



Raport științific și tehnic
Etapa a III-a, an 2020
„Resurse și tehnologii pentru dezvoltarea interfețelor om-mașină în
limba română”

Aceste rezultate au fost obținute prin finanțare în cadrul Programului PN-III Proiecte complexe realizate în consorții CDI, derulat cu sprijinul MEN – UEFISCDI,
Cod: PN-III-P1-1.2-PCCDI-2017-0818, Contract Nr. 73 PCCDI/2018

Denumirea organizației participante în proiect	Acronim organizație	Tip organizație	Rolul organizației în proiect (Coordonator/partener)
Institutul de Cercetări Pentru Inteligență Artificială “Mihai Drăgănescu”	ICIA	UNI	CO
Universitatea Tehnică din Cluj-Napoca	UTCN	UNI	P1
Universitatea Politehnică din București	UPB	UNI	P2
Universitatea "Alexandru Ioan Cuza" din Iași	UAIC	UNI	P3

Raport științific - tehnic proiect component CoBiLiRo

Date de identificare proiect	
Număr contract:	PN-III-P1-1.2-PCCDI-2017-0818, Nr. 73 PCCDI/2018
Acronim / titlu:	„COBILIRO:”
Titlu livrabil:	Raport științific și tehnic (Etapa a III-a, 2020)
Termen:	Noiembrie 2020
Editor:	Dan Cristea
Adresa de eMail editor:	danu.cristea@ gmail.com
Autori, în ordine alfabetică:	Anca Bibiri, Serban Boghiu, Dan Cristea, Daniela Gîfu, Felix Cristian Pericică, Ionuț Pistol, Mihaela Onofrei, Andrei Scutelnicu
Ofițer de proiect:	Cristian Stroe

1. Rezumatul etapei III

În această etapă a proiectului complex ReTeRom consorțiul și-a propus să consolideze și apoi să exploateze rezultatele acumulate în primii doi ani, cu obiectivele: existența unui Portal pregătit a primi și prelucra resurse bilingve românești, dezvoltarea în continuare a unei colecții de resurse care să corespundă formatului agreed de consorțiu, perfecționarea lanțurilor de prelucrări lingvistice și sonore, atât asupra componentelor textuale cât și vocale ale resurselor bimodale, care să permită alinieri între componentele vocale și textuale, recunoașterea cu minimum de erori a vocii, generarea expresivă a vocii și antamarea de aplicații bazate pe aceste tehnologii.

Conținutul repozitorului de resurse, acumulate pe portal în primii doi ani ai proiectului, a fost completat cu noi colecții bimodale, iar acestea au fost valorificate atât în cadrul platformei (prin realizarea de instrumente integrate generatoare de statistici) precum și în afara platformei, prin propunerea unei serii de proiecte ce utilizează tehnologiile dezvoltate în cadrul proiectelor partenere. Similar celorlalte etape, a fost realizată și o activitate de diseminare, în principal prin publicații. De interes special pentru platforma dezvoltată a fost respectarea drepturilor de autor și anonimizarea datelor (cu precădere cele vocale) contribuie de diverși furnizori.

2. Descrierea rezultatelor activităților etapei III

Activitatea A3.1 este intitulată „**Augmentarea corpusului vocal prin noi înregistrări vocale care dublează texte existente în corpusul CoRoLa**”.

Preocuparea membrilor echipei a fost orientată în principal în direcția achiziționării din diverse surse ori chiar a creării prin mijloace proprii de înregistrări sonore însoțite de transcrierile lor textuale și completarea metadatelor corespunzătoare. Ca urmare, **Raportul A3.1** prezintă rezultatele activității, sintetizând noile colecții de documente bimodale ce au fost urcate pe platforma CoBiLiRo în decursul acestui an, după ce au trecut prin etapele de obținere a drepturilor de autor, unele dintre ele fiind aliniate, parțial automat - parțial manual, cu înregistrările sonore. Resursele nou urcate în Portal în acest an al proiectului includ 125 de piese, aparținând următoarelor colecții (descrise în atributul **@collection** din header-ul metadatei):

- *Ghici Cine Vine La Cină?* (cu 13 interviuri, voci de bărbați);
- *Alma Mater Iassiensis* (cu 13 interviuri, voci de bărbați);
- *SoRoEs* (cu 91 de înregistrări pe teren, voci de bărbați și de femei din diferite regiuni ale României);
- *Povești* (cu 8 înregistrări de povești realizate în studio).

Toate resursele au atributul **@distribution** setat la valoarea: <https://creativecommons.org/licenses/by-nc/3.0/>. Colecțiile încărcate pe platforma CoBiLiRo respectă formatul standardizat al fișierelor vorbire-text relativ la metadata (care conțin informații header și alinieri). Formatul, agreed de toți partenerii proiectului complex ReTeRom, este difuzat în proiect în livrabilul A1.3 și descrie tipurile:

- *file*: care separă fiecare componentă <speech> de format WAV/MP3, cu includerea atributului **@speechFile**, care identifică numele fișierului conținând înregistrarea vocală;
- *start-stop*: care marchează bornele temporale de început-sfârșit ale componentelor <speech> în fișiere unice asociate unui obiect <tei>, cu includerea atributelor **@start** și **@stop**, cu valori reale,

reprezentând momentele de început și de sfârșit ale înregistrării vocale; în acest caz, numele fișierului în care sunt referite bornele **@start** și **@stop** ale elementelor <speech> sunt date într-un atribut **@speechFile** încorporat elementului <teiHeader>.

Identificarea în fișierele colecției corpusului Cobiliro a primei ori a celei de a doua opțiuni de marcarea a segmentării se face în elementele <teiHeader>, care includ un atribut **@speechSegmentation**, cu una din valorile: **“file”**, respectiv **“start-stop”**, tipul de segmentare fiind unitar în lungimea fiecărui obiect. Menționăm că majoritatea obiectelor înregistrate în Portal sunt de tip *file*.

Înregistrările sonore sunt neuniforme în privința apropierii de text, unele fiind prelucrate deja de editori (pentru îndepărtarea unor repetiții, pauze, bâlbâieli, inerente vorbirii libere) în vederea tipăririi textelor (ca de exemplu colecțiile de interviuri), altele urmărind destul de aproape idiosincraziile din vorbirea liberă.

Calitatea înregistrărilor sonore este, de asemenea, diversă: de la înregistrări realizate în condiții profesionale (microfon performant, cameră anecoidă), până la înregistrări în condiții naturale (pe stradă, cu zgomot sau muzică în surdina etc.). Părerea noastră este că fiecare dintre aceste tipuri de resurse bimodale este utilă în experimentele de antrenare a tehnologiei de recunoaștere a vocii, cât și în cea de sinteză vocală, pentru că, în condiții reale de transpunere a vocii în text ori invers, toate aceste condiții pot apărea, ca urmare tehnologiile trebuie să poată fi antrenate pe resurse care să reproducă condiții de calitate cât mai diverse.

A fost, de asemenea, agreată de către membrii colectivului CoBiLiRo o procedură de dublare în voce (*read speech*) a unei colecții de fișiere textuale aflate deja în CoRoLa (*Corpusul Reprezentativ al Limbii Române Contemporane*), corpus care a fost construit între anii 2014-2017, prin colaborarea dintre două institute de cercetare ale Academiei Române: Institutul de Inteligență Artificială “M. Drăgănescu” din București și Institutul de Informatică Teoretică din Filiala Iași a Academiei Române, cu ajutor (în cadrul proiectului DRuKoLa finanțat de Fundația Humboldt), din partea Universității București și a Institutului Limbii Germane din Mannheim. Această activitate se va continua în perioada de prelungire a proiectului, până în martie 2021. Rezultatul acestei activități are o dublă semnificație, pentru că astfel se mărește atât colecția de înregistrări din Portalul CoBiLiRo, cât și cea din CoRoLa.

Consortiul a acordat o grijă deosebită respectării cadrului legal de procurare și utilizare a resurselor în scopuri de cercetare, pentru respectarea drepturilor de autor și a securității informațiilor private. Pentru fiecare dintre resursele încărcate pe platforma CoBiLiRo, membrii Consortiului care le-au achiziționat au semnat protocoale de colaborare cu autorii ori proprietarii lor. Alte discuții s-au purtat pentru stabilirea unui punct de vedere relativ la posibilitatea de a ceda drepturi de utilizare unor terțe entități. Cum părerile au fost divergente în această privință, s-a convenit ca fiecare partener să stabilească propriii termeni de licențiere asupra resurselor procurate ori construite, iar în cazul în care resursele ce sunt solicitate de terți au intrat în posesia consortiului prin donații din partea altor autori sau furnizori să se revină asupra protocoalelor semnate cu aceștia inițial pentru a obține din partea lor acceptul de extensie a dreptului de utilizare către terți. Opinia membrilor Consortiului asupra drepturilor de exploatare a resurselor din Portalul Cobiliro este către o utilizare liberă pentru scopuri de cercetare. Atunci când rezultatele cercetării sunt însă explicit enunțate ca urmărind obținerea de profit, opiniile sunt divergente, pentru că, pe de o parte nu se poate pune problema participării la profit a unora dintre parteneri (spre exemplu, obținerea de profit nu este admisă în Statutul Academiei Române), iar pe de altă parte, e discutabil aspectul moral al obținerii de profit din resurse cedate gratuit de proprietari, chiar dacă acestor resurse le-a fost adăugată o plus-valoare rezultată în urma cercetării.

Toate obiectivele incluse în plan la această activitate au fost realizate.

Activitatea A3.2 este intitulată „**Augmentarea corpusului voce-text: completare de metadate, alinieri prin rularea algoritmilor dezvoltati în proiectele P2, P3 și P4 și adnotări manuale și semiautomate la corpusul voce/text**”.

În **Raportul A3.2** este descris modul în care resurselor bimodale aflate în Platforma CoBiLiRo li s-au asociat alinieri între secțiunile de voce și de text, majoritar automat, prin rularea tehnologiilor-servicii create în proiectele P2, P3 și P4. Pentru că unele transcrieri textuale nu redau cu acuratețe echivalentele lor sonore (cum se întâmplă, de exemplu, în colecția de interviuri, înregistrate în direct, care apoi au fost transcrise în text și voit „înfrumusețate” în vederea tipăririi), alinierea au pus în evidență doar zonele de coincidență.

Folosind aceste servicii, au rezultat două tipuri de procesări:

- Un prim proces de aliniere dintre fișierul audio și text (Figura 1), în care sunt marcate timpii de început și sfârșit al fiecărui cuvânt.
- Cel de-al doilea tip de procesare în care au fost marcate *timestamp*-uri aferente secvențelor de cuvinte. Formatul este după cum urmează: câte un fișier text pentru fiecare resursă audio. Pe fiecare linie este câte o propoziție (cu *timestamps* la nivel de cuvânt într-o variantă; fără *timestamps* în cealaltă variantă - după cum se poate observa în Figura 2). În cadrul lunilor de implementare rămase aferente proiectului se va

realiza prelucrarea acestui rezultat oferit de serviciul TADARAV, pentru a aduce acest rezultat la formatul CoBiLiRo convenit (livrabilul A3.1).

Datorită unor probleme tehnice întâmpinate pe serverul care găzduiește Portalul CoBiLiRo, o serie de înregistrări ale vechilor resurse au fost pierdute din baza de date. Pentru ca ele să redevină vizibile în Portal, s-a realizat un soft care a parcurs lista tuturor directoarelor hard-ului care găzduia back-up-urile resurselor și le-a reintrodus în baza de date a Portalului.

În momentul de față pe server există aproximativ 70 de mii de fișiere audio cu perechile textuale respective ce conțin transcrieri ale înregistrărilor sonore, precum și fișierele corespunzătoare cu metadate.

Toate obiectivele incluse în plan la această activitate au fost realizate.

Activitatea A3.3 este intitulată „**Realizarea de statistici privind corpusul bimodal voce/text**”. Raportul prezintă toate tipurile de statistici implementate în decursul acestui an și care pot fi lansate la cerere asupra resurselor voce-text existente în Portalul CoBiLiRo. Statisticile prezintă, la cererea utilizatorului, parametri calculați din exploatarea informațiilor conținute în fișierele text și audio, în alinierea dintre cele două, cât și în metadatele aferente.

Componenta de procesare lingvistică necesară obținerii de valori statistice, este realizată sub forma unui script Python, încadrat într-un web API, care poate fi accesat folosind principiile arhitecturii Restful, apelat doar la nivel intern (pe server); aici s-a folosit tehnologia **Flask Restful**.

Pentru realizarea prelucrării și procesării la nivel de limbaj natural a textelor a fost utilizată librăria open-source **Spacy**, mai exact modelul **ro_core_news_sm**; acest model atribuie vectori de tokeni specifici contextului și etichete ale părților de vorbire având o acuratețe de 95.62%. Ca alternativă poate fi folosită componenta de analiză a textului produsă de TEPROLIN; motivul pentru care nu a fost inclusă în versiunea inițială este că nu toate fișierele de pe platformă au fost procesate de acea componentă. Pentru finalul proiectului avem în plan tranziția spre tehnologia TEPROLIN.

Statisticile implementate se referă la: liste de obiecte obținute pe baza informațiilor din header-ul metadatelor (ca de exemplu: obiecte în care componentele vocale au formatul wav/mp3, obiecte cu componenta vocală segmentată la nivel de fișier/cuvânt, obiecte voce-text aliniate la nivel de frază/cuvânt/fonem etc.), statistici de natură lexicală (densitatea lexicală, numărul de token-uri, de leme diferite, de forme flexionare distincte, raportul tip-token etc.), precum și statistici la nivelul întregii colecții de resurse (numărul total de ore de vorbire, de token-uri, densitatea lexicală medie, raportul mediu al numărului de token-uri pe categorii etc.).

Oferirea unor statistici relevante pentru o colecție de resurse contribuie atât la mărirea încrederii unor potențiali noi contributory cât și la potențialul resurselor din colecție de a-și găsi utilitatea în diverse contexte. Statisticile deja incluse pe platformă precum și cele care vor fi incluse până în aprilie 2021 promit să asigure acest avantaj pentru colecția de resurse bimodale, principalul rezultat al proiectului CoBiLiRO.

Până la terminarea proiectului avem în vedere ca procesele de determinare a lungimii fișierelor să fie lansate doar asupra fișierelor nou introduse sau updatate. Serverul e capabil să urmărească fișierele asupra cărora s-a umblat și să calculeze/recalculeze lungimile doar pentru acestea. Odată aflată o lungime, ea se va înscrie în metadate, ceea ce înseamnă că nu va trebui recalculată la fiecare cerere din partea utilizatorului.

Toate obiectivele incluse în plan la această activitate au fost realizate.

Activitatea A3.4 este intitulată „**Proiectare de aplicații de exploatare a corpusului bimodal și a tehnologiilor de prelucrări textuale și voce, create în proiectele P2, P3, P4**”.

Raportul activității prezintă un număr de posibile aplicații care să demonstreze potențialul tehnologiilor dezvoltate în celelalte proiecte componente, în scopul mării vizibilității proiectului, în special după finalizarea acestuia. Tehnologiile vizate în special sunt cele descrise în rapoartele 3.7 (Definitivarea, testarea, validarea și împachetarea într-o soluție „*ready-to-use*” a platformei integrate și configurabile de prelucrare a textelor în limba română.), 3.10 (Îmbunătățirea soluției de filtrare și aliniere a transcrierilor aproximative cu semnalul de vorbire), 3.15 (Dezvoltarea unei noi tehnologii pentru adaptarea vocii sintetice la stilul și expresivitatea unui nou vorbitor) și 3.17 (Integrare tehnologie nouă și demonstrare în realizarea interfețelor om-mașină pentru sinteza text-vorbire). O parte din proiectele descrise au depășit deja faza de proiect, fiind în stadii relativ avansate de implementare (în activități cu studenții de la UAIC-FII).

Cele 6 proiecte de aplicații descrise sunt prezentate succint mai jos:

Aplicația 1: Suport pentru învățarea limbii române - PD builder. Un dicționar de pronunție pentru cuvintele unei limbi poate constitui un suport semnificativ pentru cei care doresc să învețe acea limbă. Colecția de resurse colectate pe platforma CoBiLiRO are potențialul să atingă o dimensiune suficientă pentru a include pronunția mării majorități a cuvintelor din limba română; de aici a plecat ideea de a dezvolta o

tehnologie ce caută în resursele din Portal și reproduce pronunții ale unor cuvinte, la cererea utilizatorului. Sunt utilizate adnotările produse de TEPROLIN și alinierea produsă de TADARAV, marcate ce însoțesc majoritatea resurselor existente pe platforma CoBiLiRo.

Aplicația 2: Analiza corpusurilor bimodale. Proiectul propune dezvoltarea unui sistem capabil să evalueze calitatea unui corpus bimodal voce-text din perspectiva unei alinieri automate sau manuale. O parte din posibilele erori semnalate ar putea fi corectate, lucru care ar mări semnificativ calitatea alinierii automate dar și a resursei bimodale în general. Prin corelarea statistică a mai multor parametri, pot fi detectate posibile diferențe între conținutul înregistrării și transcrierea ei în text, diferențe care pot avea cauze diverse: segmente adiționale prezente în text față de înregistrare sau invers, transcrieri aproximative, particularități ale pronunției și prozodiei vorbitorului etc.

Aplicația 3: *I listen to my speaking agent reading book fragments as I walk by.* Aplicația are la bază o colecție de texte care abundă în entități geografice, marcate XML explicit, textele fiind însoțite de metadate care descriu minimum: autorul și titlul cărții, anul de apariție și editura. Instalată pe un dispozitiv mobil, ea va semnala proximitatea telefonului față de locațiile menționate în texte și va citi acele fragmente care includ mențiunile respective. În felul acesta, o plimbare printr-un mare oraș se poate transforma într-o călătorie literară. Aplicația se adresează persoanelor care ar dori să primească sugestii de lecturi și doresc să afle lucruri noi despre locurile pe care le vizitează. Ea își propune să îmbine plăcut literatura și tehnologia.

