

# Proiect component TADARAV

---

## - Raport științific și tehnic în extenso 2020 -

*Lucian Georgescu, Alexandru Caranica, Cristian Manolache, Gheorghe Pop, Dan Oneață,  
Horia Cucu, Dragoș Burileanu, Corneliu Burileanu*

**Program:** PNCDI III - Programul 1 - Dezvoltarea sistemului național de CD

**Proiect complex:** *Resurse și tehnologii pentru dezvoltarea interfețelor om-mașină în limba română (ReTeRom)*

**Proiect component:** *Tehnologii pentru adnotarea automată a datelor audio și pentru realizarea interfețelor de recunoaștere automată a vorbirii (TADARAV)*

**Data:** 30.11.2020

**Etapa:** 3 / 2020

**Activitatea / activitățile:**

- *Activitatea 3.9 - Analiza impactului utilizării de RAV complementare pentru generarea de adnotări în contextul îmbunătățirii sistemelor de RAV*
- *Activitatea 3.10 - Îmbunătățirea soluției de filtrare și aliniere a transcrierilor aproximative cu semnalul de vorbire*
- *Activitatea 3.11 - Îmbunătățirea soluției pentru generarea de scoruri de încredere pentru RAV*
- *Activitatea 3.12 - Analiza impactului utilizării transcrierilor aproximative în vederea reantrenării sistemelor de RAV*
- *Activitatea 3.13 - Analiza impactului utilizării scorurilor de încredere pentru filtrarea transcrierilor RAV în vederea reantrenării sistemelor RAV*
- *Activitatea 3.14 - Diseminare*

**Număr contract:** 73PCCDI / 2018

**Acord de colaborare:** 30/20.02.2018 ICIA, 4726/01.03.2018 UTCN, 3950/07.03.2018 UPB, 3805/06.03.2018 UAIC

**Autoritatea contractantă:** *Unitatea Executivă pentru Finanțarea Învățământului Superior, a Cercetării, Dezvoltării și Inovării*

**Conducător proiect component:** *Universitatea POLITEHNICA din București*

**Conducător proiect complex:** *ICIA*

**Responsabil proiect component:** *Conf. Horia Cucu*

**Responsabil proiect complex:** *Prof. Corneliu Burileanu*

# Cuprins

Rezumatul etapei	4
Descrierea științifică și tehnică a activităților	6
2.1 Seturi de date	6
2.1.1 Seturi de date de vorbire adnotată	6
2.1.2 Seturi de date de vorbire brute	8
2.1.3 Seturi de date de vorbire adnotată rezultate în această etapă a proiectului	9
2.2 Activitatea 3.9 - Analiza impactului utilizării de RAV complementare pentru generarea de adnotări în contextul îmbunătățirii sistemelor de RAV	9
2.2.1 Dezvoltarea sistemului de transcriere de vorbire ESPnet pentru limba română	9
Arhitectura sistemului RAV propus	9
Arhitectura software a sistemului ESPnet	11
Setup-ul experimental pentru experimentele de RAV	12
Calibrarea hiperparametrilor sistemului de RAV	12
Rezultatele experimentale	14
2.2.2 Dezvoltarea metodei de adnotare bazate pe sisteme RAV complementare și rezultatele obținute pe parcursul proiectului	14
2.3 Activitățile 3.10 și 3.12 - Îmbunătățirea soluției de filtrare și aliniere a transcrierilor aproximative cu semnalul de vorbire și Analiza impactului utilizării transcrierilor aproximative în vederea reantrenării sistemelor de RAV	16
2.3.1 Îmbunătățirea soluției de aliniere a transcrierilor aproximative cu semnalul de vorbire	16
Metoda de bază	16
Metoda propusă	17
Evaluarea metodelor	18
Rezultate alinieri	18
Sisteme RAV ce utilizează noile corpusuri obținute	19
2.3.2 Aplicarea metodei pe setul de date CoBiLiRo-raw	19
Analiză inițială a setului de date CoBiLiRo-raw	19
Aplicare metodă de bază set de date CoBiLiRo-raw	22
2.3.3 Aplicarea metodei pentru setul de date CDep-raw	22
2.4 Activitățile 3.11 și 3.13 - Îmbunătățirea soluției pentru generarea de scoruri de încredere pentru RAV și Analiza impactului utilizării scorurilor de încredere pentru filtrarea transcrierilor RAV în vederea reantrenării sistemelor RAV	23
2.4.1 Legături cu starea artei	24
2.4.2 Metodologia	24
Estimarea scorurilor de încredere	25
Îmbunătățirea probabilităților la nivel de token	25

<b>2.4.3 Setup-ul experimental pentru experimente pe limba engleză</b>	26
<b>Baze de date</b>	26
<b>Sistemul de recunoaștere automată a vorbirii</b>	27
<b>Metrici de evaluare</b>	27
<b>2.4.4 Rezultate experimentale pentru limba engleză</b>	27
<b>Caracteristici și metode de agregare</b>	27
<b>Scalarea temperaturii și tehnica de dropout</b>	29
<b>Ansambluri de modele</b>	29
<b>2.4.5 Setup-ul experimental pentru experimente pe limba română</b>	31
<b>2.4.6 Rezultate experimentale pentru limba română</b>	32
<b>2.4.7 Utilizarea scorurilor de încredere propuse pentru generarea de date adnotate</b>	33
<b>2.5 Activitatea 3.14 - Diseminare</b>	34
<b>2.6 Crearea unui sistem de RAV îmbunătățit</b>	34
<b>2.6.1 Actualizarea modelelor de limbă pentru transcriere de vorbire</b>	34
<b>2.6.2 Utilizarea seturilor de date rezultate din proiect pentru antrenarea RAV</b>	37
<b>2.6.3 Sistemul de RAV Speed îmbunătățit</b>	40
<b>2.7 Bibliografie</b>	40
<b>Structura ofertei de servicii de cercetare și tehnologice</b>	43
<b>Locuri de muncă susținute prin program</b>	44
<b>Valorificarea și îmbunătățirea competențelor și resurselor existente la nivelul consorțiului</b>	44

## 1 Rezumatul etapei

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

1. 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.
2. 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;
3. 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.

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 12% pe vorbire citită, respectiv 16% pe vorbire spontană.
- Activitățile 3.10 și 3.12 - Îmbunătățirea soluției de filtrare și aliniere a transcrierilor aproximative cu semnalul de vorbire și 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 și 3.13 - Îmbunătățirea soluției pentru generarea de scoruri de încredere pentru RAV și 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.

În urma activităților A3.9 - A3.13 din etapa 3/2020 a proiectului TADARAV, au rezultat toate livrabilele asumate de consorțiu la începutul acestei etape:

- Raport de analiză a impactului utilizării de RAV complementare pentru generarea de adnotări în contextul îmbunătățirii sistemelor de RAV;

- 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);
- Raport de analiză a impactului utilizării transcrierilor aproximative în vederea reantrenării sistemelor de RAV;
- Raport de analiză a impactului utilizării scorurilor de încredere pentru filtrarea transcrierilor RAV în vederea reantrenării sistemelor RAV;
- Sistem de RAV actualizat (TRL5).

Diseminarea rezultatelor proiectului (activitatea A3.14) a fost realizată prin intermediul website-ului proiectului (<https://tadarav.speed.pub.ro>) și prin publicarea mai multor articole științifice după cum urmează:

- 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.
- 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.
- 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.
- 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. S, pp. S77–S89, 2020, ISSN 1453-8245, ISI IF 0.661.
- 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.
- 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.

Dintre articolele listate mai sus, al patrulea este deja indexat în Web of Science (Thompson Reuters - ISI), al doilea este deja indexat IEEE Xplore și în curs de indexare în Web of Science, primul și al treilea au apărut în volumul conferințelor respective și sunt în curs de indexare în Web of Science, iar al cincilea și al șaselea vor apărea în volumul conferinței respective și vor fi indexare în Web of Science.

## 2 Descrierea științifică și tehnică a activităților

### 2.1 Seturi de date

#### 2.1.1 Seturi de date de vorbire adnotată

Pentru antrenarea și evaluarea sistemelor de RAV, au fost folosite două seturi de date de vorbire în limba română: *Read Speech Corpus* (RSC), ce conține vorbire citită, colectată în condiții de laborator, fără zgomot de fundal și *Spontaneous Speech Corpus* (SSC), ce conține vorbire continuă, spontană, preluată de la posturi de radio și TV, uneori afectată de zgomot. Ambele corpusuri cuprind fișiere audio și transcrieri corespunzătoare și sunt divizate în seturi de antrenare și seturi de evaluare. RSC-train este setul de antrenare din RSC, ce conține 100 ore de vorbire citită, cuvinte izolate sau fraze de la 157 de vorbitori diferiți. RSC-eval este setul de evaluare din RSC; acesta conține vorbire de la 22 de vorbitori diferiți, însumând 5.5 ore de vorbire. SSC-train1+2 este setul de antrenare din SSC și conține 130 ore de vorbire spontană, majoritatea din emisiuni de știri și *talkshow*-uri. SSC-eval1 este setul de evaluare din SSC și însumează 3.5 ore de vorbire. Informații detaliate despre seturile de date RSC [Georgescu, 2020] și SSC se regăsesc în raportul TADARAV 2019 [Georgescu, 2019a], secțiunea 2.1.1.

Seturile de date de evaluare RSC-eval și SSC-eval au făcut în acest an obiectul unei analize amănunțite în vederea corectării eventualelor erori de adnotare. Această analiză a relevat faptul că o serie întregă de pronunții din setul RSC-eval nu corespund transcrierii de referință. Practic, vorbitorii care au înregistrat acele propoziții au rostit cuvinte în plus sau în minus, au pronunțat greșit anumite cuvinte etc. Câteva dintre erorile identificate și corectate sunt listate în Tabelul 2.1.a. În setul de date SSC-eval, set ce a fost obținut prin transcrierea unor materiale audio deja existente, erorile identificate și corectate constau în mare măsură în omisiuni de semne diacritice, inversiuni de litere sau cuvinte scrise greșit în transcrierea de referință. Tabelul 2.1.a prezintă și o serie de erori identificate și corectate în setul de date SSC-eval.

**Tabelul 2.1.a** Exemple de erori din seturile de evaluare RSC-eval și SSC-eval identificate și corectate

Set/ ID audio	Eroare și corectură
RSC-eval 018_12_0066	ERR: [...] meciul universitatea CLUJ cfr cluj nu trebuia oprit COR: [...] meciul universitatea CRAIOVA cfr cluj nu trebuia oprit
RSC-eval 018_12_0069	ERR: [...] cu scorul de unu LA zero [...] COR: [...] cu scorul de unu ** zero [...]
RSC-eval 003_01_0573	ERR: [...] mea culpa am fost UN dobitoc [...] COR: [...] mea culpa am fost ** dobitoc [...]
RSC-eval 003_01_0575	ERR: [...] nu mai e pe MOMENT disponibil [...] COR: [...] nu mai e pe MOMENTUL disponibil [...]
SSC-eval	ERR: jandarm aflat în POSTU de pază COR: jandarm aflat în POSTUL de pază
SSC-eval	ERR: [...] era ANAGAJAT la jandarmeria sibiu din două mii cinci [...] COR: [...] era ANGAJAT la jandarmeria sibiu din două mii cinci [...]
SSC-eval	ERR: în cele din URMA șerban ionescu a infirmat această știre COR: în cele din URMă șerban ionescu a infirmat această știre
SSC-eval	ERR: [...] răspunde tuturor ÎNTREBARILOR puse în ultimele zile COR: [...] răspunde tuturor ÎNTREBĂRILOR puse în ultimele zile
SSC-eval	ERR: pentru a o ajuta pe soția mea să ÎNCETZE teatrul zilnic fățarnicia COR: pentru a o ajuta pe soția mea să ÎNCETEZE teatrul zilnic fățarnicia
SSC-eval	ERR: POMPIERII au încercat să ajungă la el cu o barcă pneumatică COR: POMPIERI au încercat să ajungă la el cu o barcă pneumatică

După corectare, seturile de date RSC-eval și SSC-eval (versiunile v2) au fost utilizate pentru reevaluarea celui mai performant sistem de RAV pe limba română, disponibil la finalul anului 2019. Redăm în Tabelul 2.1.b rezultatele, măsurate în termeni de eroare la nivel de cuvânt (WER), acestui sistem de RAV pe versiunile v1, respectiv v2 ale seturilor de evaluare. În acest fel, dorim să punem în evidență diferența artificială de performanță rezultată nu din îmbunătățirea sistemului de RAV, ci din corectarea seturilor de date de evaluare. De această diferență artificială de performanță ar trebui să ținem cont ori de câte ori vom compara rezultatele de RAV raportate de Speed pe seturile de date RSC-eval și SSC-eval înainte vs. după această corectură (înainte de anul 2020, respectiv după și inclusiv în anul 2020).

