiuni de terminologie a domeniului, sintetizează aspectele relevante în crearea corpusurilor bimodale, incluzând modalități de achiziție, prezintă un număr de exemple notorii de astfel de corpusuri, cu sintetizarea a 10 trăsături care au putut fi revelate din literatură, și prezintă preocupări de standardizare în realizarea corpusurilor bimodale, precum și aspecte legale.

Referințe

- Bainbridge, William (2004). Berkshire Encyclopedia of Human-computer Interaction. alle
- Godfrey, John J. and Edward Hollman (1997). *Switchboard I, Release 2*, Linguistic Data Consortium, Philadelphia.
- Grifoni, Patrizia (2009). Multimodal Human Computer Interaction and Pervasive Services. IGI Global. p. 37.
- Halle, M., & Stevens, K. (1962). *Speech recognition: A model and a program for research*. IRE Transactions on Information Theory, 8(2), pp. 155-159.
- Hughes, Thad, Kaisuke Nakajima, Linne Ha, Atul Vasu, Pedro Moreno and Mike LeBeau (2010). *Building transcribed speech corpora quickly and cheaply for many languages*. In: Proceedings of INTERSPEECH 2010.
- Karray, Fakhreddine, Alemzadeh, Milad, Abou Saleh, Jamil, Nours Arab, Mo. (2008). *Human-Computer Interaction: Overview on State of the Art*. International Journal on Smart Sensing and Intelligent Systems. 1 (1).
- Kurkovsky, Stan (2009). Multimodality in Mobile Computing and Mobile Devices: Methods for Adaptable Usability. IGI Global. pp. 210–211.
- Kurosu, Masaaki (2013). *Human-Computer Interaction: Interaction Modalities and Techniques*, Springer, p. 366.
- Li, A. and Y. Zu (2006). *Corpus Design and Annotation for Speech Synthesis and Recognition*. In *Advances in Chinese Spoken Language Processing*. Lee, C., Li, H., Lee, L., Wang, R., & Huo, Q. (eds.) World Scientific Publishing Co.

Li, Ai-jun and Zhi-gang Yin (2007). *Standardization Of Speech Corpus*. In: Data Science Journal, Volume 6, Supplement, 18 November.

MacWhinney, B. (2000). *The CHILDES Project: Tools for Analyzing Talk*. 3rd Edition. Mahwah, NJ: Lawrence Erlbaum Associates, copie online aici: <http://repository.cmu.edu/cgi/viewcontent.cgi?article=1181&context=psychology>.

Nascimento, Bacelar do, M. F., E. Cresti, M. Moneglia, A. Moreno Sandoval, J. Veronis, P. Martin, K. Choucri, V. Mapelli, D. Falavigna, A. Cid and C. Blum (2002). *The C-ORAL-ROM Project. New methods for spoken language archives in a multilingual romance corpus LREC*. In: M. C. RODRIGUES e C. SUAREZ ARAUJO, Proceedings of the *Third International Conference on Language Resources and Evaluation*, Paris: ELRA, vol. 1, pp. 2-10.

Palanque, Philippe and Paterno, Fabio (2001). *Interactive Systems. Design, Specification, and Verification*. Springer Science & Business Media, p. 43.

Pitt, M.A., L. Dilley, K. Johnson, S. Kiesling, W. Raymond, E. Hume and E. Fosler-Lussier (2007). *Buckeye Corpus of Conversational Speech (2nd release)* [www.buckeyecorpus.osu.edu] Columbus, OH: Department of Psychology, Ohio State University.

Rasso, T. and H. Mello (eds.) (2014). *Spoken Corpora and Linguistic Studies*, John Benjamins.

Ruhi, Ş. (2011). *Creating a sustainable large corpus of spoken Turkish for multiple research purposes (Ulusal Konuşma ve Dil Teknolojileri Platformu Kuruluşu: Türkçede Mevcut Durum Çalıştayı'nda sunulan görüş bildirisi)*. TÜBİTAK-TÜSSİDE, TÜBİTAK-BİLGEM, 6-7 Ekim 2011, Gebze.

Sperberg-McQueen, C.M. and Lou Burnard (2018). Original editors, revised and expanded under the supervision of the Technical Council of the TEI Consortium. TEI P5: Guidelines for Electronic Text Encoding and Interchange, version 3.3.0, last update: 31st January 2018, revision: f4d8439.

Shirai, K. and Fukuoka, T. (2018). *JAIST Annotated Corpus of Free Conversation*. In: Proceedings of LREC 2018, pp. 741-748.

Williams, J. R. (1998). Guidelines for the use of multimedia in instruction. In: Proceedings of the Human Factors and Ergonomics Society 42nd Annual Meeting, pp. 1447-1451.

Ziefle, M. (December 1998). *Effects of display resolution on visual performance*. In: Human factors. 40 (4), pp. 554-680.