Aplicația 4: Sistem de sinteză text-vorbire (TTS) și clonarea vocii în limba română cu metoda învățării prin transfer. Aplicația poate fi utilă persoanelor care și-au pierdut capacitatea de a vorbi, dar există înregistrări precedente pierderii vocii. Alte utilități pot fi imaginat în învățământul la distanță, în muzee, în colecțiile de moșteniri culturale etc., unde tehnologia poate fi folosită pentru reproducerea vocii personalităților istorice. Clonarea vocii reprezintă un caz particular de TTS, în care un text este sintetizat în vocea unui vorbitor necunoscut folosindu-se un număr limitat de înregistrări audio ale acestuia acompaniate de transcrierile lor.

Aplicația 5: Asistent inteligent al ședințelor online. Pandemia Covid a obligat ca o parte neneglijabilă a activităților colectivităților umane să migreze în mediul online. Propunem realizarea unui instrument sau a unei colecții de API-uri care să se poată integra aplicațiilor de teleconferințe, în scopul extragerii de informații și generării de rapoarte din discuțiile purtate și din mesajele schimbate. Astfel, se poate imagina: transcrierea conversațiilor în text, realizarea de rezumate ale întâlnirilor, extragerea de informații punctuale, pe anumite segmente din conferință ori din intervențiile unor anumiți vorbitori, generarea automată de procese verbale, interogarea minutei ori a procesului verbal generat, pe bază de cuvinte cheie, căutări de secvențe sonore în înregistrarea sonoră etc.

Aplicația 6: *Tracking assistant.* O aplicație Android care își propune să asiste persoanele care suferă de boala Alzheimer în a-și aminti traseele pe care le-au efectuat în timpul zilei printr-o interfață ușor de folosit, în limbaj natural. Pe baza datelor receptate din mai multe canale ale unui telefon mobil, se poate reface traseul și tipul de activitate al persoanei. La sfârșitul zilei pacientul ar purta un dialog cu sistemul inteligent, care l-ar ajuta să-și rememoreze activitățile de peste zi.

Toate obiectivele incluse în plan la această activitate au fost realizate.

Activitatea A3.5 este dedicată „Diseminării, drepturilor de autor și protejării intereselor vorbitorilor (anonimizare, unde este cazul)”.

Diseminarea rezultatelor cercetărilor obținute pe parcursul acestui an s-a realizat, în principal, prin publicarea de lucrări în jurnale, conferințe și workshopuri naționale și internaționale. Publicațiile se adresează comunității cercetătorilor din domeniul prelucrării limbajului natural și al vorbirii, lingviștilor, cu care comunitatea noastră științifică are legături foarte puternice, cât și întreprinzătorilor din industria IT, interesați să exploateze în practică perspectivele pe care proiectul ReTeRom le deschide în viitor (de ex. realizarea de interfețe om-mașină care să utilizeze recunoașterea și sinteza vorbirii).

Totodată, în cadrul acestei activități au fost întreprinse măsuri de respectare a drepturilor de autor ale furnizorilor resurselor achiziționate și de protejare a intereselor contributorilor (în esență, persoanele cărora li s-a înregistrat vocea). Instituția implicată în realizarea acestei etape a fost UAIC, realizatoarea componentei CoBiLiRo, sprijinită de toți membrii consorțiului ReTeRom.

Raportul prezintă lista completă a lucrărilor publicate de parteneri în anul curent: 4 de partenerul UAIC, 5 de partenerul UPB, 5 de partenerul UTCN și 4 de partenerul ICIA.

Drepturile de autor au fost protejate prin protocoale de colaborare, semnate de furnizor și de o instituție parteneră în proiect. La achiziționarea componentelor vocale din colecția de obiecte CoBiLiRo au fost respectate prevederile legale de protecție a intereselor personale ale vorbitorilor.

Raport științific - tehnic proiect component TEPROLIN

Date de identificare proiect	
Număr contract:	PN-III-P1-1.2-PCCDI-2017-0818, Nr. 73 PCCDI/2018
Acronim / titlu:	„TEPROLIN”
Titlu livrabil:	Raport științific și tehnic (Etapa a III-a, 2020)
Termen:	Noiembrie 2020
Editor:	Dan Tufiș
Adresa de eMail editor:	tufis@racai.ro
Autori, în ordine alfabetică:	Verginica Mititelu, Elena Irimia, Radu Ion, Vasile Păiș, Horia Cucu, Lucian Georgescu, Cristian Manolache, Adriana Stan, Beáta Lórinçz
Ofițer de proiect:	Cristian Stroe

1. Rezumatul etapei III

Această activitate are rolul de a evalua erorile facute de sistemul RAV al UPB, și reantrenarea cu date corectate. De asemenea au fost evidențiate îmbunătățirile pe care lexiconul construit de noi le aduce aplicațiilor de TTS dezvoltate la UTCN, prin comparare cu rezultatele obținute cu resurse folosite anterior, și apoi prin analiza a diferite scenarii de evaluare care folosesc informația din lexicon integral sau parțial (atât ca număr de intrări, selectate aleatoriu sau sistematic, cât și ca tipuri de informație utilizată). Totodată în această etapă a fost realizată definitivarea, testarea, validarea și împachetarea într-o soluție „ready-to-use” a platformei integrate și configurabile de prelucrare a textelor în limba română.

2. Descrierea rezultatelor activităților etapei III

Activitatea 3.6: Analiza erorilor sistemelor ASR și TTS antrenate în proiectele 3 și 4 pe corpusul bimodal agregat în proiectul 1, adnotat și corectat

Recunoașterea automată a vorbirii în limba română - ASR

În faza actuală a proiectului, echipa UPB a îmbunătățit performanțele sistemului de RAV pe de o parte prin antrenarea pe noi seturi de date create și accesibile prin proiectul CoBiLiRo, iar pe de altă parte prin recrearea modelului de limbă prin adăugarea de texte noi. (Pentru date despre cantitatea de texte noi folosite, despre tipurile de informații existente în aceste corpusuri folosite la antrenare, despre noile performanțe ale sistemului de RAV, a se vedea raportul fazei întocmit de echipa UPB.)

Evaluarea la nivel de cuvânt arată, așa cum era de așteptat, o îmbunătățire a tuturor criteriilor de analiză: cuvinte recunoscute greșit, cuvinte inserate eronat, cuvinte netranscrise.

Evaluarea rezultatelor sistemului de RAV asupra corpusului din platforma CoBiLiRo a arătat că, în afară de cuvintele obișnuite, se produc o serie de erori în cazul a două tipuri de cuvinte: nume de entități și cuvinte cu cratimă. Toate acestea sunt cuvinte inexistente în lexiconul folosit de sistemul de RAV (engl. „out of vocabulary word”). Pentru fiecare categorie am adoptat o metodă de lucru considerată adecvată.

- Nume de entități: am extras din diverse surse liste separate de nume de persoane, nume de locuri, nume de organizații/firme/etc. cu scopul de a îmbogăți lexiconul sistemului de ASR. Listele create sunt disponibile la adresa https://relate.racai.ro/resources/ro_namelists_20201013.zip.
- Cuvinte cu cratimă: folosind un lexicon intern ICIA, am validat o parte dintre cuvintele cu cratimă, adică acelea corecte, dar absente din lexiconul ASR.
- Cuvinte obișnuite: s-a folosit aceeași procedură ca la cuvintele cu cratimă.

N.B. Nici una dintre aceste abordări nu rezolvă complet problema evaluării rezultatelor sistemului de RAV, atâta timp cât textele din corpusul CoBiLiRo conțin greșeli de scriere, iar rezultatele sistemului sunt comparate cu forma scrisă a corpusului. Pentru aceasta, am folosit o măsură de similaritate a tuturor tipurilor de cuvinte (Levenshtein distance) pentru a indica forma cea mai apropiată din listele de nume și lexiconul folosite. Aceste aproximări au fost validate manual și s-au creat liste de corecturi propuse pentru îmbunătățirea transcrierilor existente în corpusul din platforma CoBiLiRo.

O altă observație rezultată din lucrul cu datele este nevoia de uniformizarea textelor din corpus în ceea ce

privește tipul de litere cu diacritice folosit, i.e. utilizarea literelor ș și ț conform standardului actual.

În ceea ce privește alte tipuri de erori pe care le-am identificat printre rezultatele sistemului de RAV, unele au caracter general și considerăm că țin de procedura de evaluare, în sensul că o relaxare a restricțiilor impuse în evaluare ar arăta o îmbunătățire a performanțelor. Este vorba aici despre posibilitatea de a scrie unele pronume clitice cu sau fără cratimă, distincția dintre cele două ținând uneori doar de durata mai scurtă, respectiv mai lungă a unui sunet: ex.: *te-aștept* versus *te aștept*. Întrucât nu se poate stabili o limită de durată dincolo de care să se scrie fără cratimă, considerăm că recunoașterea oricăreia dintre variante trebuie acceptată ca fiind corectă. Similar, doar cu un sunet în plus, sunt alte clitice, care pot pierde o vocală la o rostire mai rapidă: ex.: *mă aștept* versus *m-aștept*.

În aceeași situație se află inițiala în vocala *î* a cuvintelor precedate de anumite cuvinte funcționale: ex: *la început* versus *la-nceput*, *a început* versus *a-nceput*, *să înceapă* versus *să-nceapă* etc.

Sinteza automată a vorbirii în limba română - TTS

În rapoartele anterioare aferente activităților 1.8 și 2.7 descriam în detaliu lexiconul cu informație extinsă, dezvoltat și validat în cadrul proiectului 2, TEPROLIN; pentru fiecare intrare din lexicon, informația asociată reprezintă: lema (forma de dicționar a cuvântului), eticheta morfo-sintactică în format MSD (Erjavec, 2004), împărțirea în silabe, marcarea accentului (printr-un apostrof) și transcrierea fonetică a formei ocurență. Am subliniat de asemenea importanța acestui lexicon în aplicațiile pentru recunoaștere automată și pentru sinteza vorbirii, susținând că este esențial ca o astfel de resursă să fie de bună calitate. În raportul aferent activității 2.7 am descris în detaliu procesul complex de validare prin care lexiconul ReTeRom a trecut până în acest moment. Activitatea de validare a constat atât în corectare manuală intrare cu intrare, cât și în corectare automată, prin implementarea ca expresii regulate a diverse tipuri de reguli, bazate pe cunoaștere lingvistică.

În această etapă, împreună cu echipa parteneră UTCN, am putut testa și dovedi utilitatea lexiconului în cadrul unei aplicații de sinteză a vorbirii. Pentru a simplifica procesul de evaluare, am restrâns contextul evaluat la nivel textual, deoarece calitatea vorbirii sintetizate este greu de evaluat în mod obiectiv, cuantificabil. Astfel, rețeaua neuronală folosită pentru evaluarea lexiconului are ca scop predicția (în format text) concurentă a silabificării, accentului și transcrierii fonetice pornind doar de la forma (ortografică a) cuvântului, în lipsa unui context de utilizare.

În Tabelul 1 puteți vizualiza câteva perechi de date de intrare/ieșire ale rețelei: punctul marchează limitele silabelor iar accentul este marcat prin apostrof urmând vocala accentuată.

Intrare	Ieșire
abandonarăți	a . b a n . d o . n a' . r @ t s i 0
basculantei	b a s . k u . l a' n . t e j
Ciclopul	tS i . k l o' . p u l

Tabelul 1. Exemple de date de intrare și ieșire pentru rețeaua neuronală

În familia arhitecturilor neuronale utilizate pentru a învăța funcții secvență-la-secvență (en. Sequence-to-sequence), rețelele neuronale convoluționale (CNN, Gehring et al. 2017) și rețelele bazate pe atenție (Vaswani et al., 2017) sunt cele care au demonstrat cea mai bună acuratețe în probleme de procesare a limbajului natural. De aceea, într-o primă etapă a acestui studiu, s-a realizat o evaluare CNN vs transformer (vezi (Stan, 2020) pentru detalii despre implementarea acestora), în cadrul căreia rezultatele rețelei transformer au fost evident mai bune. Arhitectura rețelei are o structură de tip encoder-decoder și este prezentată în figura 1.

Pentru selectarea parametrilor rețelei transformer ne-am bazat pe rezultatele din (Stan, 2020), rezultând o rețea cu 3 unități encoder, 4 unități decoder, 4 centre de atenție, dimensiunea stratului ascuns de 1.024 și dimensiunea vectorului embedding¹ de 128. Ponderile vectorului embedding sunt inițializate aleatoriu înainte de antrenare. Dimensiunea lotului (eng. „batch”) a fost stabilită la 512 iar pentru actualizarea ponderilor s-a folosit optimizatorul Adam cu o rată inițială de învățare de 0,0002. După 50 epoci, rata de învățare a fost redusă cu un factor de 0,2 iar procesul de antrenare a fost limitat de un criteriu de oprire bazat pe valoarea funcției de cost după 5 epoci. În procesul de evaluare, setul de date a fost divizat în loturi de de 70%-10%-20%, dedicate proceselor de antrenare, validare și, respectiv, testare. Seturile de antrenare și

¹ reprezentare codificată numeric a șirului sursă/țintă

validare variază în funcție de natura experimentelor, dar setul de testare este fix (aproximativ 80.000 de intrări).

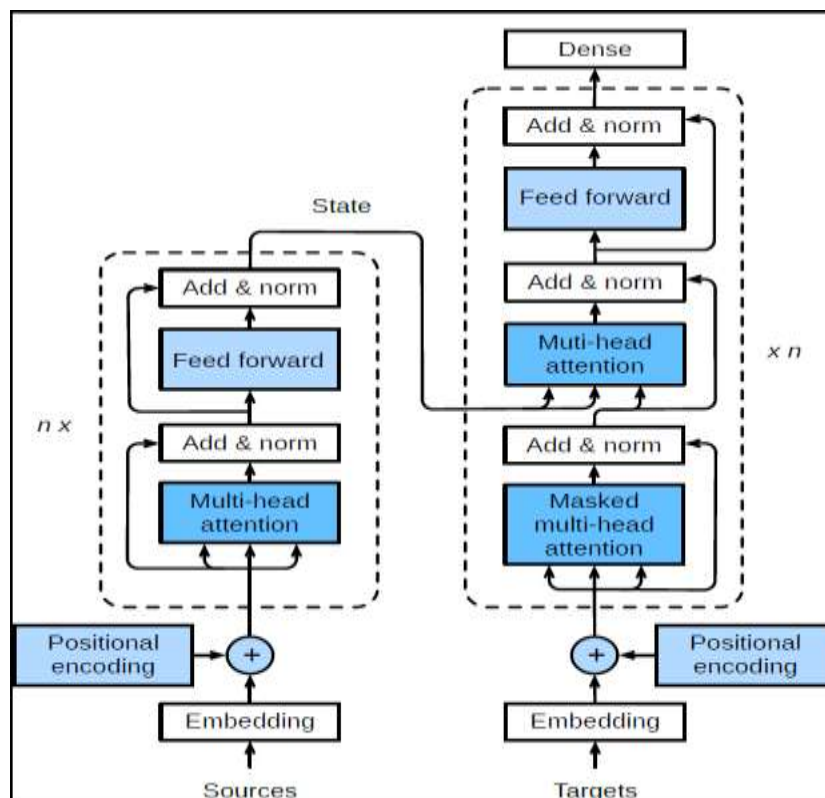


Figura 1. Arhitectura rețelei transformer folosite în experimentele noastre

În prealabil, ne-a interesat să investigăm utilitatea folosirii lexiconului ReTeRom în comparație cu folosirea concurentă a resurselor MaRePhor, RoSyllabiDict și a accentului din DEX (pe care și lexiconul nostru se bazează, extinzându-le), resurse utilizate până în acest moment pentru predicție concurentă de informație lexicală în aplicații de TTS la UTCN. Ulterior, am conceput scenarii de evaluare mai complexe, care să ne ajute să găsim modalități mai eficiente de folosire a resursei dezvoltate, după cum se va vedea în continuare.

Astfel, ne-a interesat îmbunătățirea performanței sistemului de predicție atunci cât setul de date crește gradat și identificarea momentului de platou al acestei îmbunătățiri. În acest scop, au fost selectate partiții aleatoare ale datelor pornind de la un minim de 5.000 de intrări.