**Tabelul 2.1.b** Diferența artificială de performanță rezultată în urma corectării seturilor de date RSC-eval și SSC-eval pentru cel mai performant sistem de RAV Speed la finalul anului 2019.

WER [%]	RSC-eval		SSC-eval	
	v1	v2	v1	v2
Best RAV Speed 2019	2.41	1.83	12.45	11.04

În etapa anterioară a proiectului, ca parte a activității A2.11, au fost obținute seturile de date de vorbire adnotată SSC-train3-trans-v4 și SSC-train4-trans-v4 (vezi raportul TADARAV 2019 [Georgescu, 2019a], secțiunea 2.2). Seturile de date au fost adnotate folosind metoda corelației cu transcrierile aproximative. Și celelalte metode de adnotare automată au fost utilizate pentru generare de adnotări, pornind de la seturile de date neadnotate SSC-train3-raw și SSC-train4-raw, însă s-a dovedit că aceste seturi de date de vorbire adnotată (SSC-train3-trans-v4 și SSC-train4-trans-v4) sunt cele mai utile în reantrenarea sistemelor de RAV. În consecință, aceste seturi de date au fost înglobate în setul de date *Spontaneous Speech Corpus* (SSC) și au fost folosite în continuare în etapa 2020 pentru antrenarea sistemelor de RAV. Diversele componente ale seturilor de date SSC și RSC sunt prezentate în Tabelul 2.1.c.

În etapa curentă a proiectului ReTeRom, suplimentar față de activitățile din cadrul proiectului, grupul Speed a mărit numărul de seturi de date de evaluare după cum urmează. Din setul de date de vorbire brut SSC-train3+4-raw au fost selectate și transcrise 100 de fișiere de aproximativ un minut fiecare. A rezultat astfel setul de date SSC-eval2 ce cuprinde aproximativ 91 de minute de vorbire provenind de la un post de radio și două posturi de televiziune. Mai multe informații despre achiziția setului de date brut SSC-train3+4-raw se găsesc în raportul TADARAV 2019 [Georgescu, 2019a], secțiunea 2.1.1. Din setul de date CDep-raw, ce conține discursuri susținute în Parlamentul României, au fost selectate și transcrise 300 de fișiere de aproximativ un minut fiecare. A rezultat astfel setul de date CDep-eval ce cuprinde aproximativ 5 ore de vorbire. Setul de date de vorbire CDep-raw este descris succint în secțiunea următoare.

Suplimentar, în etapa curentă a proiectului a fost utilizat pentru antrenarea sistemelor de RAV și setul de date de vorbire din corpusul CoRoLa [Barbu, 2018]. Textele citite din CoRoLa sunt în principal înregistrări profesionale din diverse surse (posturi de radio, studiouri de înregistrare) împreună cu transcrierile lor. O altă parte a corpusului oral este reprezentată de texte citite de la posturi de știri (radio) sau texte citite de vorbitori profesioniști, înregistrate în studiouri, respectiv texte extrase din Wikipedia și citite de voluntari neprofesioniști, înregistrate în medii neprofesionale. O parte semnificativă din transcrierile pentru acest set de date de vorbire conțineau adnotări inconsistente de nume de vorbitori, pronunții incomplete sau informale etc. Pentru a putea utiliza o parte cât mai mare din acest set de date de vorbire la antrenarea sistemelor de RAV, au fost realizate corecturi ortografice în mod semiautomat.

### 2.1.2 Seturi de date de vorbire brute

În etapa curentă metodele de adnotare automată au fost aplicate pe două seturi de date brute: CDep-raw și CoBiLiRo-raw (Tabelul 2.1.d)

Setul de date brut Camera Deputaților (CDep-raw), a fost achiziționat în cursul anului trecut și cuprinde înregistrările video și stenogramele ședințelor de pe site-ul Camerei Deputaților extrase între ianuarie 2003 și februarie 2019. Setul de date conține 3.510 ore de vorbire de la ~2.500 de vorbitori, însoțite de transcrieri aproximative însumând ~25M de cuvinte. Setul de date CDep-raw a fost achiziționat cu ajutorul mai unor aplicații Java ce efectuează mai multe procese: (i) extragerea fișierelor video și HTML pentru fiecare ședință;

**Tabelul 2.1.c** Seturile de vorbire adnotată folosite pentru antrenarea și evaluarea sistemelor de RAV și seturile de vorbire adnotată obținute în etapa anterioară (2/2019)

Setul de date	Subset	Durăță
<b>Spontaneous Speech Corpus (SSC)</b>	SSC-train1+2	130h, 44m
	SSC-train3-trans-v4	41h, 00m
	SSC-train4-trans-v4	250h, 10m
	SSC-eval1	3h, 29m
	SSC-eval2	1h, 31m
<b>Read Speech Corpus (RSC)</b>	RSC-train	94h, 46m
	RSC-eval	5h, 29m
<b>Contemporary Romanian Language (CoRoLa)</b>	n/a	85h, 11m
<b>Camera Deputaților (CDep)</b>	CDep-eval	5h, 00m

(ii) conversia fișierelor video în fișiere wav; (iiia) extragerea transcrierilor fiecărui vorbitor în fișiere HTML separate; (iiib) tăierea fișierelor wav în segmente de fișiere audio corespunzătoare transcrierilor vorbitorilor; (iiic) actualizarea unui fișier text ce conține lista vorbitorilor; (iv) extragerea transcrierilor în format text din fișierele HTML. Detalii suplimentare despre achiziția setului de date CDep-raw se găsesc în [Manolache, 2019], capitolul 4.

Setul de date brut CoBiLiRo-raw conține 76 de fișiere audio împreună cu transcrierile aproximative corespunzătoare, extrase din emisiuni și interviuri. Acest set de date este împărțit la rândul său în 4 subseturi ce însumează aproximativ 714k de cuvinte pronunțate într-o durată de aproximativ 70 de ore. Numărul de cuvinte și de ore pentru fiecare subset se regăsesc în Tabelul 2.1.d.

### 2.1.3 Seturi de date de vorbire adnotată rezultate în această etapă a proiectului

După aplicarea celor trei metode de adnotare automată au fost obținute seturile de date din Tabelul 2.1.e.

**Tabelul 2.1.d** Seturi de date de vorbire neadnotată (+ transcrieri aproximative) utilizate ca date de intrare pentru cele trei metode de adnotare automată.

Setul de date	Sursa	Durăță	Transcrieri aproximative	Număr de vorbitori
<b>CDep-raw</b>	Parlamentul României	3,510h, 13m	25.1M cuvinte	~2,500
<b>CoBiLiRo-raw</b>	Alma <sup>1</sup>	19h, 29m	127.5k cuvinte	n/a
	VBarbu <sup>2</sup>	7h, 30m	80.8k cuvinte	
	GCVC <sup>3</sup>	16h, 49m	324.5k cuvinte	
	Ro100 <sup>4</sup>	25h, 35m	181.8k cuvinte	

<sup>1</sup> Seria de emisiuni Alma Mater Iassiensis

<sup>2</sup> Interviu cu prof. Viorel Barbu

<sup>3</sup> Seria de emisiuni Ghici cine mai vine la cină

<sup>4</sup> Seria de emisiuni România 100: Iașul în arcul timpului



**Tabelul 2.1.e** Seturile de vorbire adnotată rezultate în urma aplicării metodelor de adnotare automată

Setul de date	Sursa	Durată		Eficiență aliniere [% ore]	
<b>CDep-trans-v4</b>	Parlamentul României	878h, 48m		25.0%	
<b>CoBiLiRo-trans-v4</b>	Alma	10h, 24m	31h, 30m	53.3%	45.4%
	VBarbu	1h, 24m		18.7%	
	GVCV	6h, 36m		39.3%	
	Ro100	13h, 6m		51.2%	

## 2.2 Activitatea 3.9 - Analiza impactului utilizării de RAV complementare pentru generarea de adnotări în contextul îmbunătățirii sistemelor de RAV

Activitatea 3.9/2020 a avut două subactivități principale. În primul rând ne-am propus dezvoltarea suplimentară a metodei de generare automată de adnotări folosind sisteme RAV complementare. Astfel am încercat adaptarea și aplicarea metodei proiectate în activitățile anterioare folosind un sistem de RAV complet nou, bazat pe platforma ESPnet, ca fiind unul dintre cele două sisteme RAV complementare. Secțiunea 2.2.1 prezintă pașii urmați pentru dezvoltarea acestui sistem.

În al doilea rând am efectuat un studiu comparativ asupra performanțelor sistemului de RAV inițial (la începutul proiectului ReTeRom) și a sistemelor de RAV rezultate în cadrul activităților A1.13/2018 și A2.13/2019, ca urmare a proiectării și aplicării metodei de generare automată de adnotări folosind sisteme RAV complementare. Metoda se bazează pe folosirea a două sisteme RAV cât mai diferite, ce produc transcrieri pentru un corpus neadnotat, considerând ulterior părțile identice din transcriere ca fiind corecte. Această presupunere este certificată prin verificarea complementarității sistemelor; cele două sisteme diferă din punct de vedere constructiv astfel încât ele produc erori diferite și necorelate. În final, setul de date adnotat este utilizat pentru reantrenarea celui mai performant sistem inițial.

### 2.2.1 Dezvoltarea sistemului de transcriere de vorbire ESPnet pentru limba română

Pentru dezvoltarea noului sistem de RAV pentru limba română, am optat să folosim ESPnet [Watanabe, 2018], un utilitar de învățare profundă modern ce permite realizarea unui sistem de RAV complet sub forma unei rețele neurale unice (arhitectură de tip end-to-end). Biblioteca de învățare profundă ESPnet poate fi folosită pentru o gamă întreagă de aplicații de inteligență artificială pornind de la recunoașterea automată a vorbirii și sinteză de vorbire pornind de la text și mergând, ceva mai recent, și spre traducerea automată a vorbirii (en: speech translation) între multiple limbi de circulație internațională, dar și conversia vorbirii (en: voice conversion).

Utilitarul include exemple de rețete (directoare de proiect) pentru limba engleză, de la care se poate porni dezvoltarea și cercetarea în domeniu, cele mai multe fiind bazate pe resurse de antrenare audio și text disponibile în licențe de tip sursă deschisă (e.g. LibriSpeech [Panayotov, 2015], CommonVoice [Ardila, 2020]), astfel încât orice cercetător să poată ușor replica rezultatele unor articole state-of-the-art.

ESPnet folosește, la rândul ei, bibliotecile Chainer [Tokui, 2019] și PyTorch [Paszke, 2019] pentru implementarea funcționalităților legate de învățarea profundă (en: deep learning), lucru care simplifică extinderea codului și arhitecturii de procesare. Ambele biblioteci sunt folosite de cercetătorii din domeniul învățării automate, pentru aplicații în procesarea imaginilor și a limbajului natural.