Am urmărit de asemenea îmbunătățirea performanței în funcție de cantitatea de date atunci când selecția unui set redus de date se face asigurându-ne ca dispunem de întreaga acoperire a lexiconului la nivel de leme. Cum limba română este o limbă cu morfologie bogată, dar are o resursă de tip lexicon sau corpus cuprinde toate variantele morfologice ale unei leme, am vrut să vedem dacă performanța scade drastic atunci când doar lema, sau doar o formă (care nu este întotdeauna egală cu lema) sau două forme asociate fiecărei leme se regăsesc în datele de antrenare. Ne-am rezumat la a reduce numărul de forme doar pentru cuvintele conținut cu morfologie bogată: adjective, verbe și substantive, iar pentru restul tipurilor de cuvinte, numărul formelor asociate s-a păstrat intact. Astfel au fost derivate trei subset-uri, *1-FORM* (30.150 intrări, cu o singură formă pentru adjective, verbe și substantive) și *2-FORMS* (55.185 intrări, două forme pentru fiecare adjectiv, verb și substantiv) și *LEMMA* (35.890 intrări, cu o singură formă pentru adjective, verbe și substantive, care este întotdeauna identică cu forma dicționar, sau lema). În procesul de selectare a formelor din subset-urile *1-FORM* și *2-FORMS*, ne-am asigurat că reprezentăm în mod balansat fiecare trăsătură morfologică asociată acestor părți de vorbire (gen, număr, caz, mod, timp, persoană, etc.) pentru a încerca să păstrăm diversitatea terminațiilor morfologice existentă în limba română. Desigur, într-un scenariu real, nu se poate controla păstrarea integrală a acestei diversități (adică nu ne putem asigura că toate terminațiile sunt reprezentate într-un corpus sau într-un lexicon extras dintr-un corpus), dar o varietate ne-exhaustivă a formelor este naturală, pe când prezența exclusivă a formelor de dicționar este improbabilă. În plus, ne-a interesat să investigăm potențialul utilizării părții de vorbire (primul simbol din eticheta morfosintactică) și a etichetei morfo-sintactice (MSD) ca trăsături suplimentare de intrare. Teoretic, această informație ar putea

ajuta sistemul sa diferențieze între omografele care nu sunt omofone. În practică, în lexiconul nostru, omografele care au rostire diferită în funcție de POS-ul sau MSD-ul asociat reprezintă aproximativ 2000 de intrări și nu ne așteptăm ca rezolvarea acestei probleme de predicție să aibă un impact mare asupra ratei de eroare. Eventuala rezolvare reprezintă în schimb un plus de acuitate lingvistică adus sistemului. De asemenea, presupunem că rețeaua poate învăța să asocieze anumite terminații morfologice (și pronunțiile lor, care sunt specifice atunci când sunetele reprezintă terminații) cu anumite părți de vorbire/etichete morfosintactice, fapt ce ar putea compensa pentru lipsa tuturor formelor asociate cuvintelor conținut. Vom verifica această presupunere în scenariile de evaluare bazate pe subset-urile LEMMA, 1-FORM și 2-FORMS. În toate aceste scenarii de evaluare, output-ul rețelei transformer a fost, așa cum am prezentat în Tabelul 1, o combinație a celor trei sarcini de predicție lexicală, iar inputul a fost fie forma ortografică a cuvântului, fie forma împreună cu partea de vorbire, fie forma împreună cu eticheta morfosintactică asociată.

Drept metrici de evaluare, am folosit *word error rate* (WER, ro. „rata de eroare la nivel de cuvânt”) și *character error rate* (CER, ro. „rata de eroare la nivel de caracter”), evaluate pe șiruri de caractere care includ marcajele de silabificare și accent. CER a fost măsurată folosind distanța Levenshtein (Levenshtein, 1966) între secvența de caractere prezisă și cea țintă. În plus, ne-am dorit să evaluăm erorile introduse de fiecare sarcină de predicție, motiv pentru care am calculat WER și CER eliminând din predicția concurentă informația furnizată de una sau mai multe dintre sarcini.

Evaluarea prealabilă, în care comparăm resursa noastră cu cele folosite anterior, ne arată că, cu lexiconul ReTeRom, putem obține o rată a erorii WER de 3,08 și o rată CER de 1,08 pentru predicția concurențială, atunci când folosim toată informația disponibilă (inclusiv eticheta morfosintactică), ceea ce reprezintă o reducere importantă a ratei erorii de la 10,47 WER și 3,3 CER, obținute cu resursele utilizate anterior la UTCN.

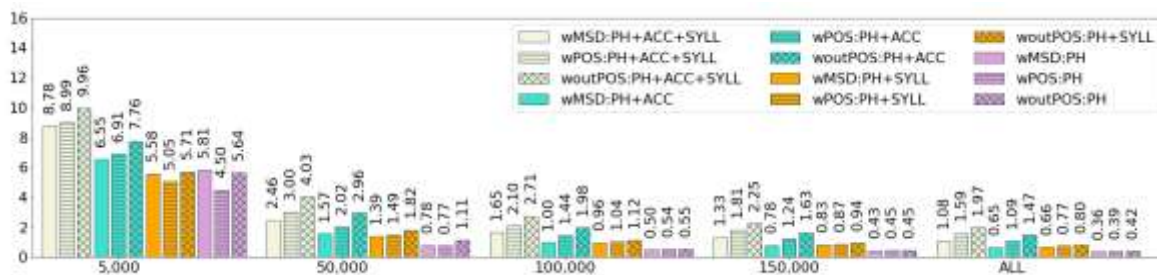


Figura 2. a

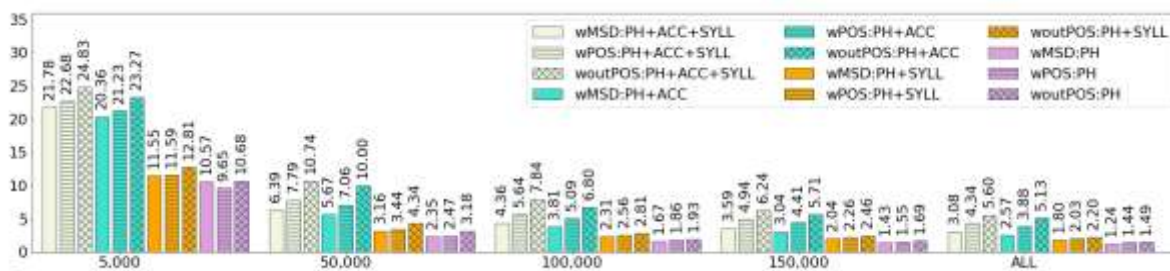


Figura 2. b

Fig 2. (a) CER și (b) WER) pentru diferite dimensiuni ale datelor de antrenare selectate în mod aleator, evaluate pentru: predicție lexicală completă (PH+ACC+SYLL), eliminând din predicție silabificare (PH+ACC), eliminând din predicție silabificarea și accentul (PH). Rezultatele indică ratele de eroare cu POS (wPOS), cu MSD (wMSD) și fără nici o informație suplimentară în afară de forma cuvântului (woutPOS).

În Figura 2, unde prezentăm rezultatele scenariului de evaluare 1 combinat cu scenariul de evaluare 3, culorile corespund combinației de informație menținută în predicție iar simbolul de hașurare discriminează între cele două tipuri de input (cu sau fără informație morfosintactică). Se poate observa că după 100.000 de intrări, creșterea acurateții atinge un platou, dar totuși există o îmbunătățire relativă de 24% începând de la 100.000 intrări și până la setul de date complet, în termeni de WER, pentru predicția întregii secvențe, atunci când se folosește informația morfosintactică. Din punct de vedere al contribuției sarcinilor de predicție la rata erorii, accentul are cea mai mare influență, lucru de așteptat pentru limba română unde accentul nu depinde de reguli predefinite (în imagine se pot vedea scăderile bruște ale ratelor de eroare de la evaluarea PH+ACC+SYLL (culoarea crem) la evaluarea PH+SYLL (culoarea cărămiziu, atunci când nu se evaluează

acentul)). Un rezultat similar a fost prezentat în (Stan și Giurgiu, 2018). Așa cum ne așteptam, în cele mai multe etape de evaluare, informația despre partea de vorbire și eticheta morfosintactică crește performanța modelului.

În Tabelul 2 se poate observa o diferență foarte mare între ratele de eroare folosind subsetul de antrenare *LEMMA* și cele folosind subsetul de antrenare *1-FORM*. De asemenea, $E(\text{ror})R(\text{ates})$ -urile subsetului *LEMMA* sunt de două ori mai mari decât ER -urile obținute cu subsetul de 5000 de intrări selectate aleator. Ne explicăm scăderea substanțială a ratelor de eroare în cazul subsetului *1-FORM* față de subsetul *LEMMA*, deoarece, așa cum spuneam atunci când am descris procesul de selecție a subset-urilor am avut grijă să ne asigurăm că păstrăm diversitatea formelor morfologice, chiar dacă am selectat o singură formă pentru fiecare leamă, în subsetul *1-FORM*. În contrast, subset-ul *LEMMA* conține doar formele de dicționar și nu învață nici o terminație morfologică. Un subset de intrări selectate aleator de 7 ori mai mic dublează probabilitatea rețelei de a învăța proprietățile terminațiilor morfologice. În continuare, în cazul subsetului *2-FORMS*, observăm că rata de eroare continuă să scadă semnificativ. Dacă comparăm performanțele acestor subseturi cu cele ale căror intrări sunt selectate aleatoriu, observăm că: *1-FORM* (30.150 intrări), pentru scenariul *wMSD*, are o $CER=4,29\%$ pentru predicție concurentă, comparabilă cu $CER=4,36\%$ a subsetului de 100.000 de intrări în scenariul *wMSD*; *2-FORM* (55.185 intrări), pentru scenariul *wMSD*, are o $CER=3,16\%$, aproape de $CER=3,08$ obținută atunci când se folosește întregul set de date.

Aceste performanțe demonstrează că o selecție strategică a intrărilor reduce foarte mult munca de corectură manuală într-un scenariu de construire incrementală - cu adnotare automată și corectare manuală pe subset-uri de date - a unui lexicon. Deși se observă și aici, în general, o îmbunătățire a performanței legată de adăugarea informației *POS* și *MSD*, acest lucru nu este valabil pentru subsetul *LEMMA*.

	LEMMA		1-FORM		2-FORMS		All Forms	
	WER	CER	WER	CER	WER	CER	WER	CER
woutPOS								
PH+ACC+SYLL	46.53%	19.03%	13.30%	5.00%	10.20%	3.86%	5.60	1.97
PH+ACC	45.02%	16.01%	12.17%	3.57%	9.28%	2.64%	5.13	1.47
PH+SYLL	22.31%	7.29%	5.47%	2.37%	4.15%	1.83%	2.20	0.80
PH	16.53%	5.30%	3.91%	1.36%	2.86%	0.86%	1.49	0.42
wPOS								
PH+ACC+SYLL	52.55%	20.69%	11.54%	4.49%	8.39%	3.24%	4.34	1.59
PH+ACC	50.87%	18.18%	10.42%	2.93%	7.49%	2.08%	3.88	1.09
PH+SYLL	27.82%	8.92%	5.71%	2.55%	4.21%	1.81%	2.03	0.77
PH	20.87%	6.50%	4.21%	1.50%	3.07%	0.96%	1.44	0.39
wMSD								
PH+ACC+SYLL	47.57%	18.38%	11.07%	4.29%	8.18%	3.16%	3.08	1.08
PH+ACC	46.35%	15.62%	9.86%	2.82%	7.21%	1.95%	2.57	0.65
PH+SYLL	24.01%	7.26%	5.46%	2.42%	4.13%	1.83%	1.80	0.66
PH	20.51%	6.90%	4.04%	1.46%	2.98%	0.95%	1.24	0.36

Tabelul 2. CER și WER pentru diferite dimensiuni ale datelor de antrenare selectate în modalitatea descrisă pentru scenariul 2, cu subset-urile *1-FORM* (30.150 intrări), *2-FORMS* (55.185 intrări) și *LEMMA* (35.890 intrări) evaluate pentru: predicție lexicală completă (PH+ACC+SYL), eliminând din predicție silabificare (PH+ACC), eliminând din predicție silabificarea și accentul (PH). Rezultatele indică ratele de eroare cu (wPOS) sau fără (woutPOS) etichete morfosintactice adăugate la intrarea în transformer; pentru comparație, pe ultima coloană reluăm rezultatele pentru lexiconul complet reprezentate în Figura 2.

Aici, informația POS și cea MSD nu poate compensa absența variantelor morfologice, pentru că nu are ocazia să învețe nici o asociere între terminații morfologice și etichete POS ale cuvintelor conținut (care reprezintă majoritatea covârșitoare a lexicului unei limbi), pe care să o aplice apoi cuvintelor noi (de exemplu: 1. la transcrierea fonetică a terminațiilor, pentru care există reguli de transcriere specială a sunetului “i” final, și a sunetelor “ce/ci/ge/gi/che/chi/ghe/ghi” la sfârșit de cuvânt; 2. la nivel de silabificare, rețeaua nu are ocazia să învețe despărțirea în silabe a terminațiilor morfologice). Dimpotrivă, faptul că ratele de eroare wPOS și wMSD sunt în acest caz mai mari decât ratele de eroare woutPOS ne face să presupunem că restricția de a prezice informație lexicală coerentă cu partea de vorbire a cuvântului (adică, de exemplu, pentru un verb, rețeaua decide în funcție de ce a învățat doar despre verbele din datele de antrenare) este de fapt nefericită, iar o predicție bazată doar pe forma morfologică este de preferat. La nivel de accent, în fișierul de antrenare din scenariul LEMMA, 18.981 din 35.890 intrări (mai mult de jumătate) au accent pe ultima silabă (spre deosebire de scenariile 1-FORM (6.329 din 30.150) și 2-FORMS (11.966 din 55.185)). Dintre acestea, verbe care au accent pe ultima silabă sunt 6.421 din 7.158. Asta înseamnă că rețeaua transformer învață că probabilitatea ca accentul unui cuvânt să fie pe ultima silabă, în general, este mare, iar pentru verbe, aceasta este foarte mare. Pentru variantele morfologice din datele de test, acest comportament al rețelei este dezavantajos: de exemplu, pentru “a mânca”, accentul cade pe ultima silabă doar pentru 20 din cele 30 de forme morfologice ale conjugării.

Procentul de cuvinte cu accent pe ultima silabă din lista de erori aferentă scenariului wPOS (24.125 intrări din 42.510, dintre care pentru 10.835 eroarea provine doar de la acest accent pe ultima silabă) este chiar mai mare decât cel din lista de erori aferentă scenariului wout POS (19.366 intrări din 37.806, dintre care pentru 10.183 eroarea provine doar de la accentul pe ultima silabă). Analizând în detaliu erorile specifice apărute cu subset-ul LEMMA wPOS vs LEMMA woutPOS, am observat că, din cele 12.085 de erori obținute cu LEMMA wPOS care nu se găsesc în lista de erori obținute cu LEMMA woutPOS, 7.365 de erori sunt asociate verbelor. În timp ce numărul erorilor distincte asociat celorlalte părți de vorbire rămâne relativ același, numărul de erori asociate verbelor crește foarte mult (de mai mult de 5 ori), când trecem de la scenariul woutPOS la cel wPOS. Aceste rezultate demonstrează că preponderența cuvintelor cu accent pe ultima silabă din setul LEMMA produce un procent mare de erori în scenariile wPOS și wMSD. Așa cum spuneam anterior, restricția cuvântului la clasa părții de vorbire căruia îi aparține - în cazul acesta faptul că forțăm rețeaua să trateze verbele din datele de test asemănător verbelor din datele de antrenare - conduce la erori care ar fi putut fi evitate dacă rețeaua se uita doar la forma cuvântului.

Activitatea 3.7: Definitivarea, testarea, validarea și împachetarea într-o soluție „ready-to-use” a platformei integrate și configurabile de prelucrare a textelor în limba română

1. Platforma de prelucrare a textelor TEPROLIN a fost îmbunătățită după cum urmează:

- *Am introdus dependențe de tip graf între operațiile de prelucrare a textelor.* Acest tip de a preciza ce operații trebuie rulate mai întâi pentru a se putea rula operația de prelucrare dorită e mult mai eficient decât tipul de rulare în secvență pe care se baza platforma. De exemplu, pentru a putea rula operația de adnotare cu etichete morfo-sintactice, nu mai sunt necesare operații precum silabificarea sau detecția accentului. Modulul Python 3 în care sunt precizate aceste operații este TeproAlgo.py iar metoda se numește `_assignAlgorithmsToOperations()`.
- *Am adăugat modulul de prelucrare a textelor UDPipe (<http://ufal.mff.cuni.cz/udpipe/1>)* ca alternativă la TTL și NLPCube. Este foarte rapid (cea mai rapidă componentă de prelucrare din cele trei) și are performanțe bune. A fost configurat ca modul implicit dacă se solicită operații precum adnotare cu etichete morfo-sintactice sau analiză cu relații de dependență sintactică.
- *Modulul de inserare a diacriticelor DiacRestore.py a fost îmbunătățit* pentru a detecta mai bine când un text este scris fără diacritice sau cu puține diacritice și a rula, astfel, automat pentru a insera diacriticele lipsă.
- *Am eliminat numărarea caracterelor din modulul de statistici* pentru că frecvența caracterelor prelucrate de TEPROLIN putea crește foarte mult când se prelucrau texte de zeci de milioane de cuvinte. Au rămas statisticile despre numărul de accesări la serviciu și despre numărul de cuvinte prelucrate pe zi.

Păiș et al., (2020) descriu un studiu de caz în care TEPROLIN rulează pe mai multe fire de execuție în RELATE și adnotează corpusul legislativ din proiectul MARCELL. În total, au fost adnotați aprox. 456 de milioane de tokeni, ceea ce demonstrează că testarea și validarea platformei s-au încheiat cu succes.

2. Soluția „ready-to-use” a platformei TEPROLIN

TEPROLIN se poate utiliza într-unul din următoarele patru moduri:

- Pentru testare cu fraze scurte (pentru evaluarea performanțelor) cu efectuarea tuturor operațiilor disponibile, se poate accesa link-ul <https://relate.racai.ro/index.php?path=teprolin/complete> și se pot vizualiza adnotările făcute;
- Pentru rularea unor operații la alegere pe fraze scurte, folosind algoritmi preferați, se poate accesa link-ul <https://relate.racai.ro/index.php?path=teprolin/custom>;
- Pentru adnotarea corpusurilor cu mai mult de 1000 de cuvinte, se poate solicita acces la platforma RELATE care rulează TEPROLIN pe mai multe fire de execuție;
- Ca modul Python 3, clonând repository-ul <https://gitlab.com/raduion/teprolin> și urmând indicațiile din fișierul README.md. Recomandăm ca toate pachetele necesare să fie instalate într-un mediu dedicat Python 3 (eng. „virtual environment”), executând comenzile:
 - a. `python3 -m venv /calea/către/mediul/dedicat/teprolin`
 - b. `pip3 install -r requirements.txt`

Concluzii

Am arătat că folosirea resursei noi dezvoltate în ReTeRom aduce, în mod evident, îmbunătățiri unui sistem de TTS, în comparație cu resursele utilizate anterior. Am presupus și demonstrat că, atunci când ne interesează să evaluăm influența pe care cantitatea de date o are asupra ratelor de eroare, selectarea datelor în mod strategic și nu aleatoriu este esențială. De exemplu, în procesul de corectare manuală a unui lexicon de tipul celui pe care l-am construit noi, se poate reduce foarte mult timpul de corectare dacă selectăm un set esențial de intrări pe care să le corectăm și pe care să antrenăm instrumentul de adnotare automată pentru a adnota intrările rămase, care ne așteptăm să fie mai corecte dacă am ales un set de antrenare potrivit. În cazul nostru, vezi diferența între setul de antrenare *1-FORM*, setul de antrenare *LEMMA* și setul de antrenare selectat aleatoriu, de dimensiuni similare, dar producând efecte foarte diferite asupra rețelei la antrenare. Recomandăm deci o selecție a datelor conform scenariului 2 de evaluare atunci când este importantă reducerea set-ului de date de antrenare. În plus, după cum s-a observat, prin acest tip de selecție putem observa și o descreștere a ratelor de eroare atunci când sunt adăugate informații morfo-sintactice.