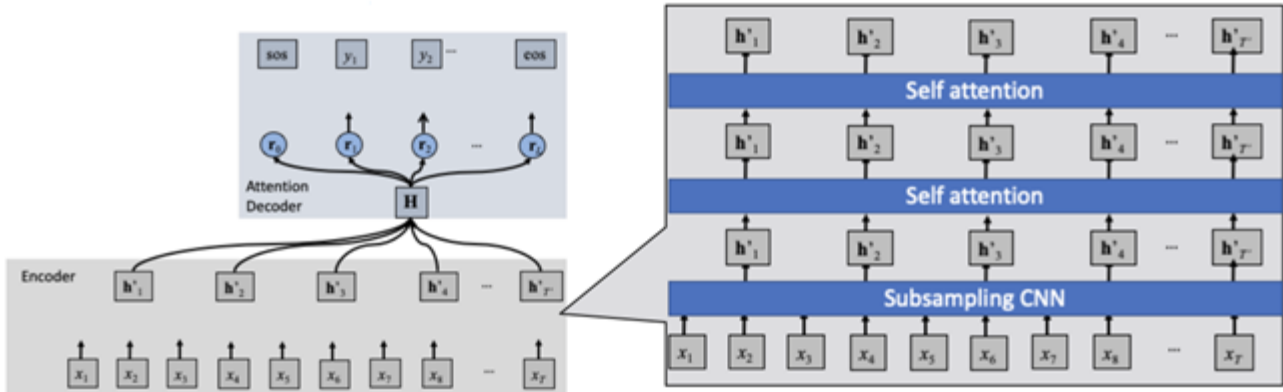
De asemenea, un alt mare avantaj al ESPnet constă în procesarea datelor audio de intrare în stil Kaldi (cel mai popular utilitar de învățare profundă pentru RAV [Povey, 2011]) astfel încât extragerea parametrilor vocali și formatul fișierelor caracteristicilor extrase se păstrează, iar portarea directoarele ce conțin resurse audio și text între cele două sisteme fiind astfel extrem de facilă.

### Arhitectura sistemului RAV propus

Sistemul RAV propus este bazat pe arhitectura end-to-end din ESPnet, folosind atât clasificarea temporală conexionistă (CTC) cât și rețeaua codor-decodor bazată pe Transformer (en: self-attention encoder-decoder) [Kim, 2017; Watanabe, 2017]. Această a doua metodă utilizează un mecanism de auto-atenție, pentru a

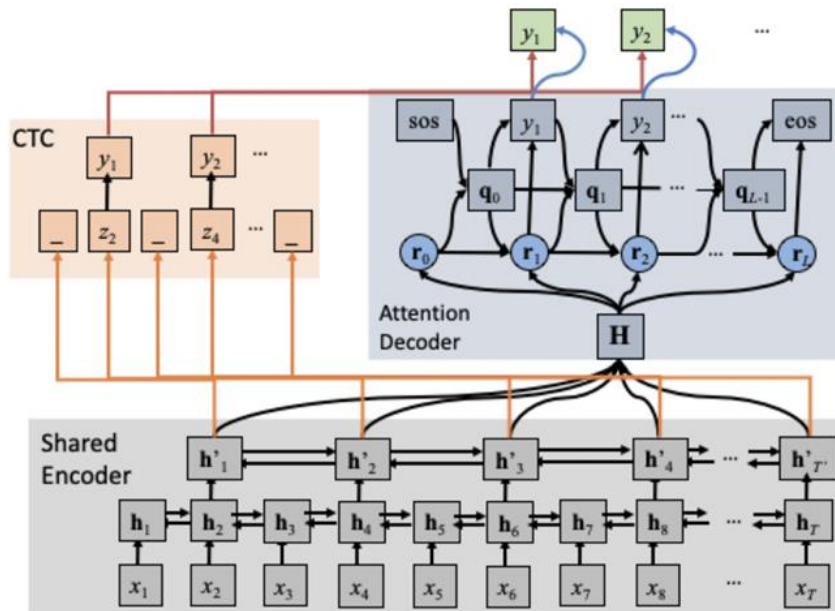
efectua alinierea între cadrele acustice și simbolurile recunoscute, în timp ce CTC folosește ipotezele Markov pentru a rezolva eficient problemele secvențiale prin intermediul programării dinamice.

Astfel, am adoptat pentru RAV o rețea hibridă de tip CTC/Transformer end-to-end (Figura 2.2.a), care utilizează în mod eficient avantajele ambelor arhitecturi în antrenare și decodare. În timpul antrenării, folosim framework-ul de învățare multi-obiectiv pentru a îmbunătăți robustețea aliniierilor problematice cât și pentru a obține o convergență mai rapidă.



**Figura 2.2.a.** Arhitectura rețelei CTC-Transformer utilizată, unde toate conexiunile recurente din codor-decodor-ul bazat pe atenție sunt înlocuite cu un bloc de auto-atenție (poate capta interdependențe pe distanțe foarte lungi). [Hori, 2019]

În timpul decodării, efectuăm inferența atât prin combinarea scorurilor bazate pe atenție, cât și a scorurilor CTC (Figura 2.2.b), într-un algoritmul de căutare a fasciculului (en: beam-search) cu o singură trecere, pentru a elimina în continuare aliniierile problematice. În plus față de arhitectura end-to-end prezentată mai sus realizată doar pe fișierele audio, am realizat o serie de experimente și cu modele de limbă, de tip Transformer.



**Figura 2.2.b.** Arhitectura rețelei CTC/attention, antrenarea fiind bazată pe învățarea multi-obiectiv iar inferența bazată pe combinarea scorurilor.

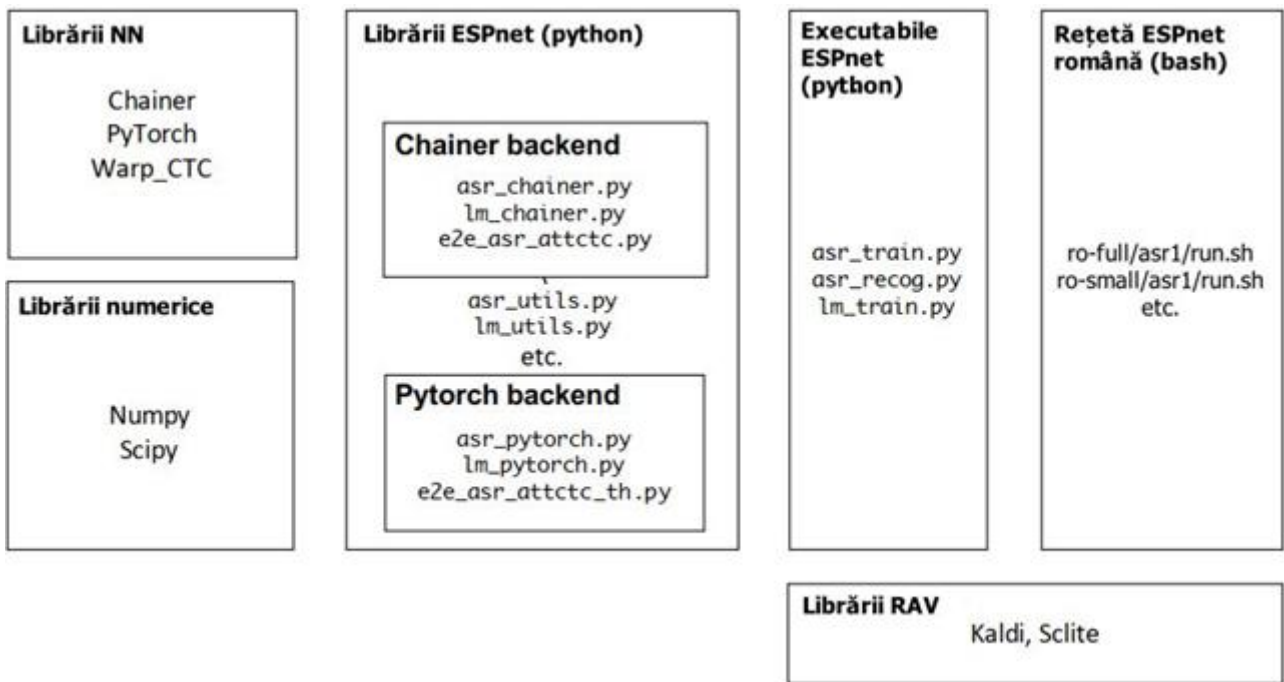


Figura 2.2.c. Arhitectura software a sistemului RAV în limba română.

### Arhitectura software a sistemului ESPnet

Figura 2.2.c prezintă arhitectura software a sistemului de recunoaștere în limba română, construit în ESPnet. Componentele principale pentru antrenarea și inferența unei rețele neuronale sunt scrise în Python, care apelează SDK-urile Chainer și PyTorch, prin comutarea backend-ului în funcție de opțiuni. Structura de directoare în care se preprocesează datele, după cum am menționat anterior, urmează filozofia Kaldi, utilizând scripturi bash pentru lansarea diferitelor binare necesare preprocesării datelor și a augmentării lor.

În continuare, Figura 2.2.d prezintă fluxul unei rețete (en: recipe) și stagiile necesare executării cu succes a unei antrenări, apoi a unei decodări, pentru sistemul în limba română.

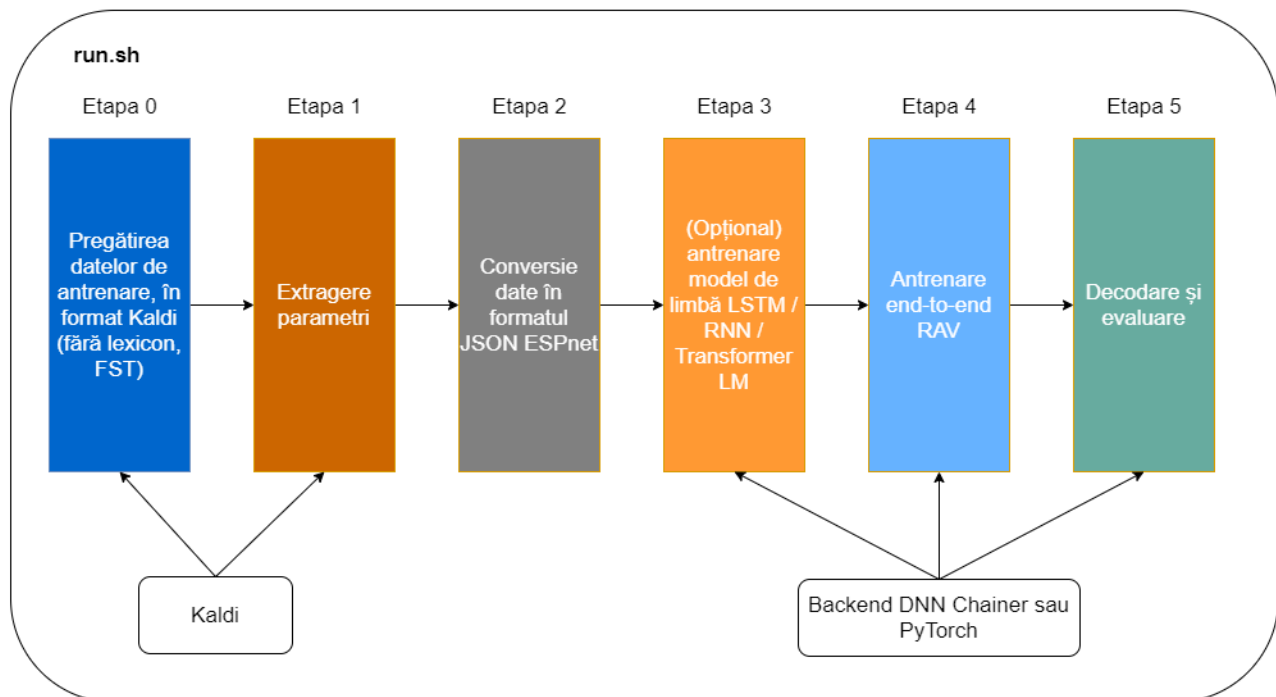


Figura 2.2.d. Fluxul de execuție al rețetei ESPnet în Limba Română.

Rețeta este semnificativ simplificată grație avantajului unui toolkit de RAV end-to-end, astfel, spre diferență de Kaldi, nu trebuie să includem un lexicon, compilarea transductorului de stare finită (FST), antrenarea / alinierea bazată pe modelarea mixturii HMM și Gaussiene, precum și generarea structurilor de tip latice pentru antrenarea secvențială discriminativă.