Platforma de prelucrare a textelor TEPROLIN se află la <https://gitlab.com/raduion/teprolin>. Pentru a avea acces, trebuie să aveți cont pe GitLab și să solicitați accesul autorului platformei.

Toate obiectivele incluse în plan la aceste activități au fost realizate.

Bibliografie

- Erjavec, Tomaz. "MULTEXT-East Morphosyntactic Specifications: Version 3.0." Supported By EU Projects Multext-East, Concede And TELRI (2004)
- Gehring, J., M. Auli, D. Grangier, D. Yarats, and Y. N. Dauphin, "Convolutional sequence to sequence learning," CoRR, vol. abs/1705.03122, 2017.
- V. Levenshtein, "Binary Codes Capable of Correcting Deletions, Insertions and Reversals," Soviet Physics Doklady, vol. 10, p. 707, 1966.
- Stan, A. and M. Giurgiu, "A Comparison Between Traditional Machine Learning Approaches And Deep Neural Networks For Text Processing In Romanian," in Proceedings of the 13th International Conference on Linguistic Resources and Tools for Processing Romanian Language (ConsILR), 2018.
- Stan, Adriana. "RECOApy: Data recording, pre-processing and phonetic transcription for end-to-end speech-based applications," in Proceedings of Interspeech, Shanghai, China, 2020.
- Vaswani, A., N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in Advances in Neural Information Processing Systems 30, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. Curran Associates, Inc., 2017, pp. 5998–6008.
- Păiș, V., Tufiș, D. și Ion, R. (2020) A Processing Platform Relating Data and Tools for Romanian Language. In Proceedings of the 12th Language Resources and Evaluation Conference (LREC), European Language Resources Association, Marseille, France, pages 81—88.

Raport științific - tehnic proiect component TADARAV

Date de identificare proiect	
Număr contract:	PN-III-P1-1.2-PCCDI-2017-0818, Nr. 73 PCCDI/2018
Acronim / titlu:	„TADARAV:” Tehnologii pentru adnotarea automată a datelor audio și pentru realizarea interfețelor de recunoaștere automată a vorbirii
Titlu livrabil:	Raport științific și tehnic (Etapa a III-a, 2020)
Termen:	Noiembrie 2020
Editor:	Horia Cucu (Universitatea Politehnica din București)
Adresa de eMail editor:	horia.cucu@upb.ro
Autori, în ordine alfabetică:	Alexandru-Lucian Georgescu, Cristian Manolache, Dan Oneață, Gheorghe Pop, Horia Cucu, Corneliu Burileanu, Dragoș Burileanu
Ofițer de proiect:	Cristian Stroe

1. Activitățile etapei de raportare în contextul general al proiectului

Activitățile realizate în etapa 3/2020 au fost următoarele:

- **Activitatea 3.9 - Analiza impactului utilizării de RAV complementare pentru generarea de adnotări în contextul îmbunătățirii sistemelor de RAV.** Această activitate a presupus un studiu comparativ asupra performanțelor sistemului de RAV inițial și a sistemelor de RAV rezultate în cadrul activităților A1.13/2018 și A2.13/2019. De asemenea, în cadrul acestei activități s-a încercat aplicarea metodei proiectate în activitățile anterioare cu un sistem de RAV complet nou, bazat pe platforma ESPnet. Sistemul de RAV rezultat s-a dovedit a fi extrem de lent în transcriere, astfel că nu a putut fi utilizat corespunzător. Cu toate acestea, per ansamblu, sistemul de RAV inițial (disponibil la începutul proiectului) a putut fi îmbunătățit folosind metoda proiectată în activitățile A1.13/2018 și A2.13/2019 cu aproximativ 50% atât pe vorbire citită, cât și pe vorbire spontană.
- **Activitățile 3.10 Îmbunătățirea soluției de filtrare și aliniere a transcrierilor aproximative cu semnalul de vorbire și 3.12 Analiza impactului utilizării transcrierilor aproximative în vederea reantrenării sistemelor de RAV.** Activitățile A3.10 și A3.12 au avut ca scop îmbunătățirea și evaluarea metodei de generare de seturi de date de vorbire adnotată folosind materiale audio brute împreună cu transcrieri aproximative. Ideea principală a acestei metode este următoarea: un sistem de RAV inițial este folosit pentru a genera transcrieri pentru materialul audio brut, iar ulterior aceste transcrieri sunt aliniate cu textele aproximative deja existente. Părțile aliniate sunt considerate corecte și sunt folosite pentru reantrenarea sistemului inițial de RAV. În urma realizării activității A2.11 din etapa anterioară am tras concluzia că între două secvențe de text aliniate corect există transcrieri aproximative care sunt în mare proporție corecte. Îmbunătățirea propusă și evaluată în cadrul activităților din anul 2020 constă în utilizarea acestor transcrieri aproximative aflate între două secvențe de text aliniate corect. Rezultatele experimentale au arătat că îmbunătățirea RAV rezultată este nesemnificativă (mai mică de 1%).
- **Activitățile 3.11 Îmbunătățirea soluției pentru generarea de scoruri de încredere pentru RAV și 3.13 Analiza impactului utilizării scorurilor de încredere pentru filtrarea transcrierilor RAV în vederea reantrenării sistemelor RAV.** În cadrul acestor activități au fost propuse o serie de noi metode de generare de scoruri de încredere pentru RAV folosind o platformă modernă de tip end-to-end: ESPnet. Metodele propuse au fost evaluate și comparate, iar metoda cea mai performantă a fost selectată pentru utilizarea ulterioară în vederea producerii de transcrieri precise pentru o parte a semnalului de vorbire. Sistemul de RAV ESPnet s-a dovedit însă a fi prea lent pentru a putea fi utilizat într-o astfel de metodă de adnotare automată de date: cu un factor de timp real de aproximativ 3, transcrierea celor 913 de ore de vorbire brută ar fi durat aproximativ 100 de zile. În consecință, noile metode de estimare a scorurilor de încredere nu au putut fi evaluate în contextul reantrenării sistemului de RAV inițial, folosind datele adnotate rezultate.

2. Gradul de realizare a obiectivelor specifice pentru Etapa a III-a, 2020

A treia etapă a proiectului TADARAV a avut trei obiective principale ce au fost realizate în proporție de 100%:

- evaluarea globală, la nivelul întregului proiect TADARAV, a metodei ce presupune utilizarea sistemelor de recunoaștere automată a vorbirii (RAV) complementare pentru generarea automată de adnotări pentru date audio și, ulterior, utilizarea acestor adnotări pentru reantrenarea unui nou sistem de RAV.
- îmbunătățirea manierei de utilizare a transcrierilor aproximative ale materialelor ce conțin vorbire, împreună cu un sistem de RAV inițial, pentru a produce în mod automat transcrieri precise pentru o parte a semnalului de vorbire;
- propunerea de noi metode de generare de scoruri de încredere pentru RAV și utilizarea scorurilor de încredere rezultate pentru a produce în mod automat transcrieri precise pentru o parte a semnalului de vorbire.

3. Rezultatele etapei și descrierea lor științifică și tehnică

Rezultatele etapei 3/2020 a proiectului TADARAV sunt listate mai jos. Primele două rezultate din lista de mai sus sunt prezentate pe larg în livrabilele anexate: [Manolache, 2020a], [Oneață, 2020]. Al treilea rezultat, sistemul RAV actualizat este descris în continuare.

- Soluție îmbunătățită de filtrare și aliniere a transcrierilor aproximative cu semnalul de vorbire (TRL4);
- Soluție îmbunătățită pentru generarea de scoruri de încredere pentru RAV (TRL4);
- Sistem de RAV actualizat (TRL5).

3.1 Actualizarea modelelor de limbă pentru transcriere de vorbire

În această secțiune sunt prezentate sumar demersurile realizate de colectivul SpeedD în vederea actualizării modelelor de limbă utilizate în sistemele de transcriere de vorbire. Etapele parcurse pentru această actualizare sunt descrise în detaliu în [Manolache, 2020b].

Primul pas întreprins a fost colectarea unui nou corpus de text. Acest corpus nou de text, denumit mai departe news2020, constă în text brut, organizat pe propoziții ce însumează aproximativ 255M de cuvinte. Cele 24 surse de știri ce formează corpusul news2020 sunt prezentate în Figura T1.

Al doilea pas întreprins a vizat actualizarea aplicației de preprocesare de text în vederea îmbunătățirii modului de tratare a cuvintelor cu cratimă [Manolache, 2020b].

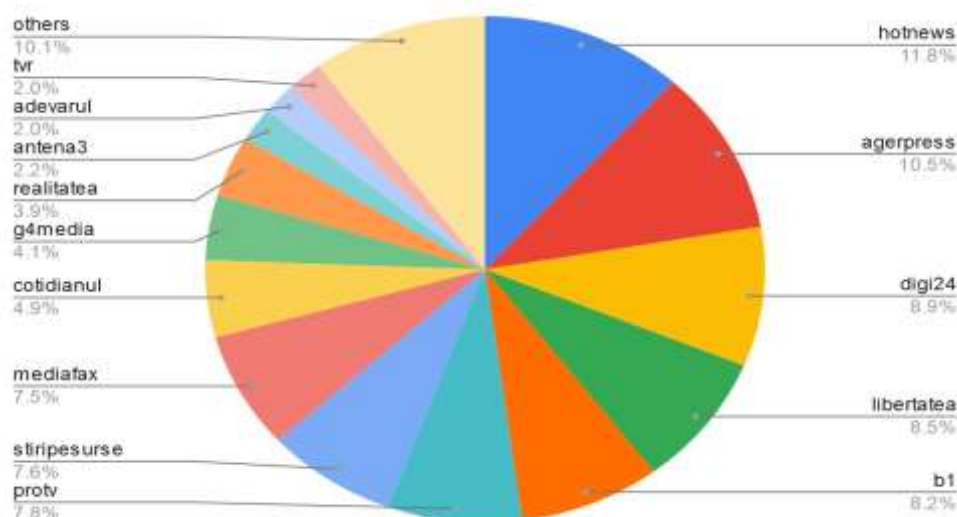


Figura T1 Sursele de știri componente ale corpusului news2020. Procentele sunt relative la dimensiunea întregului corpus (255M cuvinte)

Al treilea pas întreprins a vizat actualizarea modelelor de limbă. Nu au fost create tipuri noi de modele de limbă, ci au fost antrenat modele de limbă folosind corpusuri de text extinse și preprocesate mai bine.

Modelele de limbă de bază au fost antrenate folosind numai corpusul news002 (361M de cuvinte), pe când noile modele de limbă au fost antrenate pe corpusul news002 plus corpusul news2020 (361M + 255M de cuvinte). Textul folosit pentru antrenarea modelelor de limbă noi a fost preprocesat folosind noul procesor de limbaj natural discutat în acest subcapitol. Modelele de limbă îmbunătățite sunt evaluate în funcție de perplexitate (PPL) și cuvinte out-of-vocabulary (OOV). Modelele de limbă mai simple (2-gram și 3-gram) au fost încorporate în etapa de decodare a sistemului ASR, în timp ce modelele de limbă mai complexe (4-gram și RNN) au fost folosite în etapa de rescoring. Modelele mai complexe nu au fost folosite în etapa de decodare din cauza constrângerilor arhitecturale și de memorie.

În final, în al patrulea pas, modelele de limbă nou create au fost încorporate în sistemul de RAV. Pentru etapa de decodare am optat pentru modelele 3-gram 200k, iar pentru rescoring am ales LM-2020-4g-large-200k și LM-2020-5g-RNN. Cele mai bune rezultate s-au obținut pentru sistemul RAV ce folosește LM-2020-RNN pentru rescoring. În Tabelul T1 sunt prezentate rezultatele WER ale noilor sisteme RAV cu și fără rescoring.

Tabelul T1 Rezultatele WER ale sistemelor RAV după decodare și rescoring

Model de limbă folosit pentru rescoring	Model de limbă folosit pentru decodare	WER[%] fără rescoring			WER[%] cu rescoring		
		RSC-eval	SSC-eval1	SSC-eval2	RSC-eval	SSC-eval1	SSC-eval2
LM-2017-RNN	LM-2017-2g-large-200k	2.8	13.3	16.4	1.8	11.0	14.0
LM-2020-RNN	LM-2020-3g-large-200k	2.7	11.0	14.4	2.1	10.0	12.9
	LM-2020-3g-medium-200k	3.1	11.9	15.1	2.2	10.4	13.3
	LM-2020-3g-small-200k	3.3	12.6	15.8	2.2	10.5	13.7

În Tabelul T1 se poate observa că cel mai bun sistem RAV este cel ce folosește LM-2020-3g-large-200k cu LM-2020-RNN pentru rescoring. Acest model are o scădere relativă de 9% și 7% WER pe seturile de evaluare de vorbire spontană, anume SSC-eval1 și respectiv SSC-eval2, comparativ cu sistemul RAV de bază cu rescoring (Tabel T1 linia 1). Însă, pe setul de evaluare RSC-eval sistemul RAV a înregistrat o creștere de 16% WER.

3.2 Utilizarea seturilor de date rezultate din proiect pentru antrenarea RAV

În această secțiune analizăm impactul utilizării unor seturi de date de antrenare mai mari asupra performanței sistemelor de RAV Speed bazate pe utilitarul Kaldi. Pentru a putea compara direct îmbunătățirile rezultate strict din creșterea numărului de ore de antrenare, toți ceilalți hiperparametri ai sistemelor au fost păstrați constanți. Printre altele, toate sistemele prezentate în această secțiune folosesc pentru decodare de vorbire modelul LM-2017-2g-large-200k, iar pentru reevaluare lingvistică modelul LM-2017-RNN.

Tabelul T2 prezintă evoluția sistemului RAV Speed în perioada 2018-2020. Fiecare linie reprezintă un model acustic diferit, îmbunătățit față de cel precedent, pentru antrenarea căruia au fost adăugate seturi de date suplimentare. Fiecare sistem este evaluat pe cele 3 seturi de evaluare existente: RSC-eval, SSC-eval1, SSC-eval2. Primul model acustic, denumit în tabel *train-base*, este cel din cadrul sistemului inițial RAV #3, prezentat în secțiunea 2.2.2 a raportului [Georgescu, 2020b]. Acesta a fost antrenat folosind setul de antrenare *baseline*: RSC-train1 + SSC-train1 + SSC-train2.

Cu ajutorul sistemului inițial RAV #3, au fost decodate seturile SSC-train3-raw și SSC-train4-raw, iar transcrierile obținute, împreună cu transcrierile aproximative deja existente, au fost aliniate în etapa A2.11/2019, rezultând seturile SSC-train3-trans-v4 și SSC-train4-trans-v4. Aceste seturi au fost adăugate la setul de antrenare *baseline*, adunând 292 ore de vorbire la cele 225 ore anterioare. Prin reantrenarea sistemului, s-a obținut astfel cel de-al doilea model acustic, denumit *train3-v4*. Acest fapt a dus la o scădere relativă a ratei de eroare la nivel de cuvânt de 5% pe RSC-eval, respectiv 26%-30% pe SSC-eval1 și SSC-eval2. Faptul că pe vorbire spontană s-a obținut o îmbunătățire mai mare, se datorează caracteristicilor noilor seturi de antrenare acestea cuprinzând preponderent vorbire spontană.

Sistemul RAV #3, corespunzător modelului acustic *train-base*, a fost utilizat pentru a transcrie setul de date CoBiLiRo, prezentat în secțiunea 2.1.2 a raportului [Georgescu, 2020b]. În secțiunea 2.3.2 a raportului

[Georgescu, 2020b] este prezentat modul în care a fost aliniată transcrierea obținută cu transcrierile aproximative deja existente, rezultând în acest fel setul de date CoBiLiRo-trans-v4. Acesta a fost adăugat la seturile de antrenare anterioare, fiind obținut un nou model acustic, denumit *train14*. În comparație cu modelul *train3-v4*, modelul *train14* obține rezultate mai bune pe vorbire citită (11% scădere relativă a WER pe RSC-eval), dar rezultate mai slabe pe vorbire spontană. Acest lucru se datorează similitudinii dintre CoBiLiRo-trans-v4 și setul de date RSC. CoBiLiRo-trans-v4 conține interviuri și dialoguri într-un mediu fără zgomot. Totodată, trebuie ținut cont și de cantitatea mică a datelor nou adăugate, 31 ore, față de cele 517 deja existente.

Ultimul model acustic, *train11*, a presupus reantrenarea sistemului RAV după ce a fost adăugat setul de date CoRoLa, introdus în secțiunea 2.1.1 a raportului [Georgescu, 2020b]. Îmbunătățirile relative față de sistemul anterior sunt cuprinse între 8%-12% pe vorbire spontană, în timp ce în cazul vorbirii citite, nu se înregistrează schimbări. Acest lucru se explică prin conținutul setului CoRoLa: vorbire spontană, de la vorbitori profesioniști, în cadrul unor emisiuni TV. În total, față de modelul *train-base*, acest ultim model obține o scădere relativă a erorii de 15% pe vorbire citită, respectiv 31%-39% pe vorbire spontană.