Etapele principale necesare fluxului unui sistem în limba română sunt:

- Etapa 0. Pregătirea datelor. Pentru această etapă a fost utilizat formatul Kaldi de stocare a datelor de antrenare (același format utilizat și pentru celelalte sisteme de RAV ale Speed).
- Etapa 1. Extragere parametri. Din datele de intrare se extrag parametri de vorbire tipici, de tip MFCC, însă de înaltă rezoluție: 80 de coeficienți cepstrali; vectorul de parametri corespunzător fiecărei ferestre din semnalul de vorbire va avea 80 sau 83 de dimensiuni în funcție de cum la cei 80 de coeficienți cepstrali se adaugă sau nu și informații despre frecvența fundamentală.
- Etapa 2. Conversie formate în ESPnet. Această etapă convertește informația din structura de directoare Kaldi (transcrieri, id-uri vorbitor și limbă, lungimea vectorilor de intrare) într-un singur JSON (data.json). Caracteristicile extrase din fișierele audio rămân în format Kaldi.
- Etapa 3. Antrenare model de limbă. Pas opțional, necesar dacă se dorește utilizarea unui model de limbă suplimentar pentru reevaluare lingvistică. În experimentele noastre, am încercat ambele variante: sistem RAV end-to-end și sistem RAV cu reevaluare lingvistică.
- Etapa 4. Antrenare end-to-end. În această etapă se antrenează model codor-decodor hibrid de tip CTC-Transformer, folosind una dintre cele două biblioteci de învățare profundă PyTorch sau Chainer.
- Etapa 5. Decodarea și evaluarea. În această ultimă etapă se decodează setul de date de evaluare cu sistemul RAV nou antrenat format din modelul end-to-end (etapa 4) și, eventual, modelul de limbă suplimentar (etapa 3). După decodare, se realizează alinierea transcrierii ipotetice cu cea de referință și se calculează rata de eroare la nivel de cuvânt (WER).

### ***Setup-ul experimental pentru experimentele de RAV***

Experimentele realizate au urmărit dezvoltarea unui sistem de RAV bazat pe ESPnet în aceleași condiții în care a fost dezvoltat cel mai bun sistem de RAV Speed din anul 2019. Motivația acestei alegeri a fost bineînțeles dorința realizării unei comparații relevante între sistemul de RAV actual și noul sistem dezvoltat cu ESPNet.

Astfel, au fost utilizate seturile de date listate în continuare și descrise pe larg în secțiunea 2.1.1 a acestui raport:

- RSC-train și SSC-train1+2 - set de date de antrenare restrâns utilizat în etapa de calibrare a hiperparametrilor;
- RSC-train, SSC-train1+2, SSC-train3+4-trans-v4 - set de date de antrenare utilizat pentru antrenarea sistemului final, după calibrarea hiperparametrilor;
- RSC-eval, SSC-eval1 și SSC-eval2 - set de date de evaluare;
- news2014 și talkshows - set de date de text utilizat pentru antrenarea modelului de limbă.

Pentru toate experimentele, setul de dezvoltare necesar în antrenarea rețelei neuronale ESPnet a fost selectat aleatoriu din setul de date de antrenare. Setul de dezvoltare reprezintă 10% din setul de date de antrenare și a fost extras din acesta înainte de începerea procesului de antrenare.

### ***Calibrarea hiperparametrilor sistemului de RAV***

Utilitarul ESPNet include o serie rețete de antrenare pentru sisteme de RAV în limba engleză. Dintre acestea, cele mai populare sunt rețetele Librispeech, TED-LIUM v1, WSJ și CommonVoice, rețete ce folosesc seturile de date cu același nume pentru antrenarea modelului de RAV. În ceea ce privește calibrarea hiperparametrilor sistemului, abordarea generală pe care am ales-o a fost următoarea:

- am analizat rețetele menționate mai sus și am identificat hiperparametri ce diferă de la caz la caz;
- hiperparametri care nu diferă la diversele rețete pentru limba engleză au primit aceleași valori și în experimentul pe limba română;

- pentru hiperparametri ale căror valori diferă de la o rețetă la alta, am ales pentru limba română o valoare adaptată dimensiunii setului de date de antrenare pe care îl avem la dispoziție, ținând cont, bineînțeles, și de particularitățile limbii române.

Pentru exemplificarea procedurii de mai sus, luăm ca exemplu dimensiunea vocabularului de subcuvinte (*i.e.* Byte Pair Encoding [Sennrich, 2016]). ESPnet nu transcrie vorbirea direct în cuvinte, ci în unități lingvistice mai scurte numite informal subcuvinte. Acestea sunt pur și simplu secvențe de litere ce apar în componența cuvintelor, însă fără să aibă vreo semnificație lingvistică anume (*e.g.* nu reprezintă prefixe, sufixe etc.). Subcuvintele se extrag din corpusul de date text de antrenare pe baza frecvenței lor de apariție. Dacă vocabularul de subcuvinte este limitat la dimensiunea  $N$ , atunci, folosind algoritmul prezentat în [Sennrich, 2016] se vor alege cele mai frecvente  $N$  subcuvinte cu care se pot genera toate cuvintele care apar în respectivul set de date.

În rețetele LibriSpeech, TED-LIUM v1, respectiv CommonVoice, dicționarul de BPE-uri este generat pe baza transcrierilor fișierelor de vorbire, transcrieri ce cuprind 9.4M cuvinte, 6.6M cuvinte, respectiv 4.5M cuvinte. Dimensiunile vocabularelor de subcuvinte variază de la 5000 (LibriSpeech) la 365 (CommonVoice). Astfel, raportul (număr subcuvinte)/(număr cuvinte în setul de antrenare) variază între 0.05% pentru LibriSpeech și 0.0075% pentru TED-LIUM. Am decis ca pentru limba română să păstrăm acest raport între limitele menționate mai sus, însă am ținut cont și de faptul că limba română este o limbă cu multe forme flexionate. Astfel, am ales să folosim un vocabular de subcuvinte de dimensiune mare: 1000 de subcuvinte, raportul (număr subcuvinte)/(număr cuvinte în setul de antrenare) situându-se la limita superioară: 0.05%.

În mod similar au fost analizate valorile hiperparametrilor pentru rețetele pentru limba engleză și au fost alese valori concrete pentru următorii hiperparametri:

- pentru straturile de atenție:
  - adim (dimensiunea unităților neurale din straturile de atenție) a fost aleasă 256;
  - aheads (numărul de capete pentru straturile de atenție) a fost ales 4;
- pentru procedura de antrenare:
  - batch-size (numărul de rostiri per mini-batch) a fost limitat la 32, ca urmare a limitării memoriei RAM de doar 12 GB disponibile în procesoarele grafice Nvidia Tesla K40m, pe care s-a realizat antrenarea;
  - epochs (numărul de epoci de antrenare) a fost limitat la 60, deși ar fi fost poate util să antrenăm modelul timp de 120 de epoci. Alegerea a fost influențată de timpul mare de antrenare pe hardware-ul disponibil;
  - transformer-lr (rata de învățare a modelului de tip transformer) a fost aleasă 5, ca un compromis între ratele folosite în rețetele pentru limba engleză
- pentru optimizarea obiectivului:
  - accum-grad (numărul gradienti însumați înainte de o ajustare a ponderilor rețelei) a fost ales 3 pentru că antrenarea s-a realizat în paralel pe 3 plăci grafice, astfel că ne-am permis să ajustăm ponderile rețelei după 3 treceri forward-backward realizate în paralel.
  - batch-bins (numărul de ferestre de vorbire ce intră în componența fiecărui mini-batch) a fost setat la 12M, ca urmare a limitării memoriei RAM de doar 12 GB disponibile în procesoarele grafice Nvidia Tesla K40m, pe care s-a realizat antrenarea;
- pentru modelul de limbă (de tip Transformer):
  - att-unit (numărul de unități de atenție) a fost ales 256;
  - head (numărul de capete pentru straturile de atenție) a fost ales 2;
  - layer (numărul de straturi transformer) a fost ales 4;
- pentru procedura de antrenare a modelului de limbă:
  - batchsize (numărul de propoziții per min-batch) a fost ales 64;
  - epoch (numărul de epoci de antrenare) a fost limitat la 20, deși ar fi fost poate util să antrenăm modelul timp de 50 de epoci. Alegerea a fost influențată de timpul mare de antrenare pe hardware-ul disponibil.

După antrenarea folosind parametri evidențiați anterior, au fost realizate un set de experimente pentru stabilirea valorilor optime pentru hiperparametri procesului de decodare:

- LMW (language model weight – ponderea scorului dat de modelul de limbă);
- CTCW (CTC weight – ponderea scorului dat de CTC).

## Rezultatele experimentale

Rezultatele experimentale obținute sunt evidențiate în Tabelul 2.2.a. Putem observa că cele mai bune rezultate se obțin dacă ponderile date scorurilor provenind de la modelul de limbă suplimentar și obiectivul CTC sunt egale (și egale cu 0.5). De asemenea, putem observa că acest prim sistem RAV ESPnet se apropie de rezultatele obținute cu sistemul RAV dezvoltat cu Kaldi, însă nu este încă la fel de performant. Bineînțeles, trebuie luat în considerare și faptul că atât modelul end-to-end de transcriere de vorbire, cât și modelul lingvistic suplimentar au fost antrenate un număr de epoci relativ mic (mai puțin de 50% din cât ar fi trebuit) din cauza timpului lung de antrenare pe sistemele hardware avute la dispoziție. Într-un viitor apropiat, urmează să fie efectuate mai multe experimente cu sistemul ESPnet, iar aceste rezultate preliminare sunt îmbucurătoare și ne permit să fim încrezători că vom putea egala performanțele sistemului dezvoltat cu Kaldi.

Din păcate, nu se poate afirma același lucru și despre durata proceselor de antrenare și decodare. Procesul de antrenare al sistemului RAV bazat pe ESPnet durează de aproximativ 5 ori mai multe decât un proces de antrenare similar pentru un sistem RAV Kaldi. În ceea ce privește procesul de decodare, experimentele noastre au arătat că sistemul ESPnet decodează o oră de vorbire în 2 până la 7 ore, în funcție de cât de similară este acustica respectivei ore de vorbire raportat la setul de date de antrenare. De partea cealaltă, decodarea unei ore de vorbire cu sistemul RAV bazat pe Kaldi durează doar 1 minut. Acest lucru ne-a permis să folosim sistemele bazate pe Kaldi în cadrul proiectului ReTeRom pentru a transcrie seturi de date mari (peste 500 de ore) de vorbire neadnotată și apoi să proiectăm și să aplicăm diverse metode pentru a obține un subset de date transcris corect (cu o acuratețe de peste 99%). Din cauza timpului extrem de lung de decodare pentru sistemul RAV bazat pe ESPnet acest sistem nu poate fi utilizat în mod similar în adnotarea automată a seturilor de date de vorbire.

**Tabelul 2.2.a.** Comparatie între sistemul RAV ESPnet și sistemul RAV Kaldi. Optimizarea hiperparametrilor de decodare pentru sistemul RAV ESPnet în limba română.

Set de antrenare de vorbire adnotată	Model lingvistic	Parametri decodare		WER [%]		
		LMW	CTCW	RSC_eval	SSC_eval1	SSC_eval2
RSC-train + SSC-train1+2	Nu	n/a	n/a	13.3	25.1	78.0
	Da	0.5	0.5	8.8	<b>21.3</b>	<b>38.9</b>
	Da	0.6	0.4	<b>8.7</b>	21.6	46.5
	Da	0.7	0.3	9.0	23.0	64.3
	Da	0.8	0.2	11.2	25.8	84.6
RSC-train + SSC-train1+2 + SSC-train3-trans-v4 + SSC-train4-trans-v4	Da	0.6	0.4	3.4	15.3	23.5
Sistem RAV baseline bazat pe Kaldi (antrenat pe același set de date ca mai sus, folosind model de limbă pentru reevaluare lingvistică)				1.8	11.0	14.0

### 2.2.2 Dezvoltarea metodei de adnotare bazate pe sisteme RAV complementare și rezultatele obținute pe parcursul proiectului

În această secțiune prezentăm un studiu comparativ asupra performanțelor sistemului de RAV inițial (la începutul proiectului ReTeRom) și a sistemelor de RAV rezultate în cadrul activităților A1.13/2018 și A2.13/2019, ca urmare a proiectării și aplicării metodei de generare automată de adnotări folosind sisteme RAV complementare.