În continuare ne-am pus problema contribuția fiecărui nou set de date de vorbire la îmbunătățirea totală a sistemului de RAV. Astfel, am antrenat sisteme de RAV pe seturi de date formate din setul baseline și un singur alt set de date de vorbire. Tabelul T3 prezintă rezultatele. Rezultatele obținute pe setul de date RSC-eval diferă foarte puțin, dacă nu chiar deloc, între diversele experimente. Tragem concluzia că modelul din adăugarea de date noi numai pentru creșterea performanței pe vorbire spontană. Deoarece seturile nou adăugate conțin în mare măsură vorbire spontană, îmbunătățirile substanțiale au rezultat pe seturile de evaluare SSC-eval1 și SSC-eval2 (ce conțin astfel de vorbire). Este foarte probabil ca în urma reantrenărilor, modificări ale parametrilor rețelei neuronale să fi apărut numai în acele zone din rețea care se ocupă de prelucrarea vorbirii spontane.

O alta concluzie se poate trage cu privire la dimensiunea seturilor noi de date. De exemplu, chiar dacă modelul *train3-v4* conține mai mult decât dublul datelor folosite la antrenarea lui *train-base*, îmbunătățirile obținute nu sunt de două ori mai mari. Nu există o dependență direct proporțională între cantitatea de date folosite la antrenare și performanța sistemului. Cu cât acuratețea crește, este nevoie de mult mai multe date pentru a avea parte de îmbunătățiri mici.

Tabelul T2. Evoluția RAV Speed în perioada 2018 - 2020. Îmbunătățirile rezultate prin adăugarea la setul de date de antrenare BAS (RSC-train + SSC-train1 + SSC-train2) a seturilor de date rezultate din proiect SSC (SSC-train3-trans-v4 + SSC-train4-trans-v4), COB (CoBiLiRo-trans-v4) și COR (CoRoLa)

Cod model acustic	Set de antrenare					Set de evaluare (WER[%])		
	BAS 225h	SSC 292h	COB 31h	COR 84h	CDP 879h	RSC-eval	SSC-eval1	SSC-eval2
train-base	x					1.9	15.0	20.0
train3-v4	x	x				1.8	11.0	14.0
train14	x	x	x			1.6	11.3	14.4
train11	x	x	x	x		1.6	10.3	12.2

Tabelul T3. Contribuția fiecărui set de date nou la îmbunătățirea RAV Speed 2018 - 2020. Seturile de date de antrenare folosite sunt: BAS (RSC-train + SSC-train1 + SSC-train2), SSC (SSC-train3-trans-v4 + SSC-train4-trans-v4), COB (CoBiLiRo-trans-v4) și COR (CoRoLa)

Cod model acustic	Set de antrenare					Set de evaluare (WER[%])		
	BAS 225h	SSC 292h	COB 31h	COR 84h	CDP 879h	RSC-eval	SSC-eval1	SSC-eval2
train-base	x					1.9	15.0	20.0
train3-v4	x	x				1.8	11.0	14.0

train15	x		x			1.8	14.0	21.1
train17	x			x		1.8	11.9	15.4

Modelul acustic *train15* prezintă influența adăugării datelor din setul CoBiLiRo-trans-v4. Dată fiind cantitatea redusă a datelor din acest set, în comparație cu datele folosite la antrenarea modelului *train-base*, performanțele sistemului nou nu sunt cu mult diferite față de cel inițial.

Modelul acustic *train17*, antrenat prin adăugarea setului de date CoRoLa, însumând 84 de ore, ajunge să obțină performanțe asemănătoare cu modelul *train3-v4*, ce totalizează 292 ore. Deși ambele seturi, atât CoRoLa, cât și SSC-train3-trans-v4 + SSC-train4-trans-v4, conțin vorbire spontană, diferența este dată de lungimea rostirilor din cele două seturi. CoRoLa a fost adnotat manual, iar rostirile au dimensiuni mai mari decât în cazul lui SSC-train3-trans-v4 + SSC-train4-trans-v4, care a fost adnotat automat, prin metoda alinierii transcrierilor aproximative, unde au fost selectate doar segmente comune de lungimea a câtorva cuvinte. Rețeaua neuronală folosită la antrenarea modelului acustic beneficiază în prima situație de un context temporal mult mai larg, în timp ce în a doua situație, rețeaua are un context temporal insuficient pentru a învăța interdependențele fonemelor.

Tabelul T4 prezintă contribuția noilor seturi de date introduse în 2020 în mod comparativ. Unul din cele 3 seturi noi din 2020, pe lângă CoBiLiRo-trans-v4 și CoRoLa, este cdep-trans-v4, prezentat în secțiunea 2.3.3 a raportului [Georgescu, 2020b]. Setul de date cdep-raw a fost transcris folosind modelul acustic *train3-v4*, iar apoi, prin alinierea cu transcrierile aproximative deja existente, a rezultat setul cdep-trans-v4. Surprinzător este faptul că deși setul de date cdep-trans-v4 însumează 879 ore, fiind de 28 de ori mai mare decât CoBiLiRo-trans-v4, respectiv de 10 ori mai mare decât CoRoLa, adăugarea acestui set la antrenare nu produce rezultate mai bune decât adăugarea celorlalte două seturi amintite. O pistă de investigat în acest sens este lungimea rostirilor din setul cdep-trans-v4, fiind posibil ca acestea să fie foarte scurte (e posibil ca transcrierile aproximative și transcrierile date de sistemul RAV să se suprapună într-o mică măsură).

Tabelul T4. Contribuția noilor seturi de date introduse în 2020. Sistemul inițial a fost antrenat pe seturile de date BAS (RSC-train + SSC-train1 + SSC-train2) și SSC (SSC-train3-trans-v4 + SSC-train4-trans-v4). În etapa 3/2020 au fost introduse seturile de date COB (CoBiLiRo-trans-v4), COR (CoRoLa) și CDP (cdep-trans-v4)

Cod model acustic	Set de antrenare					Set de evaluare (WER[%])		
	BAS 225h	SSC 292h	COB 31h	COR 84h	CDP 879h	RSC- eval	SSC- eval1	SSC- eval2
train-base	x					1.9	15.0	20.0
train3-v4	x	x				1.8	11.0	14.0
train14	x	x	x			1.6	11.3	14.4
train18	x	x		x		1.8	10.5	12.5
train13	x	x			x	1.7	12.3	15.4

3.3 Sistemul de RAV Speed îmbunătățit

În concluzie, selectând cea mai performantă configurație de model de limbă pentru decodare, model de limbă pentru reevaluare lingvistică (conform secțiunii 3.1) și model acustic (conform secțiunii 3.2), cel mai performant sistem de RAV Speed are actualmente performanțele listate în Tabelul T5.

Tabelul T5. Sistemul RAV Speed 2020

Model acustic		Model lingvistic		WER [%]			
Set de antrenare	Tip model	Corpus antrenare	Tip model	RSC- eval	SSC- eval1	SSC- eval2	CDep- eval
RSC-train + SSC-train +SSC-train3-trans-v4 +SSC-train4-trans-v4 +CoBiLiRo-trans-v4 +CoRoLa	HMM-DNN (TDNN3) Kaldi toolkit	news2017 +news2020	Decod.: LM-2020-3g-large-200k Reev. lingv.: LM-2020-RNN	1.9	9.4	11.4	7.0

4. Oferta de servicii de cercetare, locuri de muncă susținute și valorificarea resurselor

Laboratorul de cercetare *Speech and Dialogue* (Speed) din cadrul Universității Politehnica din București (UPB), reprezentantul UPB în proiectul TADARAV, oferă pe platforma ERRIS (<https://erris.gov.ro/Speed--UPB>) serviciile de cercetare și tehnologice enumerate în Tabelul T6.

Tabelul T6. Servicii de cercetare și tehnologice oferite de Laboratorul de cercetare *Speech and Dialogue*

Serviciu	Detalii
Serviciu și aplicație web de transcriere de documente ce conțin vorbire în limba română	https://transcriptions.speed.pub.ro
Serviciu și aplicație web de identificare de cuvinte cheie în documente ce conțin vorbire în limba română	https://keywords.speed.pub.ro
Serviciu și aplicație web de restaurare de diacritice în limba română	https://diacritics.speed.pub.ro
Proiectarea și implementarea de aplicații personalizate de transcriere a vorbirii continue	La cerere
Proiectarea și implementarea de aplicații personalizate de identificare de cuvinte și termeni de interes	La cerere
Proiectarea și implementarea de aplicații personalizate de sinteză de vorbire pornind de la text	La cerere
Proiectarea și implementarea de sisteme de recunoaștere de pattern-uri folosind inteligență artificială	La cerere

Echipa de cercetare a Universității Politehnica din București pentru proiectul component TADARAV este prezentată în Tabelul T7.

Tabelul T7. Echipa de cercetare UPB

Nr.	Nume	Calitatea	Poziția	Normă
1	Horia CUCU	Conf. Univ.	Responsabil proiect component	Parțială
2	Corneliu BURILEANU	Prof. Univ.	Membru cercetător	Parțială
3	Dragoș BURILEANU	Prof. Univ.	Membru cercetător	Parțială
4	Alexandru-Lucian GEORGESCU	ACS	Membru cercetător	Parțială
5	Dan Theodor ONEAȚĂ	CS	Membru cercetător nou	Întreagă
6	Gheorghe POP	ACS	Membru cercetător nou	Întreagă
7	Cristian MANOLACHE	ACS	Membru cercetător nou	Întreagă

În această etapă proiectul TADARAV nu a avut fonduri la capitolul bugetar CEC-uri.

5. Concluzii

Activitățile de cercetare desfășurate în etapa a III-a de implementare a proiectului (2020) au condus la obținerea rezultatelor așteptate și ele sunt în concordanță cu obiectivele specifice ale etapei. Astfel, rezultatele raportate în acest document și descrise detaliat în cele 3 livrabile aferente perioadei de raportare, asigură testarea și evaluarea finală a tehnologiei dezvoltate în etapa finală a proiectului (2021).

6. Referințe

- [Georgescu, 2020a] A.-L. Georgescu, H. Cucu, A. Buzo, C. Burileanu, “RSC: A Romanian Read Speech Corpus for Automatic Speech Recognition,” in the Proceedings of The 12th Language Resources and Evaluation Conference (LREC), pp. 6606-6612, 2020, Marseille, France.
- [Georgescu, 2020b] A.-L. Georgescu, A. Caranica, C. Manolache, G. Pop, D. Oneață, H. Cucu, C. Burileanu, D. Burileanu, „Proiect component TADARAV: Raport științific și tehnic în extenso 2020”.
- [Georgescu, 2020c] A.-L. Georgescu, A. Caranica, H. Cucu, D. Burileanu, C. Burileanu, „Raport de analiză a impactului utilizării de RAV complementare pentru generarea de adnotări în contextul îmbunătățirii sistemelor de RAV,” Livrabil proiect ReTeRom, 2020.
- [Oneață, 2020] D. Oneață, A. Caranica, H. Cucu, C. Burileanu, D. Burileanu, „Raport de analiză a impactului utilizării scorurilor de încredere pentru filtrarea transcrierilor RAV în vederea reantrenării sistemelor RAV,” Livrabil proiect ReTeRom, 2020.
- [Manolache, 2020a] C. Manolache, A.-L. Georgescu, H. Cucu, C. Burileanu, D. Burileanu, „Raport de analiză a impactului utilizării scorurilor de încredere pentru filtrarea transcrierilor RAV în vederea reantrenării sistemelor RAV,” Livrabil proiect ReTeRom, 2020.
- [Manolache, 2020b] C. Manolache, A.-L. Georgescu, H. Cucu, V. B. Mititelu, C. Burileanu, “Improved text normalization and language models for Speed’s Automatic Speech Recognition System”, in Proc. 15th International Conference “Linguistic Resources and Tools for Processing the Romanian Language” (ConsILR) 2020.

Raport științific - tehnic proiect component SINTERO

Date de identificare proiect

Număr contract:	PN-III-P1-1.2-PCCDI-2017-0818, Nr. 73 PCCDI/2018
Acronim / titlu:	„SINTERO: Tehnologii de realizare a interfețelor om-mașină pentru sinteza text-vorbire cu expresivitate”
Titlu livrabil:	Raport științific și tehnic (Etapa a III-a, 2020)
Termen:	Noiembrie 2020
Editor:	Mircea Giurgiu (Universitatea Tehnică din Cluj-Napoca)
Adresa de eMail editor:	Mircea.Giurgiu@com.utcluj.ro
Autori, în ordine alfabetică:	Mircea Giurgiu, Beata Lorincz, Maria Nuțu, Adriana Stan
Ofițer de proiect:	Cristian Stroe

1. Activitățile etapei de raportare în contextul general al proiectului

În a treia etapă (2020) a proiectului SINTERO „Dezvoltarea unei noi tehnologii de realizare a interfețelor om – mașină pentru sinteza text – vorbire cu expresivitate”, s-au desfășurat activitățile: **a)** dezvoltarea unei noi tehnologii de adaptare a vocii sintetizate la stilul și expresivitatea unui nou vorbitor (Livrabil D3.15), **b)** dezvoltarea unei metode de adaptare rapidă a vocii sintetizate folosind date audio atipice (Livrabil D3.16), **c)** integrarea diferitelor tehnologii într-un flux de prelucrări automatizat (Livrabil D3.17), respectiv activități de testare, validare, diseminare și demonstrare online (Livrabil D3.18).

2. Gradul de realizare a obiectivelor specifice pentru Etapa a III-a, 2020

Ob. Pr4.3.15: *Dezvoltarea unei tehnologii pentru adaptarea vocii sintetizate la stilul și expresivitatea unui vorbitor* (Obiectiv realizat integral). *Rezultate:* ● 1 nou corpus cu 29 de vorbitori și aproximativ 65 de ore de vorbire ● 1 aplicație software (RecoApy) pentru automatizarea înregistrărilor ● 1 vocoder de tip WaveGlow antrenat cu vocile din noul corpus ● 1 sistem de sinteză de tip DC-TTS pentru adaptare la noi vorbitori folosind o nouă metrică pentru calculul erorii de antrenare ● 1 metodă de adaptare de la date ne-expresive la date expresive ● 2 articole publicate la 2 conferințe internaționale ● un livrabil (D3.15) cu titlul „Tehnologie de adaptare a vocilor sintetice la noi vorbitori”.

Ob. Pr4.3.16: *Dezvoltarea unei metode de adaptare rapidă a vocii sintetizate folosind date audio atipice* (Obiectiv realizat integral). *Rezultate:* ● 1 metodă de predicție simultană cu rețele neuronale a informației lexicale ● 1 lexicon RoLEX utilizat pentru antrenarea modelelor de text ● câte 5 sisteme de sinteză text-vorbire folosind informații lexical precise, evaluate atât pentru limba română, cât și engleză ● 1 demonstrator pe bază de Flowtron pentru transferul stilului de vorbire, respectiv pentru antrenarea cu date atipice ● pagini web cu demonstrarea calității semnalului audio sintetizat ● 2 articole la conferințe internaționale ● 1 livrabil (D3.16) cu titlul „Metodă de adaptare rapidă a sistemului de sinteză cu date atipice”.

Ob. Pr4.3.17: *Integrarea tehnologiilor și demonstratoare online* (Obiectiv realizat integral). *Rezultate:* ● integrarea tehnologiilor dezvoltate ● 20 de sisteme de sinteză text vorbire disponibile pentru demonstrare online în interfața //speech.utcluj.ro/ronna ● 1 livrabil (D3.17) cu titlul „Tehnologie de realizare a interfețelor om – mașină pentru sinteza text - vorbire”.

Ob. Pr4.3.18: *Diseminarea rezultatelor* (Obiectiv realizat integral). *Rezultate:* ● realizarea și actualizarea web site-ului ● pagini web cu demonstratoare cu vocile sintetizate ● 1 livrabil referitor la activitățile de diseminare (D3.18).

3. Rezultatele etapei și descrierea lor științifică și tehnică

3.1. Tehnologie de adaptare a vocilor sintetice la noi vorbitori

Adaptarea vorbirii sintetice la stilul și expresivitatea unui nou vorbitor a fost realizată prin intermediul a trei noi metode: 1) adaptarea de la un sistem cu vorbitori multipli la un sistem cu un singur vorbitor folosind un set redus de date, 2) adaptarea de la un sistem cu voce ne-expresivă la un sistem cu voce

expresivă, 3) adaptarea sistemului de sinteză la un nou vorbitor printr-o procedură de adaptare bazată pe post-filtrare. Descriere în extenso în D3.15.

(1) Crearea unui nou corpus (SWARA 2.0) pentru adaptarea sistemului de sinteză la noi vorbitori. Acest corpus a fost înregistrat de către studenți voluntari în condiții ambientale pentru o mai mare variabilitate a datelor audio și folosind aplicația RecoApy². A rezultat un set de înregistrări ce conține 51.839 de segmente audio de la 29 de vorbitori: 14 masculini și 15 feminini și o durată totală a înregistrărilor de aproximativ 65 de ore (vezi D3.15). Noul corpus a fost utilizat pentru a antrena un vocoder de tip *WaveGlow* utilizat împreună cu Tacotron2, precum și o serie de noi voci folosind metoda *Transfer Learning* pentru a valida tehnologia de adaptare a sistemului de sinteză la noi vorbitori. Mostre audio sunt disponibile la adresa: <https://gitlab.utcluj.ro/speech/tts-samples/-/wikis/home>.