Activitatea A1.13 din etapa 1/2018 a presupus utilizarea a două sisteme RAV inițiale pentru obținerea în mod automat de transcrieri cu grad ridicat de încredere, prin metoda sistemelor complementare. Sistemul inițial RAV #1 [Georgescu, 2017; Georgescu, 2018] a fost dezvoltat cu utilitarul CMU Sphinx [Huggins-Daines, 2006]. Modelele sistemului sunt probabilistice, de tip HMM-GMM (modelul acustic), respectiv 2-gram (modelul lingvistic). Acest sistem avea un vocabular de 64k cuvinte. Sistemul inițial RAV #2 [Georgescu, 2018] a fost dezvoltat cu utilitarul Kaldi [Povey, 2011], unde modelul acustic este unul de tip rețea neuronală cu întârziere în timp (time-delay - TDNN) specifică implementării NNET2. Modelul lingvistic este unul probabilistic, de tip 2-gram, folosit la decodare, în timp ce un model superior, de tip 4-gram, a fost folosit pentru reevaluarea lingvistică, corectând transcrierea inițială.

Cele două sisteme au fost utilizate pentru a transcrie seturile de vorbire neadnotată SSC-train3-raw, ce însumează 136 ore, și SSC-train4-raw, cu o durată de 777 ore. În urma transcrierii și aplicării metodei, s-au obținut seturile de date SSC-train3-compl-2018, cu o dimensiune de 49 ore, și SSC-train4-compl-2018, cu o durată de 280 ore. Calitatea acestor noi seturi obținute a fost verificată prin aplicarea metodei asupra unor seturi de date pentru care există transcriere de referință: RSC-eval și SSC-eval. S-a constatat că selecțiile rezultate sunt greșite într-o foarte mică proporție, numai 1.0% - 1.3%. Prin extrapolare, putem considera că și noile date sunt corecte în mare măsură (aproximativ 99%), fiind astfel posibil ca acestea să fie folosite mai departe ca seturi de antrenare adiționale.

Tabelul 2.2.a prezintă pe primele două linii caracteristicile și performanțele sistemelor RAV inițiale, în timp ce a treia linie prezintă caracteristicile și performanțele sistemului RAV #2 îmbunătățit, pentru a cărui antrenare au fost adăugate noile date obținute. Acest nou sistem a obținut o îmbunătățire relativă de 3.81% pe setul de evaluare cu vorbire citită, respectiv 8.87% pe setul de evaluare cu vorbire spontană.

Activitatea 2.13 din etapa 2/2019 a presupus aplicarea metodei sistemelor complementare, de data aceasta, cele două sisteme inițiale fiind sistemul RAV #2 prezentat în Activitatea 1.13 din 2018 și sistemul RAV #3, având modelul acustic de tip time-delay (TDNN), creat cu utilitarul Kaldi folosind implementarea NNET3. Modelul de limbă folosit pentru decodare este un 2-gram, în timp ce modelul de limbă folosit pentru reevaluarea lingvistică este unul de tip rețea neuronală recurentă (RNN), cu un istoric de 5 cuvinte. Ambele modele de limbă folosesc un vocabular de 200k cuvinte.

Metoda a fost aplicată pe aceleași seturi de date de vorbire neadnotată: SSC-train3-raw și SSC-train4-raw, cu o durată de 777 ore. În urma aplicării metodei, s-au obținut seturile de date SSC-train3-compl-2019 cu o dimensiune de 79 ore, și SSC-train4-compl-2019, cu o durată de 452 ore. Calitatea acestor noi date a fost verificată prin aplicarea metodei asupra seturilor de referință RSC-eval și SSC-eval. Selecțiile rezultate au fost greșite în proporție de 2.6% - 2.7%, mai mult decât în cazul sistemelor RAV #1 și RAV #2. Acest fapt este pus pe seama diferențelor mai mici dintre cele două sisteme, ambele fiind dezvoltate cu ajutorul utilitarului Kaldi, iar modelele acustice sunt asemănătoare din punct de vedere al tehnologiei folosite.

**Tabelul 2.2.a** Performanța sistemelor RAV inițiale și a sistemului RAV îmbunătățit din A1.13, etapa 1.

Model acustic		Model lingvistic	WER [%]		Îmbunătățire relativă a WER [%]	
Corpus antrenare	Tip model		RSC-eval	SSC-eval	RSC-eval	SSC-eval
RSC-train + SSC-train	HMM-GMM	Decodare RAV: 64k cuvinte, 3-gram	9.56	28.64	-	-
RSC-train + SSC-train	HMM-DNN (TDNN2)	Decodare RAV: 200k cuvinte, 2-gram Reev. lingv.: 200k cuvinte, 4-gram	3.41	17.91	-	-
RSC-train + SSC-train + SSC-train3-compl-2018 + SSC-train4-compl-2018	HMM-DNN (TDNN2)	Decodare RAV: 200k cuvinte, 2-gram Reev. lingv.: 200k cuvinte, 4-gram	3.28	16.32	3.81	8.87

**Tabelul 2.2.b** Performanța sistemelor RAV inițiale și a sistemului RAV îmbunătățit din A2.13, etapa 2.

Model acustic		Model linvistic	WER [%]		Îmbunătățire relativă a WER [%]	
Corpus antrenare	Tip model		RSC-eval	SSC-eval	RSC-eval	SSC-eval
RSC-train + SSC-train	HMM-DNN (TDNN2)	Decodare RAV: 200k cuvinte, 2-gram Reev. lingv.: 200k cuvinte, 4-gram	3.41	17.91	-	-
RSC-train + SSC-train	HMM-DNN (TDNN3)	Decodare RAV: 200k cuvinte, 2-gram Reev. lingv.: RNN 5-gram	1.88	14.96	-	-
RSC-train + SSC-train + SSC-train3-compl-2018 + SSC-train4-compl-2018	HMM-DNN (TDNN3)	Decodare RAV: 200k cuvinte, 2-gram Reev. lingv.: RNN 5-gram	1.80	11.70	4.26	21.79
RSC-train + SSC-train + SSC-train3-compl-2019 + SSC-train4-compl-2019	HMM-DNN (TDNN3)	Decodare RAV: 200k cuvinte, 2-gram Reev. lingv.: RNN 5-gram	1.87	11.78	0.53	21.26

Noile date obținute au fost adăugate la setul de antrenare deja existent, fiind creat sistemul RAV #3 îmbunătățit. Tabelul 2.2.b prezintă pe primele 2 linii sistemele inițiale RAV #2 și RAV #3. Următoarele două linii prezintă două versiuni ale sistemului RAV #3 îmbunătățit: (i) una antrenată după ce au fost adăugate datele SSC-train3-compl-2018 și SSC-train4-compl-2018 la setul de date inițial, respectiv (ii) una antrenată folosind seturile SSC-train3-compl-2019 și SSC-train4-compl-2019 suplimentar față de setul inițial. Urmărind rezultatele din acest tabel, putem trage două concluzii interesante:

- seturile de date obținute în 2018 sunt mai mici, însă probabil mai precise (sistemele inițiale din 2018 erau mai diferite decât sistemele din 2019) astfel că sistemul RAV reantrenat folosind seturile din 2018 este puțin mai bun decât sistemul RAV reantrenat folosind seturile generate în 2019.
- seturile de date nou obținute sunt mult mai utile pentru dezvoltarea unui sistem ce se dorește a transcrie corect vorbire spontană: îmbunătățirea relativă a WER-ului este de peste 20% pentru vorbire spontană (SSC-eval) și oarecum nesemnificativă pentru vorbire citită (RSC-eval).

Deși activitățile din etapa 3/2020 nu prevedeau acest lucru, în cursul acestei etape s-a încercat dezvoltarea suplimentară a metodei utilizând ca sisteme inițiale sistemul RAV #3 (2019) și un nou sistem de RAV dezvoltat folosind utilitarul ESPnet (vezi secțiunea anterioară). În acest demers am întâmpinat însă o problemă ce nu a putut fi încă depășită. Sistemul RAV bazat pe ESPnet s-a dovedit a fi foarte lent în transcriere, astfel încât nu a putut fi utilizat pentru transcrierea seturilor de date neadnotate SSC-train3-raw și SSC-train4-raw. Transcrierea acestor seturi de date, cu o durată totală de 913 ore, ar fi necesitat aproximativ 2700 de ore, adică peste 100 de zile. Pentru comparație, menționăm că transcrierea aceluiași seturi de vorbire cu sistemul RAV bazat pe Kaldi a necesitat numai 16 ore. În urma acestei experiențe am tras concluzia că pentru aplicarea cu succes a metodelor proiectate în cadrul proiectului nu este suficient să existe sisteme de RAV inițiale performante și complementare, ci, în plus, acestea trebuie să transcrie vorbirea suficient de rapid.

## 2.3 Activitățile 3.10 și 3.12 - Îmbunătățirea soluției de filtrare și aliniere a transcrierilor aproximative cu semnalul de vorbire și Analiza impactului utilizării transcrierilor aproximative în vederea reantrenării sistemelor de RAV

În această secțiune vom aborda următoarele subiecte: (i) îmbunătățirea metodei de utilizare a transcrierilor aproximative pentru generarea de adnotări pentru seturi de date de vorbire și (ii) aplicarea metodei pe două noi seturi de date brute: corpusul CoBiLIRo și corpusul CDep, urmată de evaluarea impactului seturilor de date nou adnotate asupra sistemului de RAV inițial.

### 2.3.1 Îmbunătățirea soluției de aliniere a transcrierilor aproximative cu semnalul de vorbire

#### Metoda de bază

Metoda de bază de filtrare și aliniere a transcrierilor aproximative cu semnalul de vorbire a fost introdusă, prezentată și evaluată în activitatea 2.11 din etapa 2019 a proiectului (vezi raport TADARAV 2019 [Georgescu, 2019a], secțiunea 2.2). În cele ce urmează este prezentată o descriere sumară a acestei metode. Obiectivul principal al metodei de bază de aliniere este obținerea adnotărilor precise pentru o parte a unui



corpus de vorbire în mod automat. Noul corpus obținut va fi folosit pentru reantrenarea sistemelor RAV existente crescând astfel variabilitatea modelului acustic, ceea ce ar trebui să implice mai departe o îmbunătățire în acuratețea transcrierilor. Metoda de adnotare se bazează pe folosirea unui corpus de vorbire brut împreună cu transcrieri aproximative și transcrierile produse de un sistem RAV inițial. Cele două seturi de transcrieri ale corpusului de vorbire sunt aliniate folosind distanța Levenstein. În urma alinierii, considerăm părțile identice dintre cele două transcrieri ca fiind corecte și le selectăm ca parte a noului corpus.

Părțile de vorbire și de transcrieri aproximative au fost extrase din mediul online mass-media (știri, interviuri, rapoarte). Deși mediul online mass-media este o sursă bogată de vorbire și text, transcrierile aproximative extrase sunt diferite de cele generate de sistemul RAV în ceea ce privește formatul. Transcrierile aproximative conțin titluri, nume proprii, semne de punctuație, numere, abrevieri etc., pe când transcrierile RAV sunt practic secvențe de cuvinte cu litere mici însoțite de etichete de timp. Ca și exemplu:

- Transcriere aproximativă: Bărbatul de 36 de ani povestește că muncise toată noaptea.
- Transcriere RAV: bărbatul(3.71, 4.02) de(4.02, 4.87) treizeci(4.87, 5.11) și(5.11, 5.74) șase(5.74, 6.02) de(6.02, 6.35) ani(6.35, 6.71) povestește(6.71, 6.94) că(6.94, 7.58) muncise(7.58, 7.89) toată(7.89, 8.16) noaptea(8.16, 8.32)

Astfel, transcrierile aproximative trebuie supuse unor operații de procesare de text, anume: restaurare de diacritice, înlocuirea URL-urilor și adreselor de email cu forma lor vorbită, înlocuirea numerelor cu text, expandarea abrevierilor, înlocuirea caracterelor speciale, transformarea literelor mari în litere mici.

După procesul de aliniere, segmentele de text obținute sunt filtrate pe baza anumitor criterii cu scopul de a înlătura segmente foarte mici (în ceea ce privește durata vorbirii și numărul de caractere) și segmente ce conțin zone lungi de liniște (durata dintre cuvinte consecutive). În final, procesele de aliniere și filtrare generează un set de transcrieri aliniate însoțite de etichete de timp. Etichetele de timp vor fi folosite mai departe pentru tăierea segmentelor audio corespunzătoare din materialele audio. Astfel, noul corpus de vorbire adnotată va conține transcrieri aliniate și segmentele audio corespunzătoare, ce pot fi folosite mai departe pentru reantrenarea sistemelor RAV.

### Metoda propusă

Metoda descrisă anterior este foarte strictă în ceea ce privește corectitudinea transcrierilor aliniate, în sensul că sunt considerate a fi corecte doar secvențele de cuvinte unde transcrierile aproximative și transcrierile RAV se potrivesc perfect. Acest fapt duce la selectarea multor secvențe scurte (*i.e.* ce cuprind puține cuvinte). În urma unei analize manuale a transcrierilor s-a constatat că există multe secvențe de câteva cuvinte pentru care procesul de aliniere a eșuat și care se află între secvențe lungi de cuvinte care s-au aliniat cu succes. Câteva astfel de exemple sunt prezentate în Tabelul 2.3.a.

Aceste secvențe nu au fost aliniate din cauza mai multor motive precum zgomot de fundal în materialul audio, cuvinte care nu se regăsesc în vocabularul sistemului RAV (de regulă nume proprii), cuvinte adiționale în fișierul audio care nu sunt pronunțate în materialul audio. Cel de-al treilea caz se întâlnește de obicei în interviuri, unde transcrierea de pe website conține text adițional pentru a preciza cine vorbește.

**Tabelul 2.3.a.** Exemple de secvențe de cuvinte pentru care procesul de aliniere a eșuat și care se află între secvențe lungi de cuvinte care s-au aliniat cu succes. Secvențele pentru care procesul de aliniere a eșuat sunt marcate cu roșu/verde având în vedere dacă transcrierea aproximativă este incorectă/corectă.

Nr .	Secvența anterioară (alinată)	Secvență nealinată	Secvența următoare (alinată)	Observații
1	... se îndreaptă către susținătorii săi	o dronă	filmează evenimentul ...	confuzat cu "un pod" de către sistemul RAV din cauza zgomotului de fundal
2	... în zona	Egiptului	unde există ...	confuzat cu "edituri" de către sistemul RAV deoarece acest cuvânt nu se află în vocabular
3	... a obținut fondurile	administrator bloc	la doi ani ...	secvența de text nu este pronunțată în materialul audio, scopul acesteia în transcrierea aproximativă este aceea de a specifica cine vorbește

În urma unei analize a 50 de astfel de situații, am ajuns la concluzia că dacă am lua în considerare secvențe de mai puțin de 6 cuvinte pentru care procesul de aliniere a eșuat, în aproximativ 68% din cazuri transcrierea aproximativă este corectă, pe când transcrierea RAV este incorectă. Astfel, în aceste cazuri am decis să unim cele 3 secvențe de cuvinte. Adăugând aceste secvențe de text sperăm să generăm un corpus de vorbire mai mare ce va fi folosit pentru reentrenarea sistemelor RAV. Noi credem că aceste secvențe vor fi mai valoroase pentru reentrenare, deoarece sistemul RAV inițial nu a reușit să le transcrie. Desigur, există și o parte negativă la acest compromis: 32% din secvențele adăugate sunt incorecte și pot afecta în mod negativ procesul de antrenare.

### **Evaluarea metodelor**

O metodă de adnotare automată poate fi evaluată în funcție de dimensiunea corpusului nou generat și în funcție de calitatea transcrierilor, exprimată în funcție de rata de eroare la nivel de cuvânt (WER). Dimensiunea corpusului nou generat se dorește a fi pe cât de mare posibilă. Desigur, limita superioară este dimensiunea corpusului audio brut (materialele audio). Astfel, o metrică mai potrivită este eficacitatea metodei, definită ca procentajul de material audio ce a fost selectat ca făcând parte din corpusul nou generat și pentru care metoda a generat transcrieri corecte. Calitatea adnotării exprimată în rata de eroare la nivel de cuvânt (WER) și/sau rata de eroare la nivel de caracter (ChER) nu poate fi măsurată pentru metoda de bază din cauza absenței unei referințe pentru transcrierile aproximative. Totuși, evaluarea se poate face pe alt criteriu care se bazează pe această metrică. O adnotare de calitate înaltă ar trebui să ducă în mod normal la o performanță mai bună a sistemului RAV reentrenat pe noul corpus de vorbire generat.

### **Rezultate alinieri**

Procedura de aliniere și filtrare prezentată în secțiunea “Metoda de bază” a fost aplicată pe seturile de date brute SSC-train3-raw și SSC-train4-raw. În etapa 2/2019 am folosit ca punct de pornire transcrierile acestor două seturi de date generate cu un sistem RAV inițial mai slab, și anume cel prezentat în [Georgescu, 2017]. Au rezultat astfel seturile de date de vorbire adnotată SSC-train3-trans-v3 și SSC-train4-trans-v3. Dimensiunile acestora, exprimate în număr de ore de vorbire și eficiența procesului de adnotare automată, exprimată sub forma procentului de date brute ce au putut fi adnotate, raportat la dimensiunea datelor brute sunt prezentate în Tabelul 2.3.b.

Aceeași metodă de aliniere și filtrare a fost aplicată în cadrul activității curente folosind ca punct de pornire transcrierile acestor două seturi de date generate cu un sistem RAV îmbunătățit, și anume cel prezentat în [Georgescu, 2019b]. A rezultat astfel versiunea v4 a acestor seturi de date adnotate, versiune ce cuprinde mai multe ore de vorbire, conform Tabelului 2.3.b.

În urma aplicării procedurii de aliniere și filtrare prezentată în secțiunea “Metoda propusă” pe seturile de date brute SSC-train3-raw și SSC-train4-raw, au fost obținute versiunile v5 ale seturilor de date adnotate. În acest caz a fost utilizat pentru transcriere tot sistemul de RAV îmbunătățit prezentat în [Georgescu, 2019b].

**Tabelul 2.3.b.** Statistici pentru seturile de date SSC-train3-trans-v4 și SSC-train4-trans-v4

Set de date	Durăță [# ore]	Eficiență aliniere [% ore]
SSC-train3-trans-v3	37,5	27,4%
SSC-train3-trans-v4	41,0	30,0%
SSC-train3-trans-v5	46,1	33,7%
SSC-train4-trans-v3	228,8	29,4%
SSC-train4-trans-v4	250,2	32,2%
SSC-train4-trans-v5	281,6	36,2%

Așa cum era de așteptat, seturile de date v4 sunt mai mari (cu aproximativ 9%) decât seturile de date v3. Eficiența alinierii a crescut de la 27.4% la 30.0% pentru primul set de date, respectiv de la 29.4% la 32.2% pe cel de-al doilea set de date. Ne așteptăm ca această creștere a seturilor de date să contribuie pozitiv la reentrenarea sistemului de RAV inițial, întrucât creșterea se bazează în totalitate pe niște transcrieri RAV inițiale mai precise.

De asemenea, seturile de date v5 sunt mai mari (cu aproximativ 12%) decât seturile de date v4. Eficiența alinierii a crescut de la 30.0% la 33.7% pentru primul set de date, respectiv de la 32.3% la 36.2% pe cel de-al doilea set de date. Urmează să vedem, în secțiunea următoare, dacă aceasta creștere a seturilor de date contribuie pozitiv sau negativ la reentrenarea sistemului de RAV inițial.

### **Sisteme RAV ce utilizează noile corpusuri obținute**

Sistemele RAV obținute în urma reentrenării cu noile seturi de date sunt prezentate în Tabelul 2.3.d, alături de sistemul RAV inițial (cel ce a fost utilizat pentru transcrierea datelor brute în anul 2020). Comparativ cu raportul precedent, evaluarea s-a realizat pe seturile de evaluare corectate RSC-eval-v2 și SSC-eval-v2 (vezi secțiunea 2.1.1 pentru mai multe detalii).

**Tabelul 2.3.d. Performanța sistemelor RAV după reentrenare**

Sistem RAV	Set antrenare	WER [%]		Îmbunătățire relativă a WER [%]	
		RSC-eval	SSC-eval1	RSC-eval	SSC-eval1
Sistem RAV inițial	RSC-train + SSC-train	1.88	14.96	n/a	n/a
RAV reentrenat 2019 (metoda de bază)	RSC-train + SSC-train + SSC-train3-trans-v3 + SSC-train4-trans-v3	1.83	11.50	2.66%	23.1%
RAV reentrenat 2020 (metoda de bază)	RSC-train + SSC-train + SSC-train3-trans-v4 + SSC-train4-trans-v4	1.83	11.04	2.66%	26.2%
RAV reentrenat 2020 (metoda propusă)	RSC-train + SSC-train + SSC-train3-trans-v5 + SSC-train4-trans-v5	1.83	11.13	2.66%	25.6%

În ceea ce privește rezultatele de transcriere de vorbire continuă (setul RSC-eval), observăm că indiferent de metoda de aliniere și filtrare utilizată și indiferent de calitatea transcrierilor inițiale, sistemele RAV reentrenate pe seturile inițiale și seturile nou adnotate prezintă aceeași îmbunătățire față de sistemul RAV inițial: WER de 1.83% vs. 1.88%.

Pentru vorbire spontană, metoda de bază aplicată anul trecut, folosind transcrieri inițiale de calitate mai slabă, conduce la obținerea unei îmbunătățiri relative de WER de aproximativ 23%. Aceeași metodă, aplicată în acest an, pornind de la transcrieri inițiale mai precise, conduce la obținerea unei îmbunătățiri relative de WER de aproximativ 26%. Se observă astfel importanța calității transcrierilor inițiale.

Metoda nou propusă în acest an conduce la obținerea unei îmbunătățiri relative de WER de aproximativ 25% față de sistemul RAV inițial. Se observă astfel că propunerea de a utiliza în antrenarea RAV și a secvențelor de câteva cuvinte (maxim 6) pentru care procesul de aliniere a eșuat și care se află între secvențe lungi de cuvinte care s-au aliniat cu succes, nu conduce la rezultate RAV mai bune. Din nou rezultă faptul că este extrem de importantă calitatea transcrierilor pentru datele de antrenare, în detrimentul cantității acestor date.

### **2.3.2 Aplicarea metodei pe setul de date CoBiLiRo-raw**

#### **Analiză inițială a setului de date CoBiLiRo-raw**

Metoda de bază, descrisă în secțiunea 2.3.1, a fost aplicată și pe setul de date CoBiLiRo-raw (vezi Tabelul 2.1.d din secțiunea 2.1.2), cu scopul îmbunătățirii sistemului RAV. Această metodă a fost preferată față de metoda propusă în secțiunea anterioară deoarece cea din urmă a fost evaluată cu rezultate mai slabe. În primă fază, s-a aplicat metoda pe doar 2 fișiere din setul de date, în vederea analizării eventualelor probleme. Rezultatele alinierii se pot observa în Tabelul 2.3.e, unde sunt prezentate numărul de cuvinte aliniate, numărul de cuvinte ce se regăsesc în transcrierea aproximativă, numărul de cuvinte ce se regăsesc în transcrierea RAV, precum și procentul de cuvinte aliniate raportat la numărul de cuvinte din transcrierea aproximativă.

**Tabelul 2.3.e. Transcrieri RAV obținute cu sistem RAV baseline [Georgescu, 2019b]**

Fișierul audio	Nr. cuvinte aliniate	Nr. cuvinte transcrieri aproximative	Nr. cuvinte transcrieri RAV	Aliniere [% cuvinte]
alma24	3.696	5.151	5.985	71,75%
interviu2	2.065	18.095	11.225	18,40%

După analiza rezultatelor din Tabelul 2.3.e, precum și a fișierelor în cauză, s-au constatat următoarele: (i) în cazul fișierului alma24, transcrierea RAV conține mai multe cuvinte decât cea aproximativă, deoarece în fișierul audio sunt pronunțate mai multe cuvinte decât există în transcrierea aproximativă, (ii) în cazul fișierului interviu2, transcrierea aproximativă conține foarte mult text care nu este pronunțat în fișierul audio. Există reformulări, adăugiri, precum și precizări suplimentare.