(2) Aplicația RecoApy³ permite înregistrarea secvențială a propozițiilor din corpusul SWARA 2.0 cu salvarea unei copii în caz de eroare, respectiv afișarea transcrierii fonetice a textului pe bază de rețele neuronale convoluționale (CNN - Convolutional Neural Networks) și de tip Transformer antrenate cu date din Wiktionary⁴ în opt limbi: română, engleză, spaniolă, franceză, germană, italiană, cehă și poloneză. Performanța transcrierii fonetice cu această aplicație este la nivelul metodelor de dată recentă din domeniu.

(3) Adaptarea unui sistem de sinteză cu vorbitori multipli pe baza unei funcții de cost suplimentare. Pentru îmbunătățirea învățării identității vorbitorului și adaptarea unui sistem de vorbitori multipli la un nou vorbitor, s-a propus folosirea unui sistem de identificare a vorbitorului, combinat cu un sistem de sinteză text-vorbire (sistemul adoptat este DC-TTS), prin adăugarea unei funcții de cost suplimentare bazată pe reprezentările interne ale sistemului de identificare a vorbitorului. Pornind de la dezvoltările software realizate de grupul de cercetare CSTR⁵, s-a implementat rețeaua DC-TTS cu posibilitatea de concatenare a unei reprezentări vectoriale a identității vorbitorului (embedding) la cele 3 module ale componentei Text2Mel, respectiv prin utilizarea reprezentării vectoriale într-o strategie de învățare a contribuției reprezentării vectoriale la canalele de informație din rețea. Astfel, este posibilă inserarea identității vorbitorului în mai multe puncte ale rețelei. Sistemul cu vorbitori multipli a fost antrenat în 3 scenarii:

Scenariul 1 (**B**): sistem cu vorbitor unic sau vorbitori multipli considerat sistemul de bază.

Scenariul 2 (**B+CS**): sistem cu vorbitor unic sau vorbitori multipli și adăugarea unei funcții de cost suplimentare obținută prin calculul similarității spectrale.

Scenariul 3 (**B+E**): sistem cu vorbitor unic sau vorbitori multipli cu adăugarea unei funcții de cost suplimentare obținută cu rata de eroare egală (EER - Equal Error Rate) folosind sistemul de verificare de vorbitor. Acest sistem a fost antrenat cu corpusul SWARA folosind o implementare⁶ de la Clova AI Research. Cele 3 scenarii au fost antrenate cu diferite seturi de date:

- **ALL** - 18 vorbitori, între 1.000 și 1.500 de propoziții / vorbitor, 21.302 pronunții.
- **RND1** – 18 vorbitori, aproximativ 500 de pronunții / vorbitor, 8.932 de pronunții.
- **RND1-100** – 18 vorbitori, aproximativ 100 de pronunții / vorbitor, 1.787 de pronunții.
- **RND1-SAM** – 1 vorbitor, în total 500 de pronunții.

Evaluare. Sistemele au fost evaluate obiectiv calculând valoarea EER. Fiecare semnal sintetizat este comparat cu o mostră naturală de la un același vorbitor și de la un alt vorbitor ales aleator. Tabelul 1 prezintă rezultatele obținute. De asemenea, rezultatele obținute prin aplicarea metodei de codare a vorbitorilor prin embedding pot fi vizualizate cu ajutorul metodei de t-Distributed Stochastic Neighbour Embedding (t-SNE) (Fig 1, 2).

Tabel 1. EER pentru sistemele Baseline, CosSim și EER, antrenate cu diferite seturi de date

² <https://gitlab.utcluj.ro/sadriana/recoapy/>

³ https://www.isca-speech.org/archive/Interspeech_2020/pdfs/1184.pdf

⁴ <https://www.wiktionary.org/>

⁵ <http://www.cstr.ed.ac.uk/>

⁶ https://github.com/clovaai/voxceleb_trainer

Sistem	ALL (EER)	RND1 (EER)	RND1-100 (EER)	RND1-SAM (EER)
<i>B</i>	6.94	4.86	8.33	2.43
<i>B+CS</i>	6.25	4.66	6.25	2.43
<i>B+E</i>	4.66	8	6	2.43

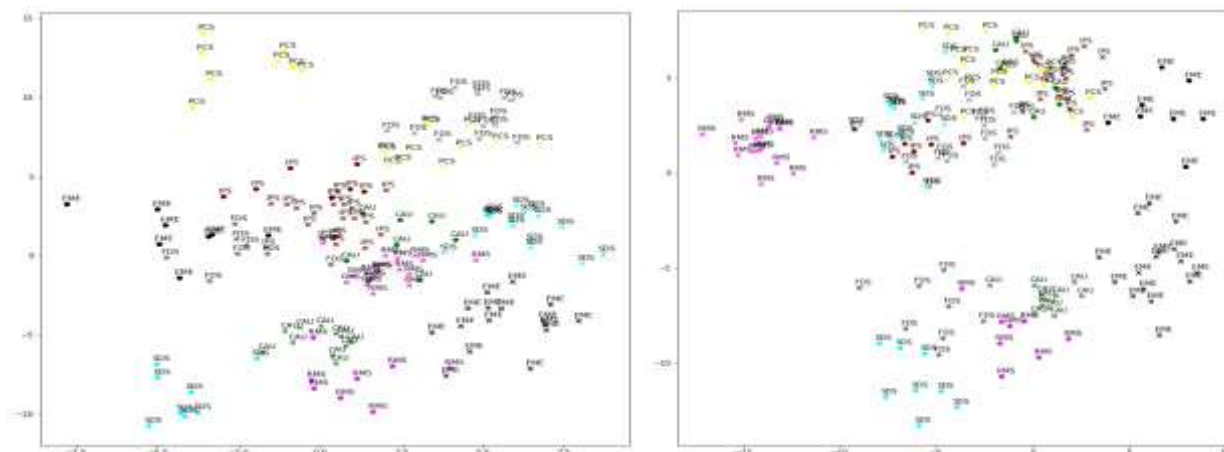


Fig. 1. Embedding pentru sistemul *B* (Baseline) Fig. 2. Embedding pentru sistemul *B + E*

(4) **Adaptarea vorbirii din date audio ne-expresive la date audio cu expresivitate.** Antrenarea sistemelor a folosit implementarea⁷ CSTR bazată pe Tensorflow 1.13. Am efectuat două tipuri de experimente.

Date audio pentru Experiment tip A. Datele folosite pentru antrenarea primului tip de experimente aparțin vorbitorului SAM din corpusul SWARA și au fost împărțite în două categorii, în funcție de evaluarea subiectivă a gradului lor de expresivitate. Cele două categorii selectate sunt bazate pe selectare manuală: date neexpresive: ● rnd1, rnd2, rnd3, diph1, diph2 (2.476 propoziții) ● date expresive: citirea nuvelor “Ivan Turbincă”, „Stan pășitul” (704 de propoziții)

Date audio pentru Experiment tip B. Pentru al doilea tip de experimente am adăugat date neexpresive de la vorbitorii BEA, EME și IPS din corpusul SWARA. De la acești vorbitori am utilizat câte 40 sau 100 de pronunții selectate din rnd1, fiind considerate date neexpresive.

Date text. Textele corespunzător datele audio au fost folosite în trei forme diferite:

(1) Forma ortografică	<i>Pe de altă parte, conform rezultatelor obținute</i>
(2) Forma transcrisă fonetic	<i>p e <> d e <> a l t @ <> p a r t e <,> k o n f o r m <> r e z u l t a t e l o r <> o p t s i n u t e</i>
(3) Forma transcrisă fonetic cu accent	<i>p e 0 <> d e 0 <> a l l t @ 0 <> p a l r t e 0 <,> k o 0 n f o l r m <> r e 0 z u 0 l t a l t e 0 l o 0 r <> o 0 p t s i 0 n u l t e 0</i>

Rezultate Experiment A. După 3.000 de epoci antrenarea a fost continuată cu date expresive de la vorbitorul SAM, textul de intrare fiind același cu cel folosit în modelul antrenat pe date neexpresive. Mostre audio pot fi ascultate aici: <https://gitlab.utcluj.ro/speech/tts-samples/-/wikis/home>. Sistemele antrenate sunt listate în Tabelul 2.

Tabel 2. Descrierea experimentelor din prima categorie (A)

⁷ <https://github.com/CSTR-Edinburgh/ophelia>

Forma text	Date audio	Epoci	Număr de propoziții	Date audio suplimentare	Epoci suplimentare	Număr de propoziții
Ortografică	neexpresiv	3.000	2.476	Expresiv	3.000	704
Fonetică	neexpresiv	3.000	2.476	Expresiv	7.000	704
Fonetic + accent	neexpresiv	3.000	2.476	Expresiv	6.000	704

Rezulate Experiment B. Textul de intrare a fost furnizat sistemului în forma transcrisă fonetic. Pentru a învăța mai mulți vorbitori, ID-ul de vorbitor a fost concatenat textului de intrare și sistemul a fost antrenat cu vorbitori multipli. Identitatea mai multor vorbitori în sistem a fost învățată cu ajutorul metodei LCC (Learning Channel Contribution). Aceste categorii de experimente au avut rezultate modeste, identitatea vorbitorului fiind păstrată sau învățată, dar caracterul de expresivitate al vorbirii a fost pierdut în procesul de învățare. Mostre audio: <https://rb.gy/6pbnnv>. Sistemele antrenate sunt listate în Tabelul 3.

Tabel 3. Descrierea experimentelor din a doua categorie (B)

Forma text	Date audio	Epoci	Număr de propoziții	Vorbitor	Date suplimentare	Epoci adiționale	Nr propoziții	Vorbitor
Fonetică	Expresiv	2.000	704	SAM	neexpresiv	1.500	40	BEA
Fonetică	Expresiv	2.000	704	SAM	neexpresiv	1.500	100	BEA
Fonetică	Expresiv	2.000	704	SAM	neexpresiv	1.500	40	EME
Fonetică	Expresiv	2.000	704	SAM	neexpresiv	1.500	100	EME
Fonetică	Expresiv	2000	704	SAM	neexpresiv	1.500	40	IPS
Fonetică	Expresiv	2000	704	SAM	neexpresiv	1.500	100	IPS

5) O nouă metodă de adaptare la vorbitor folosind post-filtrarea⁸. Extinzând diferitele experimente realizate pentru sistemele de sinteză text - vorbire bazate pe rețele neuronale, am ajuns la concluzia că există o serie de limitări în adaptarea acestora la un nou vorbitor. În special, atunci când se dorește adaptarea rețelei cu un set redus de date, problemele sunt și mai complexe. Astfel, noi am propus o nouă metodă de adaptare a sistemului de sinteză la un nou vorbitor folosind post-filtrarea semnalului vocal sintetizat cu o rețea pre-antrenată cu date puține, tot cu o structură de rețea neuronală care este antrenată să mapeze semnalul vocal sintetizat în semnal vocal natural.

Sistemele de sinteză implementate prin această metodă în laborator au fost validate atât prin metode obiective (MSD), cât și subiective (MOS). Vezi și înregistrarea, <https://youtu.be/OLAJGaqmjqA>.

3.2. Metodă de adaptare rapidă a sistemului de sinteză folosind date atipice (vezi D3.16)

(1) Augmentarea datelor de intrare de tip text prin predicția simultană a informației lexicale de nivel înalt. Pornind de la resursele existente la partenerul ICIA: RoSyllabiDict, MaRePhor, și DEX online, s-au antrenat o serie de rețele neuronale de tip LSTM, BLSTM, CNN, CNN cu atenție, ce au ca obiectiv predicția simultană a transcrierii fonetice, silabificării, poziționării accentului, lematizării, respectiv părții de vorbire (rezultatele au fost prezentate la conferința KES 2020⁹). Experimentul s-a realizat comparativ și pentru limba engleză, iar detaliile aferente sunt prezentate doar în livrabilul D3.16. (vezi Tabel 4).

Tabelul 4. Rezultate ale predicției simultane a informației lexicale din textul de intrare

⁸ <https://youtu.be/OLAJGaqmjqA> (Prezentare la conferință)

⁹ <https://youtu.be/-iLf2ZvTeKY> (Prezentare conferință)

Arhitectură	Limba	Acuratețe [%]			
		3 informații lexicale	fără silabificare	fără accent	la nivel de caracter
<i>CNN cu atenție</i>	<i>română</i>	86.64	88.83	93.84	-
<i>LSTM cu atenție</i>	<i>română</i>	86.26	88.68	92.87	-
<i>LSTM</i>	<i>română</i>	84.60	86.89	91.19	-
<i>BLSTM</i>	<i>română</i>	86.10	88.13	92.73	-
<i>CNN cu atenție</i>	<i>engleză</i>	58.96	59.70	64.00	85.53
<i>LSTM cu atenție</i>	<i>engleză</i>	53.79	54.43	57.24	81.15
<i>LSTM</i>	<i>engleză</i>	52.81	53.41	56.33	90.30
<i>BLSTM</i>	<i>engleză</i>	56.02	56.72	59.71	94.94

Aceste rezultate promițătoare au determinat atât extinderea corpusului de intrare (corpus de 330.000 cuvinte adnotate, RoLEX – Romanian Lexicon), cât și a arhitecturii de rețea neuronală la una de tip Transformer antrenată cu subseturi de 5.000, 50.000, 100.000, 150.000 și 330.000 de cuvinte. Vezi și livrabilul D3.6 (ICIA) pentru o analiză și discuție mai detaliată.

Adițional, pentru a îmbogăți descriptorii pentru textul de intrare s-a dezvoltat și un modul de predicție automată a **lemei unui cuvânt** folosind rețele neuronale recurente LSTM și convoluționale CNN. Sistemele au fost antrenate folosind setul de date DEX online precum și un subset din setul de date CoRoLa, fie doar prin perechi cuvânt – leamnă, fie introducând și atributul de parte de vorbire (POS – Part of Speech Tagging) la nivel de trigram (doar pentru CoRoLa). Pentru setul de date colectat din DEX, sistemele bazate pe rețele LSTM ating o acuratețe de 99.43% la nivel de caracter și cu adnotare POS, iar pentru subsetul CoRoLa cea mai mare acuratețe este de 99.69% pentru CNN.

(2) Aplicarea informației lexicale pentru augmentarea reprezentării textului în sisteme de sinteză text – vorbire cu date atipice. S-au antrenat mai multe sisteme de sinteză text-vorbire, atât pentru limba română (doar pe acestea le prezentăm aici), cât și pentru limba engleză (vezi livrabil D3.16 pentru detalii), toate bazate atât pe arhitectura DC-TTS cât și Tacotron2 și antrenate cu date audio, respectiv text adnotat cu informațiile lexicale prezise conform cu descrierea anterioară. **Evaluare:** în cadrul etapei de evaluare s-a analizat în ce măsură corectitudinea informațiilor lexicale suplimentare influențează sinteza sistemelor în următoarele scenarii de testare: a) silabificare și/sau accent corecte pentru întreaga propoziție, b) silabificare/accent incorecte pentru un singur cuvânt, c) silabificare/accent absente pentru un singur cuvânt, d) silabificare/accent incorecte pentru toate cuvintele (alocate aleator), e) silabificare/accent absente pentru toate cuvintele (alocate aleator)

Sistemele descrise mai sus au fost validate folosind metrica distanța cepstrală MCD (*Mel Cepstral Distortion*). Pentru calcularea valorilor MCD au fost sintetizate câte 20 de propoziții cu maximum 20 de cuvinte și 20 de propoziții scurte cu maxim 5 cuvinte, pentru fiecare sistem, în toate cele 5 scenarii de testare. Rezultatele sunt prezentate în figurile alăturate, iar mostre audio pentru aceste sisteme pot fi accesate și ascultate aici: <https://gitlab.utcluj.ro/speech/tts-samples/-/wikis/home> .

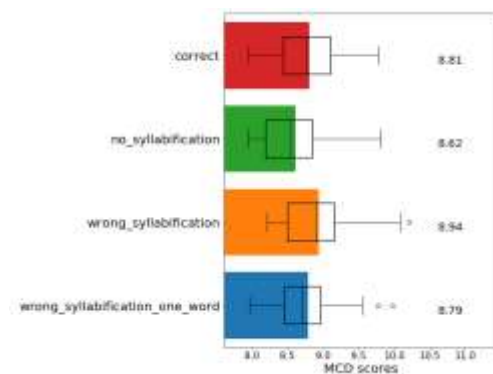


Fig. 3 MCD natural vs sintetic

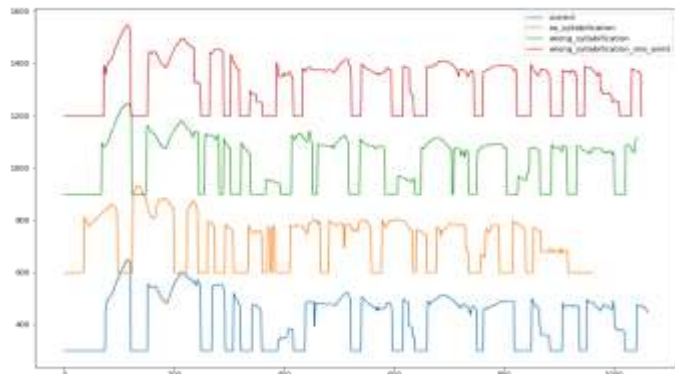


Fig. 4 Variația F0, natural vs sintetic, un vorbitor

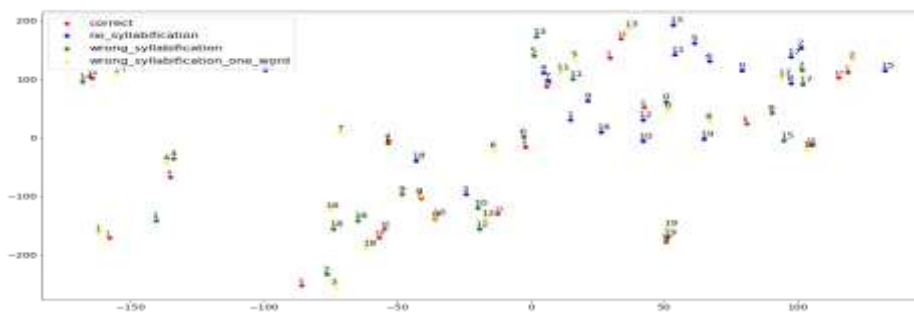


Fig. 5 Proiecțiile t-SNE pentru embedding vorbitor la cele 20 de propoziții

(3) Adaptarea sistemului de sinteză cu date atipice prin transfer de stil de vorbire și Flowtron.