În Tabelele 2.3.f și 2.3.g sunt prezentate câteva exemple de probleme întâlnite pentru cele 2 seturi de date. În Tabelul 2.3.f, unde sunt prezentate exemple de text lipsă pentru setul de date alma24, se poate observa pe prima coloană text care este întâlnit în transcrierea aproximativă dar care este și pronunțat în fișierul audio, urmat de eticheta de timp la care se încheie pronunția textului respectiv pe coloana a 2-a. Acest text este situat anterior față de cel de pe coloana a 3-a, text care este pronunțat în fișierul audio dar nu se regăsește în transcrierea aproximativă. Pe coloanele 4 și 5 din Tabelul 2.3.f se regăsește text similar cu cel din primele 2 coloane cu excepția faptului că acest text este situat ulterior față de cel de pe coloana 3, iar eticheta de timp marchează momentul de timp la care se începe pronunția textului. În final, pe ultima coloană se regăsește diferența dintre cele 2 etichete de timp, care marchează durata de timp pentru care lipsește transcrierea de pe coloana 3. Analizând exemplele din Tabelul 2.3.f, se poate observa că există porțiuni din fișierul audio de până la aproximativ 20s pentru care nu există transcriere aproximativă.

**Tabelul 2.3.f. Analiză probleme set de date alma24**

Context anterior (există atât în audio cât și în transcrierea aproximativă)	Etichetă timp context anterior	Există în fișierul audio, dar nu există în transcrierea aproximativă	Context ulterior (există atât în audio cât și în transcrierea aproximativă)	Etichetă timp context ulterior	Delta timp fără transcriere [s]
de când realizăm această emisiune	29,28	universală de dar în colaborare cu tvr iași sau tvr iași doar cu universitatea pentru a marca împlinirea a o sută cincizeci de ani de la crearea universității moderne a fost trimisă de variate până habar nici despre istoria universității emisiuni despre facultății și cu diverse personalități	spuneam într-o discuție anterioară	49,77	20,49
nouăzeci de ani	133,26	pentru știu că sunt program complex	am publicat	142,8	9,54
eu	311,43	sunt un pic mai sus	cred că	315,45	4,02
Curtea Constituțională	413,76	nu o să vă întreb despre legea pensiilor discutăm despre facultatea de drept	Domnule profesor	420	6,24

În Tabelul 2.3.g sunt prezentate exemple de text care se regăsește în transcrierea aproximativă dar nu este pronunțat în fișierul audio. După cum se poate observa, există atât secvențe scurte de text nepronunțat cât și secvențe mai lungi.

**Tabelul 2.3.g. Analiză fișier interviu2**

Context anterior (există atât în audio cât și în transcrierea aproximativă)	Etichetă timp context anterior	Există în transcrierea aproximativă, dar nu există în fișierul audio
universitatea din freiburg	110	în toate domeniile și la toate palierele vasta sa activitate de construcție a relației dintre iași și freiburg a început
foarte importante	165,99	pentru noi profesorii și deplasări în germania pe atunci relația cu freiburg
realizările lui miron	205	poate cea mai importantă
își pot imagina	218,34	ce șansă enormă însemna pentru acei oameni simpla posibilitate de a ieși fie și pentru perioade scurte din închisoarea generalizată care era românia comunistă da acest lucru a însemnat ceva excepțional
îmi amintesc	223	că prin o mie nouă sute optzeci și cinci a avut loc în germania un mare
o proprietate	337,08	unde au construit o casă unde continuă să își primească prietenii în plus au ctitorit acolo o frumoasă bisericuță
înmormântat	341	în cimititul din imediata apropiere decizia de a-și găsi odihna la malul mării nu poate fi desprinsă cred de simbolistica exilului atât de apropiată fostului mare exilat de altfel paul miron a și scris o piesă de teatru despre ovidiu alt mare exilat dar pe țărnițele pontului eucin

După analiza acestor probleme, s-a încercat o actualizare a modelelor de limba folosind transcrierile aproximative ale celor 2 seturi de date ca text adițional, obținând rezultatele din Tabelul 2.3.h.

**Tabelul 2.3.h. Transcrieri RAV obținute cu sistem cu model de limbă actualizat**

Set de date	Nr. cuvinte alinate	Nr. cuvinte transcrieri aproximative	Nr. cuvinte transcrieri RAV	Aliniere [% cuvinte]
alma24	4.137	5.151	5.621	80,31%
interviu2	33	18.095	8.830	0,18%

Comparând rezultatele din Tabelul 2.3.h cu cele inițiale din Tabelul 2.3.e, se poate observa o creștere a numărului de cuvinte alinate pentru setul alma24, dar o scădere drastică pentru setul interviu2. După realizarea unor seturi de experimente, în care s-au aliniat anumite părți de text din transcrierea aproximativă cu partea corespunzătoare a transcrierii RAV a setului interviu2 (prima jumătate a transcrierii aproximative cu prima jumătate a transcrierii RAV, a doua jumătate, primele 1000 de cuvinte, primele 2000 de cuvinte, ... , primele 8000 de cuvinte) s-a constatat că modulul de aliniere din aplicația JavaNLP2 (NISTAlign, prezentat în raportul anterior, secțiunea 2.2.3) întâlnește probleme când primește la intrare cantități mari de text. Ca soluție pentru această problemă, partea de aliniere a fost mutată în exteriorul aplicației JavaNLP2 și executată cu toolkit-ul sclite, urmând însă ca operațiile de filtrare să fie executate tot de javaNLP2. Rezultatele actualizate sunt prezentate în Tabelul 2.3.i.

Analizând rezultatele din Tabelul 2.3.i, putem observa o ușoară creștere a numărului de cuvinte alinate în cazul setului alma24 comparativ cu rezultatele din Tabelul 2.3.h. De asemenea, pentru setul interviu2, problema generată de către aliniatorul NISTAlign din toolkit-ul CMU Sphinx a fost rezolvată folosind modulul alternativ de aliniere din toolkit-ul sclite. Deși această problemă a fost rezolvată, procentul de cuvinte alinate rămâne relativ mic.

**Tabelul 2.3.i. Transcrieri RAV obținute cu sistem cu model de limbă actualizat + aliniere cu sclite**

Set de date	Nr. cuvinte alinate	Nr. cuvinte transcrieri aproximative	Nr. cuvinte transcrieri RAV	Aliniere [% cuvinte] din transcrierea aproximativă	Aliniere [% cuvinte] din transcrierea RAV
alma24	4.157	5.151	5.621	80,70%	73,95%
interviu2	3.811	18.095	8.830	21,06%	43,16%

### **Aplicare metodă de bază set de date CoBiLiRo-raw**

În urma testării metodei de bază pe cele 2 fișiere de test, metoda de bază a fost aplicată pe întreg setul de date CoBiLiRo-raw, cu mențiunea că procesul de aliniere a fost efectuat de toolkit-ul sclite și nu de CMU Sphinx (NISTAlign) pentru a evita eventualele probleme similare cu cele întâlnite la setul de date interviu2, menționate în subcapitolul anterior. Sistemul RAV folosit pentru transcrierea întregului set de date CoBiLiRo-raw a fost sistemul RAV baseline [Georgescu, 2019b]. Sistemul RAV cu model de limbă actualizat putea conduce, în final, la obținerea unui set de date adnotat mai mare, dar fără garanția că acel text ar fi fost era corect. Din acest motiv sistemul RAV cu model de limbă actualizat nu a fost folosit. În urma aplicării metodei de bază s-a obținut setul de date CoBiLiRo-trans-v4 ce conține secvențele de cuvinte aliniate pentru fiecare subset de date din CoBiLiRo-raw, împreună cu etichete de timp pentru fiecare cuvânt din secvență ce indică poziția și durata lor în fișierul audio corespunzător. Detalii legate de setul de date CoBiLiRo-trans-v4 se găsesc în Tabelul 2.1.e.

Sistemul RAV reantrenat cu seturile de date inițiale și setul de date nou adnotat (CoBiLiRo-trans-v4) a obținut o rată de eroare la nivel de cuvânt de 1.8% pe setul RSC-eval ce conține vorbire citită, respectiv 14.0% pe setul SSC-eval1 ce conține vorbire spontană. Astfel, adăugarea setului CoBiLiRo-trans-v4 de 31 ore la setul inițial de antrenare de 225 ore, a condus la o îmbunătățire relativă față de sistemul inițial cu 5.2% pe vorbire citită, respectiv 6.6% pe vorbire spontană.

### **2.3.3 Aplicarea metodei pentru setul de date CDep-raw**

Metoda de bază prezentată în secțiunea 2.3.1 a fost aplicată și pe setul de date CDep-raw (vezi Tabelul 2.1.d). Fișierele audio din setul de date CDep-raw au fost segmentate la 1 minut, obținând setul de date CDep-raw-1min. Setul de fișiere audio CDep-raw-1min a fost transcris utilizând sistemul RAV care a obținut cele mai bune rezultate în analiza făcută în 2019 [Georgescu, 2021]. A rezultat astfel setul de transcrieri RAV necesare pentru efectuarea procesului de aliniere din metoda de bază. După transcrierea fișierelor audio de 1 minut, 300 dintre acestea au fost separate, împreună cu transcrierile lor, în vederea formării unui set de date de evaluare CDep-eval. Transcrierile RAV ale celor 300 de fișiere audio au fost apoi aliniate cu transcrierile aproximative. Deoarece secvențele aliniate nu conțin întreg textul vorbit din fișierul audio a fost necesară intervenția manuală prin analiza și corectarea fișierelor text în vederea obținerii unor transcrieri de referință ce pot fi folosite mai departe pentru evaluarea viitoarelor sisteme RAV. Setul de transcrieri obținute, exceptând cele din setul CDep-eval, au fost alipite astfel încât timpii să se refere la fișierul audio original și nu la fișierul audio tăiat. După efectuarea procesului de aliniere între transcrierile RAV obținute după alipire și transcrierile aproximative CDep s-a obținut setul de date CDep-trans-v4, ce conține 21M de cuvinte și 878.8 ore de vorbire (84,2% din cuvinte, respectiv 25% din numărul de ore). Analizând rezultatele se poate observa o posibilă neconcordanță între numărul mare de cuvinte aliniate și numărul mic de ore, ceea ce a dus spre efectuarea unor analize suplimentare. Aceste analize au constatat în evaluarea fișierelor în funcție de energie și putere. Prin efectuarea acestor analize, ne-am așteptat să întâlnim un număr mare de fișiere cu un nivel de putere sub un anumit prag. Ne-am fi așteptat să găsim fișiere ce conțin liniște pe întreaga sau pe majoritatea duratei fișierului audio, ceea ce ar fi justificat numărul mic de ore raportat la numărul mare de cuvinte aliniate. După efectuarea evaluărilor și a unei analize manuale a mai multor fișiere audio, am observat că există și fișiere cu nivel relativ mic de putere ce conțin în mare parte vorbire decât liniște, însă nu am putut obține încă niște rezultate cantitative care să justifice neconcordanța menționată mai sus. Această problemă urmează să mai fie analizată în cadrul etapei 2021 a proiectului.

Sistemul RAV care a fost utilizat pentru generarea transcrierilor setului CDep-raw a fost reantrenat cu seturile de date inițiale și noul set adnotat CDep-trans-v4. Sistemul RAV reantrenat a fost evaluat pe seturile de evaluare RSC-eval și SSC-eval1 și SSC-eval2, iar rezultatele sunt prezentate în Tabelul 2.3.j alături de cele ale sistemului RAV inițial.