Cercetările din această etapă au avut în vedere și transferul stilului de vorbire folosind un set redus de date prin intermediul metodei *Transfer Learning*. Am adoptat soluția de antrenare *One-Shot Learning*. Această soluție este posibilă datorită fluxurilor de normalizare implementate la nivelul acestei rețele neuronale.

Practic, s-au antrenat 2 sisteme: unul ce folosește doar corpusul Mara (vorbitor unic) și unul ce utilizează corpusul SWARA (vorbitori multipli). Pentru transferul stilului au fost folosite date audio provenite din buletine de știri, precum și date preluate de la alți vorbitori din corpusul SWARA. Deși Flowtron este capabil să realizeze transferul stilului pe baza modificărilor spațiului latent¹⁰, este nevoie ca sistemele antrenate să folosească un set de date extins și de foarte bună calitate. Se poate observa în figura alăturată că alinierea text-audio nu este realizată în mod ideal, ceea ce introduce erori de pronunție și artefacte audio. Rezultate audio ale acestor experimente sunt disponibile la adresa: <https://gitlab.utcluj.ro/speech/tts-samples/-/wikis/home>.

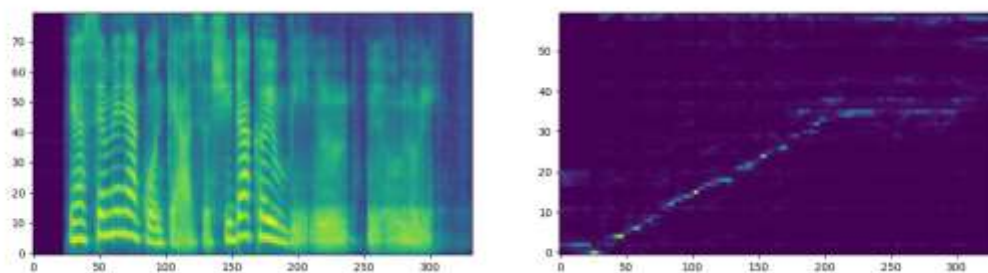


Fig. 6. Exemplu Mel-spectrogramă și aliniere text-audio în sistemul Flowtron-Mara

(4) **Rezultate experimentale ale metodei de adaptare folosind date atipice.** Datele atipice se referă la date imperfecte precum: transcrieri inexacte, lipsa segmentărilor la nivel de propoziție, condiții de zgomot diferite, etc. Toate aceste probleme se transpun în calitate scăzută a ieșirii sistemului de sinteză. Astfel, este importantă dezvoltarea unor metode ce reduc influența acestor date atipice în vocea sintetizată. Pentru acest

¹⁰ <https://nv-adlr.github.io/Flowtron>

scop, s-au antrenat 2 sisteme ce utilizează fie o voce a unei prezentatoare de știri și pentru care erau disponibile doar 20 de minute de audio, fie o voce masculină extrasă dintr-o resursă colectată în P1 și transcrisă în P3. Pentru vocea din urmă, transcrierile nu sunt 100% corecte, astfel că se poate analiza influența erorilor de transcriere. Mostre audio ale acestor voci sunt disponibile în pagina <https://gitlab.utcluj.ro/speech/tts-samples/-/wikis/home>, sistem FEM și sistem COB. Rezultatele indică o foarte bună flexibilitate a arhitecturii sistemului în ceea ce privește adaptarea la date atipice.

3.3. Tehnologie de realizare a interfețelor om-mașină pentru sinteza text-vorbire¹¹ (D3.17)

Au fost extinse tehnologiile de sinteză text-vorbire în limba română bazate pe modele Markov și disponibile în website-ul www.romaniantts.com, la tehnologiile bazate pe rețele neuronale multistrat și utilizarea informațiilor lingvistice suplimentare în transcrierea textului de intrare. Deoarece serverele pe baza de GPU care deserveșc demonstrația online sunt folosite în mod curent de grupul de cercetare pentru diferite experimente, utilizatorii au nevoie de o cheie de API ce poate fi obținută de la coordonatorii P4.

Tabel 5. Sistemele de sinteză disponibile în platforma API:

Tehnologie pentru sistemul de sinteză	Voci ¹² disponibile pentru demonstrare online
Sistem bazat pe rețele neuronale convoluționale (DC-TTS)	Mara, BAS, EME, CAU, DDM, DCS, FDS, IPS, PCS, SAM
Sistem bazat pe rețele neuronale recurente (Tacotron 2) și vocoder WaveGlow	Mara, Mara + informații lexicale, EME, NL, DOL, DLL, ZPL, SRL, FEM
Sistem bazat pe fluxuri de normalizare (Flowtron)	Mara

4. Oferta de servicii de cercetare, locuri de muncă susținute și valorificarea resurselor

Tabel 6. Sintează privind oferta de servicii, locuri de muncă și valorificarea resurselor în UTCN

Oferta de servicii în UTCN	<ul style="list-style-type: none"> oferta unor noi tehnologii de sinteză text-vorbire cu expresivitate, în limba română, bazată pe rețele neuronale și aliniată la standardele internaționale (Tacotron GST, DC-TTS, Flowtron) servicii de adnotare automată a resurselor de date text pornind de la corpusurile CoRoLa și RoLEX servicii de înregistrare audio prin aplicația RecoApy servicii de procesare paralelă a datelor folosind tehnici de învățare automată pe noile echipamente achiziționate în anul 2019 din proiect servicii software pentru dezvoltarea modelelor bazate pe învățare automată. <p>ERRIS: https://erris.gov.ro/speech.utcluj.ro</p>
Locuri de muncă susținute în UTCN	1 x CS I, 1 x CS III, 1 x Tehnician 2 x ACS nou angajați începând cu luna ianuarie 2019
Resursa umană nou angajată în UTCN	Conform acordului de grant au fost angajate 2 ACS, doctoranzi, începând cu 1 ianuarie 2019.
Valorificare resurse în parteneriat	<ul style="list-style-type: none"> UTCN a preluat de la ICIA resurse de date text (CoRoLa, RoLEX) pentru predicția simultană a informațiilor lexicale de nivel înalt UTCN a preluat de la UAIC un set de înregistrări audio / interviuri radio pentru a demonstra antrenarea cu date atipice, respectiv transferul de stil și expresivitate UTCN a furnizat pentru ICIA o arhitectură de tip Transformer pentru a realiza adnotări automate a textului UTCN a preluat de la UPB un set de transcrieri text a unor interviuri furnizate de UAIC UAIC a furnizat pentru UTCN acces la o platformă online pentru stocarea corpusurilor bimodale.
Cecuri	<ul style="list-style-type: none"> Nu au fost planificate cecuri pentru UTCN în această etapă de raportare.

5. Management și comunicare

Activitățile de management au fost orientate în special către managementul proiectului complex în vederea integrării grupurilor de cercetare și a schimbului de resurse de date sau software între partenerii din proiect. S-a solicitat în martie 2020 de către Autoritatea Contractantă întocmirea unui act adițional pentru redistribuirea unei părți din buget pe anul 2021. S-au organizat mai multe conferințe Skype pentru prezentarea rezultatelor intermediare. Din punct de vedere financiar s-au primit 3 tranșe de avans cu o

¹¹ <http://speech.utcluj.ro/ronna/>

¹² Acronimele indică vorbitori anonimizati

regularitate adecvată. Resursele financiare alocate UTCN pentru anul 2019 au fost utilizate în majoritate, cu excepția unor sume în categoria salarii pentru tinerii cercetători (concediu de maternitate pentru o nou angajată), respectiv deplasări (participarea la conferințe s-a realizat online și doar taxele de participare au fost plătite).

6. Diseminarea rezultatelor

O preocupare în UTCN și în această etapă de raportare a fost implementarea și îndeplinirea cu succes a obiectivelor stabilite în strategia de diseminare a rezultatelor elaborată în cadrul propunerii de proiect. Astfel, adecvat acestei etape inițiale s-a acționat pe următoarele direcții:

a) actualizarea paginii web a proiectului SINTERO (<http://speech.utcluj.ro/sintero/>),

b) crearea de pagini Wiki interne grupului de cercetare UTCN (mostre audio cu semnal sintetizat - <https://gitlab.utcluj.ro/speech/tts-samples/-/wikis/home>, aplicația RecoApy - <https://gitlab.utcluj.ro/sadriana/recoapy>), pagini web dedicate pentru demonstrarea online a tehnologiilor dezvoltate în această etapă (interfața online [//speech.utcluj.ro/ronna](http://speech.utcluj.ro/ronna)), prezentări video ale unor articole susținute online la conferințe internaționale (<https://youtu.be/OLAJGaQmjQ4>, http://adrianastan.com/docs/AStan_RecoApy_Interspeech2020.mp4, <https://youtu.be/-iLf2ZvTeKY>,).

c) publicații științifice cu rezultatele cercetărilor la conferințe internaționale în domeniu

d) implicarea studenților în stagii de practică, proiecte de diplomă sau disertații cu tematică apropiată de cea a proiectului: *5 studenți în stagii de practică în perioada iulie – august 2020, 5 studenți la proiect de diplomă susținut în Iulie 2020, 1 student la lucrarea de disertație susținută în Iulie 2020.*

e) acorduri de parteneriat cu mediul economic: 1 acord de parteneriat (în lucru) cu firma SC DiktaCom pentru cesionarea reciprocă de date audio în vederea dezvoltării unor tehnologii pe bază de prelucrare automată a semnalului vocal.

7. Concluzii

Activitățile de cercetare desfășurate în etapa a III-a de implementare a proiectului (2020) au condus la obținerea rezultatelor așteptate și ele sunt în concordanță cu obiectivele specifice ale etapei. Astfel, rezultatele raportate în acest document și descrise detaliat în cele 4 livrabile aferente perioadei de raportare (vezi Secțiunea 8 a acestui raport), asigură testarea și evaluarea finală a tehnologiei dezvoltate în etapa finală a proiectului (2021).

8. Referințe la livrabilele aferente etapei 2020 (Anexe la raport)

[1] Livrabil D3.15:	„Dezvoltarea unei noi tehnologii pentru adaptarea vocii sintetice la stilul și expresivitatea unui nou vorbitor”, Noiembrie 2020.
[2] Livrabil D3.16:	„Dezvoltarea unei noi metode de adaptare rapidă a vocii sintetice folosind date audio atipice”, Noiembrie 2020.
[3] Livrabil D3.17	„Integrare tehnologie nouă și demonstrarea în realizarea interfețelor om-mașină pentru sinteza text – vorbire”, Noiembrie 2020.
[4] Livrabil D3.18:	„Diseminare”, Noiembrie 2020.

Diseminarea rezultatelor proiectului complex ReTeRom

Diseminarea rezultatelor proiectului a fost realizată prin intermediul website-ului proiectului complex http://www.racai.ro/p/reterom/index_en.html care include trimiteri la paginile proiectelor componente: http://www.racai.ro/p/reterom/index_en.html#pr2, <http://cobiliro.info.uaic.ro:3083/>, <https://tadarav.speed.pub.ro>, <https://speech.utcluj.ro/sintero/#>).

De asemenea rezultatele proiectului au fost diseminate prin publicarea mai multor articole științifice:

1. V. Păiș, R. Ion, D.Tușiș. A Processing Platform Relating Data and Tools for Romanian Language. In: Proceedings of the 1st International Workshop on LanguageTechnology Platforms (IWLTP 2020), European Language Resources Association (ELRA), Georg Rehm et al. (eds.), pp. 81-88 - indexed by DBLP and ISI Web of Science.
2. V. Păiș, D. Tușiș, R. Ion. MWSA Task at GlobalLex 2020: RACAI's Word Sense Alignment System using a Similarity Measurement of Dictionary Definitions. In Proceedings of The 12th Language Resources and Evaluation Conference. European Language Resources Association, Marseille, France, pp. 69-75, - indexed by DBLP and ISI Web of Science, May 2020
3. V. Păiș, R. Ion. TermEval 2020: RACAI's automatic term extraction system. In Proceedings of the 6th International Workshop on Computational Terminology. European Language Resources Association, Marseille, France, pp. 101-105, indexed by DBLP and ISI Web of Science, May 2020
4. G. Rehm, ... D. Tușiș,..., F. Yvon. The European Language Technology Landscape in 2020: Language-Centric and Human-Centric AI for Cross-Cultural Communication in Multilingual Europe. In Proceedings of The 12th Language Resources and Evaluation Conference. European Language Resources Association, Marseille, France, pp. 3315-3325, May 2020, indexed by DBLP and ISI Web of Science
5. Ș. Boghiu, D. Gîfu. Event Detection from the Spatial-Temporal Perspective at the 3rd Conference of Doctoral Schools from the UNIVERSITARIA Consortium, see, CSCDCU-MIF2020, Oct. 22-24, 2020, Iași, Romania.
6. Ș., D. Gîfu. A Spatial-Temporal Model for Event Detection in Social Media. In: Procedia Computer Science, Vol. 176, Matteo Cristani, Carlos Toro, Cecilia Zanni-Merk, Robert J. Howlett, Lakhmi C. Jain (eds.), ELSEVIER, 2020, pp. 541-550, ISSN 1877-0509
7. D. Cristea, I. Pistol, Ș.Boghiu, A. Bibiri, D. Gîfu, A. Scutelnicu, M. Onofrei, D. Trandabăț, G. Bugeag. CoBiLiRo: a Research Platform for Bimodal Corpora. In: Proceedings of the 1st International Workshop on LanguageTechnology Platforms (IWLTP 2020), European Language Resources Association (ELRA), George Rehm et al. (eds.), pp. 22-27 - indexed by DBLP and ISI Web of Science.
8. A.Iftene, D. Gîfu, A. R. Miron, M. Șt. Dudu. A Real-Time System for Credibility on Twitter. In: Proceedings of the 12th Conference on Language Resources and Evaluation (LREC 2020), Nicoletta Calzolari et al. (eds.), European Language Resources Association (ELRA), pp. 6168-6175 - indexed by DBLP and ISI Web of Science - Rank C.
9. A.-L. Georgescu, H. Cucu, A. Buzo, C. Burileanu, "RSC: A Romanian Read Speech Corpus for Automatic Speech Recognition," in the Proceedings of The 12th Language Resources and Evaluation Conference (LREC), pp. 6606-6612, 2020, Marseille, France.
10. C. Manolache, A.-L. Georgescu, A. Caranica, H. Cucu, "Automatic Annotation of Speech Corpora using Approximate Transcripts," in the Proceedings of the 43rd International Conference on Telecommunications and Signal Processing (TSP), 2020, Milano, Italy.
11. D. Oneață, A.-L. Georgescu, H. Cucu, D. Burileanu, C. Burileanu, "Revisiting SincNet: An Evaluation of Feature and Network Hyperparameters for Speaker Recognition," in the Proceedings of the 28th European Signal Processing Conference (EUSIPCO), Amsterdam, The Netherlands, 2020.
12. G. Pop, H. Cucu, D. Burileanu, C. Burileanu, "Cough Sound Recognition in Respiratory Disease Epidemics," in Romanian Journal of Information Science and Technology, vol. 23, no. 5, pp. S77-S89, 2020, ISSN 1453-8245, ISI IF 0.661.
13. A.-L. Georgescu, C. Manolache, D. Oneață, H. Cucu, C. Burileanu, "Data-filtering methods for self-training of automatic speech recognition systems," in the Proceedings of the IEEE Spoken Language Technology Workshop (SLT), Virtual, 2021.
14. D. Oneață, A. Caranica, A. Stan, H. Cucu, "An evaluation of word-level confidence estimation for end-to-end automatic speech recognition," in the Proceedings of the IEEE Spoken Language Technology Workshop (SLT), Virtual, 2021.
15. B. Lőrincz, M. Nutu, A. Stan, M. Giurgiu. "An Evaluation of Postfiltering for Deep Learning Based Speech Synthesis with Limited Data", IEEE 10th International Conference on Intelligent Systems (IS), Bulgaria, 2020
16. B. Lőrincz. "Concurrent phonetic transcription, lexical stress assignment and syllabification with deep neural networks", Proceedings of the 24th International Conference on Knowledge-Based and Intelligent Information & Engineering Systems KES2020, 2020.
17. A.Stan, "RECOApy: Data Recording, Pre-Processing and Phonetic Transcription for End-to-End Speech-Based Applications", In Proceedings of the Interspeech, Shanghai, China, 2020
18. K. M Scott, S. Ashby, A. Stan ."Designing a Synthesized Content Feed System for Community Radio", Proceedings of the 11th Nordic Conference on Human-Computer Interaction: Shaping Experiences, Shaping Society, Estonia, 2020

Stadiul elaborării și asumării Programului comun de CDI, cu evidențierea modului de colaborare ulterioară între parteneri și atragerea de noi fonduri naționale/internaționale.

Planul comun este prezentat în anexa acestui raport.

Detalii privind angajarea și menținerea noilor cercetători

Nr. posturi asumate de noi cercetători	7
Nr. posturi ocupate de noi cercetători	8
Nr. posturi ocupate de noi cercetători (în prezent)	7

Lista noi cercetători								
Nr. crt.	Instituție	Nume	Prenume	Poziția ocupată în cadrul proiectului	Data angajare în proiect	Perioada implicare în proiect	Costuri salariale alocate	Costuri salariale plătite
1.	ICIA	Păiș	Vasile	Doctorand/ CS	01.09.2018	1,09.2018 - prezent	211.463	193.033
2.	UAIC	Pădurariu	Constantin Cristian	Membru doctorand	03.12.2018	03.12.2018 – 30.11.2019	110.871	110.436
3.	UAIC	Boghiu	Șerban Paul	Membru doctorand	01.03.2020	01.03.2020 - prezent	102.471	106.761
4.	UPB	Oneață	Dan Theodor	CȘ	24.04.2018	24.04.2018 - prezent	206.066	186.926
5.	UPB	Pop	Gheorghe	AȘC	24.04.2018	24.04.2018 - prezent	74.610	67.586
6.	UPB	Manolache	Cristian	AȘC	24.04.2018	24.04.2018 - prezent	103.871	93.686
7	UTCN	Beáta	Lőrincz	AȘC	03.01.2019	03.01.2019 13.10.2020	116.320	77.696
8	UTCN	Maria	Nuțu	AȘC	03.01.2019	03.01.2019 - prezent	116.320	105.188

ICIA va menține cercetătorul angajat pe proiect cărui îi va oferi un contract de muncă în institut pe perioada nedeterminată, activitatea sa fiind foarte apreciată. El va fi finanțat atât prin fonduri bugetare (1/2 normă) cât și din proiecte extrabugetare (1/2 normă: proiectele europene MARCELL, CURLICAT).

Universitatea "A.I. Cuza" din Iași va menține în cadrul proiectului noul cercetător angajat și îl va plăti din fonduri proprii până în luna februarie 2021, când i se va oferi un post de ASC în cadrul Facultății de Informatică pentru o perioadă de 2 ani.

Universitatea Politehnica din București va menține în vigoare posturile de noi cercetători oferind granturi interne instituției. Modalitatea de finanțare pentru personalul nou angajat în proiecte PCCDI se va face prin următoarele mecanisme:

- Cristian Manolache va fi încadrat cu normă întreagă în proiectul de cercetare "Climatologia aerosolului - de la masuratori de teledetectie la invatare profunda structurata" (proiect CLARA), contract de finanțare nr. 295PED/ 2020; el are contract până pe 16.08.2022.
- Dan Theodor ONEAȚĂ va fi încadrat cu normă întreagă în proiectul de cercetare Recunoașterea automată a vorbirii din semnale multi-modale (VORBIS), contract de finanțare nr. PD 97/ 2020; el are contract până pe 31.08.2022
- Gheorghe POP va aplica la granturile interne UPB de tip Proof of Concept (PoC) destinate tinerilor cercetatori care au participat la proiecte de tip PCCDI, în vederea finanțării activității lor de cercetare în perioada de sustenabilitate de 2 ani din fondurile UPB.

Dupa aceste date, cei doi vor putea accesa și ei un proiect UPB-POC pentru a-și completa cele 3 luni ramase din perioada de sustenabilitate a proiectului ReTeRom.

Universitatea Tehnică Cluj Napoca va menține posturile de noi cercetători și va oferi granturi interne instituției pe bază de competiție. Modalitatea de finanțare pentru personalul nou angajat în proiecte PCCDI se va face prin următoarele mecanisme: a) încadrare cu norma întreaga pe un singur proiect sau cumulativ până la o normă întreaga din mai multe proiecte de cercetare din universitate în care sunt necesare competențele noilor cercetători; b) încadrare cu norma întreaga pe posturi sustenabile financiar din Fondul de Sustinere a Cercetării cu un salariu care să maximizeze perioada de sustenabilitate (Cf. HCA 126/17.11.2020).

ANEXĂ

Stadiul elaborării și asumării Programului comun de CDI, cu evidențierea modului de colaborare ulterioară între parteneri și atragerea de noi fonduri naționale/internaționale.

- a) Considerăm că este oportun să formăm (consorțiul ReTeRom) un pol de cercetare la nivel național în domeniul tehnologiei vorbirii. Aceasta implică agregare de resurse începând de la HW, SW și resurse umane (le avem) + un plan de cercetare (de construit împreună). Considerăm că doar așa putem intra în marile proiecte ale CE, dar și la nivel național.
- b) Derivat din (a), acțiuni punctuale pentru: formare/educație, colaborare cu mediul economic, atragere de fonduri pentru sustenabilitate.
Acest proiect nu este neapărat pentru noi, cei din lista de proiect, cât mai degrabă pentru cei care vin după noi și în care am investit timp, credință, onestitate.

ICIA are în vedere aprofundarea cercetărilor și continuarea dezvoltării de aplicații de IA

A. Ne vom concentra o parte din eforturi pentru realizarea unor sisteme ASR care să poată fi folosite pentru transcrierea unor înregistrări de lungă durată. Sistemele realizate în cadrul acestui proiect nu permit transcrierea unei înregistrări mai lungi de aproximativ 20 secunde. Această restricție apare din cauza volumului de memorie RAM necesară prelucrării semnalului vocal. Pentru depășirea acestei limitări avem în vedere studiul și dezvoltarea unui flux de prelucrare inteligentă presupunând segmentarea înregistrării audio și prelucrarea succesivă a segmentelor audio cu concatenarea transcrierilor acestor segmente. Pentru că niciun sistem ASR nu este perfect, fluxul de prelucrare menționat trebuie încorporat unui mediu inteligent și prietenos de corectare a erorilor de transcriere. Problema nu este trivială și o astfel de aplicație este extrem de necesară și solicitată intens de parteneri academici dar și din media. Nu avem cunoștință de un astfel de sistem operațional pentru limba română.

B. Valorificând rezultatele obținute de ICIA într-un proiect anterior respectiv dezvoltarea unui sistem de traducere automată EN↔RO a limbajului scris și rezultatele comune cu UPB și UTCN obținute în cadrul acestui proiect în domeniul prelucrării limbajului vorbit, avem în vedere un proiect de traducere a limbajului vorbit EN↔RO. Pentru limba română vom folosi sistemele noastre ASR (ICIA și UPB) și TTS (ICIA și UTCN) iar pentru limba engleză vom folosi sisteme open-source existente și vom dezvolta noi modulele corespunzătoare. Această soluție în cascadă, deși au apărut sisteme neuronale ce realizează traducerea automată direct din semnalul vocal în semnal vocal-fără a mai trece prin fazele de ASR și TTS, este relativ ușor de realizat. Eventual, într-o fază ulterioară am putea să ne propunem realizarea unui sistem de traducere automată EN↔RO direct din semnalul vocal fără a mai trece prin etapa textuală. Rezultatele direcției de cercetare/dezvoltare de la punctul A, vor fi direct relevante pentru crearea de resurse de antrenare pentru sisteme dezvoltate în direcția B de cercetare/dezvoltare.

C. Valorificând rezultatele cercetărilor curente ale ICIA în ce privește prelucrarea textelor românești cu rețele neuronale complexe și cercetări recente în reprezentarea cunoștințelor în grafuri de relații semantice, ne interesează construirea automată a unui graf de cunoștințe cuprinzător din situl enciclopedic Wikipedia, utilizând rețele neuronale de tip BERT pentru a descoperi relațiile semantice care se stabilesc între entități. Cu un astfel de graf disponibil, care este de așteptat să conțină sute de mii de relații semantice, vom construi un sistem de întrebare-răspuns (de ex. „Câte planete sunt în sistemul solar?”) cu care se va putea interoga Wikipedia în limba română, utilizând sistemele de ASR și TTS descrise la punctul B.

UAIC are în vedere un număr de posibile aplicații care să demonstreze potențialul tehnologiilor dezvoltate în celelalte proiecte componente, în scopul măririi vizibilității proiectului, în special după finalizarea acestuia. Tehnologiile vizate în special sunt cele descrise în rapoartele 3.7 (Definitivarea, testarea, validarea și împachetarea într-o soluție „ready-to-use” a platformei integrate și configurabile de prelucrare a textelor în limba română.), 3.10 (Îmbunătățirea soluției de filtrare și aliniere a transcrierilor aproximative cu semnalul de vorbire), 3.15 (Dezvoltarea unei noi tehnologii pentru adaptarea vocii sintetice la stilul și expresivitatea unui nou vorbitor) și 3.17 (Integrare tehnologie nouă și demonstrare în realizarea interfețelor om-mașină pentru sinteza text-vorbire). O parte din proiectele descrise au depășit deja faza de proiect, fiind în stadii relativ avansate de implementare (în activități cu studenții de la UAIC-FII). Cele 6 proiecte de aplicații descrise sunt prezentate succint mai jos:

Aplicația 1: Suport pentru învățarea limbii române - PD builder. Un dicționar de pronunție pentru cuvintele unei limbi poate constitui un suport semnificativ pentru cei care doresc să învețe acea limbă.

Colecția de resurse colectate pe platforma CoBiLiRO are potențialul să atingă o dimensiune suficientă pentru a include pronunția mării majorității a cuvintelor din limba română; de aici a plecat ideea de a dezvolta o tehnologie ce caută în resursele din Portal și reproduce pronunții ale unor cuvinte, la cererea utilizatorului. Sunt utilizate adnotările produse de TEPROLIN și alinierea produsă de TADARAV, marcate ce însoțesc majoritatea resurselor existente pe platforma CoBiLiRo.

Aplicația 2: Analiza corpusurilor bimodale. Proiectul propune dezvoltarea unui sistem capabil să evalueze calitatea unui corpus bimodal voce-text din perspectiva unei alinieri automate sau manuale. O parte din posibilele erori semnalate ar putea fi corectate, lucru care ar mări semnificativ calitatea alinierii automate dar și a resursei bimodale în general. Prin corelarea statistică a mai multor parametri, pot fi detectate posibile diferențe între conținutul înregistrării și transcrierea ei în text, diferențe care pot avea cauze diverse: segmente adționale prezente în text față de înregistrare sau invers, transcrieri aproximative, particularități ale pronunției și prozodiei vorbitorului etc. Un astfel de sistem ar fi util atât dezvoltatorilor și utilizatorilor tehnologiilor de aliniere automată cât și celor care colectează și utilizează resurse bimodale.

Aplicația 3: *I listen to my speaking agent reading book fragments as I walk by.* Aplicația are la bază o colecție de texte care abundă în entități geografice, marcate XML explicit, textele fiind însoțite de metadate care descriu minimum: autorul și titlul cărții, anul de apariție și editura. Instalată pe un dispozitiv mobil, ea va semnala proximitatea telefonului față de locațiile menționate în texte și va citi acele fragmente care includ mențiunile respective. În felul acesta, o plimbare printr-un mare oraș se poate transforma într-o călătorie literară. Aplicația se adresează persoanelor care ar dori să primească sugestii de lecturi și doresc să afle lucruri noi despre locurile pe care le vizitează. Ea își propune să îmbine într-un mod plăcut literatura și tehnologia.

Aplicația 4: Asistent inteligent al ședințelor online. Pandemia Covid a obligat ca o parte neneglijabilă a activităților colectivității umane să migreze în mediul online. Propunem realizarea unui instrument sau a unei colecții de API-uri care să se poată integra aplicațiilor de teleconferințe, în scopul extragerii de informații și generării de rapoarte din discuțiile purtate și din mesajele schimbate. Astfel, se poate imagina: transcrierea conversațiilor în text, realizarea de rezumate ale întâlnirilor, extragerea de informații punctuale, pe anumite segmente din conferință ori din intervențiile unor anumiți vorbitori, generarea automată de procese verbale, interogarea minutei ori a procesului verbal generat, pe bază de cuvinte cheie, căutări de secvențe sonore în înregistrarea sonoră etc.

Aplicația 5: *Tracking assistant.* O aplicație Android care își propune să asiste persoanele care suferă de boala Alzheimer în a-și aminti traseele pe care le-au efectuat în timpul zilei printr-o interfață ușor de folosit, în limbaj natural. Pe baza datelor receptate din mai multe canale ale unui telefon mobil, se poate reface traseul și tipul de activitate al persoanei. La sfârșitul zilei pacientul ar purta un dialog cu sistemul inteligent, care l-ar ajuta să-și rememoreze activitățile și traseele de peste zi.

UPB are în vedere dezvoltarea continuă a sistemelor de transcriere automată a vorbirii și integrarea lor cu alte sisteme inteligente de procesare de text, respectiv sinteză de vorbire. Suplimentar față de participarea în proiectele expuse de ceilalți parteneri, UPB intenționează să demareze următoarele proiecte, implicând bineînțeles și partenerii din proiectul ReTeRom, în funcție de ariile lor de expertiză:

Chat-bot vocal. În ultima perioadă aplicațiile de chat-bot text au luat o amploare remarcabilă. Majoritatea site-urilor web ale marilor prestatori de servicii de telefonie, banking, etc. din România pun la dispoziție pe pagina principală un sistem interactiv de transmitere/recepție de mesaje text cu automatizarea răspunsurilor pentru întrebările frecvent apărute (chat-bot text). Tehnologiile dezvoltate sau îmbunătățite în proiectul ReTeRom pot face obiectul extinderii acestor aplicații de tip chat-bot prin adăugarea de funcționalități pentru interacțiunea prin voce. Astfel, **sistemul de recunoaștere automată a vorbirii** dezvoltat de **UPB** poate fi folosit pentru a transcrie mesajele vocale ale utilizatorilor, **tehnologiile de procesare de limbaj natural** dezvoltate de **ICIA** și **UAIC** pot fi folosite pentru înțelegerea mesajelor și generarea de răspunsuri text adecvate, iar **sistemul de sinteză de vorbire** dezvoltat de **UTCN** poate fi utilizat pentru sinteza și redarea răspunsurilor vocale către utilizator.

Sistem de IVR inteligent. Sistemele de IVR (interactive voice response) implementate în România folosesc în acest moment meniuri prin care se poate naviga prin voce, utilizatorul fiind foarte restricționat din punctul de vedere a ceea ce poate să rostească (e.g. „*Vă rugăm să ne spuneți care este problema dumneavoastră rostind: abonamentul de voce, abonamentul de internet, cartela telefonului, telefonie fixă*”). Un sistem IVR inteligent ar trebui să se adreseze utilizatorului în limbaj natural, similar cu un operator (e.g. „*Cu ce vă putem ajuta*”) și să poată să răspundă oricăror solicitări vocale ale utilizatorului. Pentru a implementa un astfel de sistem IVR inteligent pot fi folosite sistemele implementate în cadrul proiectului

ReTeRom. Astfel, **sistemul de recunoaștere automată a vorbirii** dezvoltat de **UPB** poate fi folosit pentru a transcrie solicitările vocale ale utilizatorilor, **tehnologiile de procesare de limbaj natural** dezvoltate de **ICIA** și **UAIC** pot fi folosite pentru înțelegerea solicitărilor și generarea de răspunsuri text adecvate, iar **sistemul de sinteză de vorbire** dezvoltat de **UTCN** poate fi utilizat pentru sinteza și redarea răspunsurilor vocale către utilizator.

UTCN are în vedere aprofundarea cercetărilor realizate în cadrul ReTeRom în domeniul Sintezei din Text a Vorbirii (STV) în limba română și consolidarea colaborării fructuoase pe care am avut-o cu partenerii în acest proiect, în vederea realizării unui pol de cercetare foarte puternic în domeniul tehnologiei vorbirii. Ne propunem o agendă de cercetare comună, care se dezvoltă pe următoarele direcții:

A. Colaborare cu partenerii din proiect și cu mediul economic pentru dezvoltarea resurselor audio și de text.

(A1) *Extinderea de către UTCN a corpusului SWARA* de la 65 de ore de vorbire, la mai mult de 100 de ore de vorbire, în următorul an, respectiv îmbogățirea lui cu date audio expresive și cu diversificarea modelului de limbaj pentru a putea deservi și crearea de sisteme de Recunoaștere Automată a Vorbirii (RAV).

(A2) *Extinderea adnotărilor automate* pentru resursele de semnal vocal și de text colectate în ReTeRom prin aplicarea tehnologiilor de prelucrare a textului dezvoltate de ICIA și a datelor audio dezvoltate de către UPB și UTCN.

(A3) *Cooperarea cu mediul economic* pentru extinderea resurselor de date audio și text, atragerea de fonduri pentru a asigura sustenabilitatea cercetărilor, respectiv pentru transfer tehnologic. UTCN și UAIC au dezvoltat un protocol de colaborare cu o firmă privată în vederea cesionării reciproce de resurse necesare sistemelor de antrenare (RAV, STV)

B. Dezvoltarea în parteneriat a tehnologiilor rezultate din proiect și aducerea lor pe un nivel apropiat de funcționarea în condiții reale.

(B1) *Dezvoltarea unui sistem integrat de interacțiune om – mașină care să includă RAV și STV.* Tehnologiile de RAV disponibile la UPB și ICIA, respectiv de STV, disponibile la UTCN, creează premisele unui proiect comun privind comunicarea om – mașină în limbaj natural. O provocare în cercetare, rămâne interpretarea semantică a mesajului vocal.

(B2) *Dezvoltarea tehnologiei de STV* prin creșterea naturaleței semnalului sintetizat cu apropierea de vorbirea spontană, respectiv adaptare la noi vorbitori prin aplicarea noilor tehnologii de conversie de voce.

(B3) *Optimizarea sistemului de STV*, respectiv creșterea gradului de paralelizare a procesărilor, pentru a răspunde unor aplicații de timp real.

C. Atragerea de noi fonduri pentru cercetare și susținerea noilor cercetători.

(C1) *UTCN a aplicat pentru mai multe proiecte de cercetare în anul 2020, iar în prezent este în negociere pentru a participa într-un proiect de tip Horizon 2020 în apelul Green Deal.*

(C2) *UTCN a aplicat ca și partener într-o propunere de tip COST (Cooperation in Science and Technology) cu titlul „Speech Technology in the Age of Security and Privacy Legislation” și coordonată de EURECOM (Franța).*

(C3) *UTCN are în strategia instituțională obiective ample privind atragerea de fonduri pentru cercetare, respectiv susținerea noilor cercetători prin realizarea unor competiții interne de proiecte de cercetare.*

Director Proiect Complex ReTeRom

Acad. Dan Tufis