*Tabelul 2.3.j Sistem RAV baseline vs sistem RAV*

Sistem RAV	Seturi de antrenare	WER [%]		
		RSC-eval	SSC-eval1	SSC-eval2
Sistem RAV inițial [Georgescu, 2021]	RSC-train + SSC-train1+2 + SSC-train3+4-trans-v4	1.8	11.0	14.0
Sistem RAV reantrenat	RSC-train + SSC-train1+2 + SSC-train3+4-trans-v4 + CDep-trans-v4	1.7	12.3	15.4

Analizând rezultatele din tabelul 2.3.j, putem observa o ușoară scădere a WER de la 1.8% la 1.7% în cazul vorbirii citite (RSC-eval) și o degradare de aproximativ 10% a performanței pentru vorbire spontană (seturile SSC-eval1 și SSC-eval2). Din aceste rezultate putem trage concluzia că, în ciuda adăugării setului de date CDep-trans-v4 (879h) la antrenarea modelului acustic, noul sistem RAV prezintă o performanță mai scăzută per total față de sistemul inițial. O posibilă explicație pentru această performanță mai scăzută este faptul că prelegerile din ședințele camerei deputaților, ce se regăsesc în setul CDep-trans-v4, conțin în mare parte vorbire citită (similară cu cea din RSC și diferită față de SSC) și cuvinte dintr-un domeniu diferit față de cel întâlnit în seturile de evaluare (știri, interviuri, radio).

## **2.4 Activitățile 3.11 și 3.13 - Îmbunătățirea soluției pentru generarea de scoruri de încredere pentru RAV și Analiza impactului utilizării scorurilor de încredere pentru filtrarea transcrierilor RAV în vederea reantrenării sistemelor RAV**

În această secțiune descriem soluția îmbunătățită pentru estimarea scorurilor de încredere. Scopul nostru este de a înzestra sistemele automate cu estimări cât mai fiabile a încrederii în predicțiile pe care acestea le generează. Ideal, cu cât sistemul este mai încrezător într-o predicție (adică un scor de încredere mai mare), cu atât ne așteptăm ca rezultatul prezis să fie corect. Alternativ, putem rezolva problema complementară a estimării incertitudinii — în acest caz, cu cât este mai mare scorul de incertitudine, cu atât este mai probabil ca predicția să fie eronată.

În contextul proiectului de față, și anume, al recunoașterii automate a vorbirii (RAV) și al adnotării automate a datelor audio, estimarea încrederii este de o importanță crucială pentru că ne semnalează corectitudinea datelor transcrise automat și ne permite o filtrare pe baza acestui scor. Dar metodele pentru estimarea scorurilor de încredere pentru RAV au aplicații multiple, care merg dincolo de sarcina noastră de interes: îmbunătățirea robusteții sistemelor în sarcini critice de siguranță, evitarea erorilor în sistemele de dialog om-mașină sau facilitarea corecțiilor manuale în sarcinile de transcriere audio prin semnalarea erorilor. Mai mult, lucrări anterioare de specialitate au valorificat estimările scorurilor de încredere pentru o serie de sarcini care depind de RAV: selectarea predicțiilor cu grad ridicat de încredere pentru reantrenarea sistemului de bază [Sperber, 2017], propagarea incertitudinilor în traducerea automată a vorbirii [Vesely, 2013], adnotarea manuală a predicțiilor mai puțin sigure pentru învățarea activă [Yu, 2010].

Metoda pe care o propunem vizează estimarea încrederii pentru sistemele RAV de la un capăt la altul (en., end-to-end) [Hadian, 2018] — spre deosebire de soluția de bază care consideră sistemele RAV de tip hibrid, DNN-HMM. Modelele RAV de tip end-to-end au câștigat masiv în popularitate în ultima perioadă nu doar datorită performanței lor (care o egalează și chiar depășește pe cea a RAV-urilor de tip clasic), dar și pentru avantajele suplimentare de a fi simple din punct de vedere conceptual și de a permite un proces de antrenare unitar [Lüscher, 2019; Tüske, 2019; Karita, 2019]. Cu toate acestea, surprinzător de puține lucrări tratează sarcina de estimare a încrederii în sistemele RAV de la un capăt la altul.

Metoda noastră adresează două provocări principale ale dezvoltării metodelor de estimare a încrederii pentru sistemele RAV: (i) ieșirea structurată (textul) și (ii) predicțiile granulare (la nivel de grafem sau token în loc de cuvinte). Sistemele RAV sunt modele structurate (care mapează secvențe la secvențe) spre deosebire de rețelele obișnuite de recunoaștere (cum ar fi clasificarea imaginilor) a căror ieșire este o singură etichetă. Natura secvențială a ieșirii impune o etapă de decodare, ceea ce complică nu numai predicția, cât și algoritmul de estimare a încrederii, deoarece acesta trebuie să opereze într-un context auto-regresiv (folosind secvența deja prezisă). Din acest motiv, obținem predicțiile pe baza unui RAV preantrenat și apoi aplicăm metodele de estimare a încrederii peste probabilitățile următorului simbol, care sunt condiționate de transcrierea fixă.

Pentru a evita constrângerile impuse de utilizarea unui dicționar de cuvinte fix și pentru a permite predicții de cuvinte noi, sistemele RAV de tip end-to-end de obicei folosesc la ieșire token-uri de subcuvinte. Dar având în vedere că token-urilor le lipsește semantica, pentru multe aplicații finale suntem interesați să estimăm încrederea la nivelul cuvintelor și nu a token-urilor. În acest scop, explorăm metode de agregare a măsurilor de incertitudine de la nivel de token la unități mai mari, corespunzând cuvintelor. Dar tehnicile prezentate în continuare pot fi aplicate și la structuri chiar mai ample, cum ar fi nivelul propoziției sau al enunțului.

### 2.4.1 Legături cu starea artei

În această subsecțiune facem o scurtă prezentare a lucrărilor de specialitate relevante pentru sarcina și metoda propusă. Considerăm două direcții: metode pentru scoruri de încredere pentru recunoașterea automată a vorbirii (RAV) și metode pentru scoruri de încredere pentru sisteme de tip end-to-end.

**Scoruri de încredere pentru recunoașterea automată vorbirii.** Cele mai multe lucrări anterioare privind scorurile de încredere pentru RAV vizează sisteme clasice, bazate pe paradigma HMM-GMM. Aceste metode extrag mai întâi un set de caracteristici din rețeaua de decodare, modelul acustic sau de limbă, și apoi antrenează un clasificator pentru a prezice dacă transcrierea este corectă sau nu. Exemple tipice de caracteristici includ probabilitatea realizării acustice, scorul modelului de limbă, durata cuvântului, numărul de alternative din rețeaua de confuzie [Kemp, 1997; Weintraub, 1997; Hazen, 2002]. Mai recent, Swarup *et al.* [Swarup, 2019] a mărit setul de caracteristici folosind embedding-uri profunde ale semnalului acustic de intrare și ale textului prezis, în timp ce Errattahi *et al.* [Errattahi, 2018] a arătat că adaptarea la domeniu a caracteristicilor extrase aduce beneficii de performanță. Clasificatorii folosiți de metodele de estimare a scorurilor de încredere variază de la conditional random fields [Seigel, 2013; Cortina, 2016] și perceptron cu mai multe straturi [Kalgaonkar, 2015] la rețele neuronale recurente bidirecționale [Ogawa, 2017; Del-Agua, 2018; Li, 2019].

**Scoruri de încredere în sistemele end-to-end.** Metoda de bază pentru estimarea încrederii în rețelele neuronale este de a utiliza în mod direct probabilitatea predicției celei mai probabile [Hendrycks, 2016]. S-a observat că rețelele neuronale tind să fie prea sigure și estimările probabilității pot fi îmbunătățite prin scalarea temperaturii [Hinton, 2015], ceea ce duce de obicei la o mai bună calibrare [Guo, 2017; Ashukha, 2020]. Cea mai promițătoare direcție în ceea ce privește simplitatea și utilitatea implică estimarea Monte Carlo: Gal și Ghahramani [Gal, 2016] folosesc tehnica dropout la momentul inferenței pentru a obține predicții multiple, care sunt apoi mediate, în timp ce Lakshminarayanan *et al.* [Lakshminarayanan, 2017] fac media predicțiilor unui ansamblu de rețele de obicei antrenate cu inițializări diferite. Această ultimă metodă s-a dovedit a fi foarte fiabilă pentru condițiile dificile în care avem date din afara domeniului [Ovadia, 2019], dar este costisitoare din punct de vedere computațional, deoarece implică antrenarea a multiple modele de RAV [Ashukha, 2020]. O abordare diferită a estimării încrederii este de a învăța un clasificator (de obicei o altă rețea neuronală) direct deasupra activărilor rețelei de RAV [Corbière, 2019; Chen, 2019].

La intersecția acestor două direcții de cercetare, menționăm lucrarea foarte recentă a lui Malinin și Gales [Malinin, 2020], care similar cu noi abordează sarcina de estimare a încrederii pentru sistemele RAV end-to-end. Cu toate acestea, ei sunt preocupați de estimarea incertitudinii la nivel de token și frază, în timp ce noi suntem interesați de estimarea la nivel de cuvânt, și, în consecință, oferim mai multă atenție asupra tehnicilor de agregare. Mai mult, ei folosesc ansambluri ca metodă principală de estimare a încrederii, în timp ce noi evaluăm și metodele de reducere a temperaturii și de scădere. Deși tehnica de dropout a fost folosită anterior pentru obținerea scorurilor de încredere pentru RAV [Vyas, 2019], metoda este diferită de abordarea noastră. În acea lucrare, autorii generează mai multe ipoteze prin dropout și apoi atribuie confidențe cuvintelor pe baza frecvenței aparițiilor lor în ipotezele aliniate. În schimb, noi agregăm probabilitățile posterioare și nu ipotezele, ceea ce simplifică procedura, deoarece evită pasul de aliniere.

### 2.4.2 Metodologia

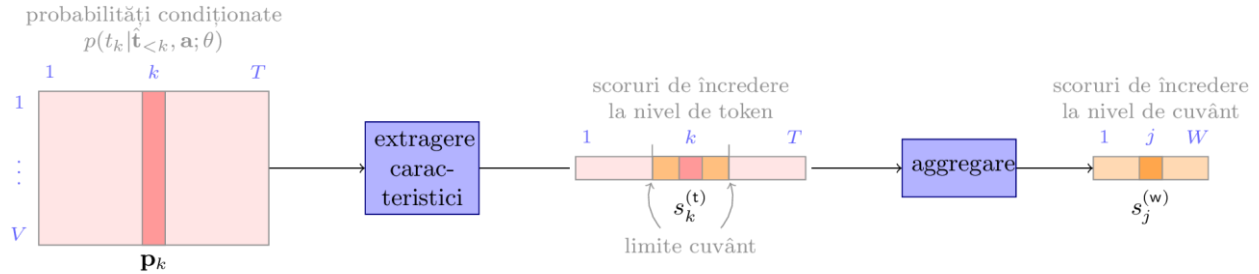
În această subsecțiune prezentăm efectiv metodologia propusă pentru estimarea scorurilor de încredere și modalități propuse de îmbunătățire a acestora. Începem mai întâi cu o descriere generală a metodei și a notației implicate.

Considerăm un model de tip secvență-la-secvență (en. sequence-to-sequence) care mapează o secvență audio  $\mathbf{a}$  la una de token-uri  $\mathbf{t} = (t_1, \dots, t_T)$ . Modelul este specificat de parametri  $\theta$ , care sunt învățați prin minimizarea pe setul de antrenare a unei funcții de pierdere, cum ar fi funcția *connectionist temporal classification* (CTC) sau divergența Kullback-Leibler. La inferență, modelul generează probabilități pentru următorul simbol  $k$  într-o manieră autoregresivă  $p(t_k | \hat{t}_{<k}, \mathbf{a}; \theta)$ , bazat pe token-urile deja prezise  $\hat{t}_{<k}$ . Aceste probabilități sunt utilizate pentru efectuarea decodării prin metoda beam search pentru a obține cea mai probabilă secvență de token-uri. Având în vedere că probabilitatea de ieșire condiționată este o distribuție peste  $V$  token-uri din vocabular, o notăm cu un vector  $V$ -dimensional  $p_k$ .



## Estimarea scorurilor de încredere

Scopul nostru este de a obține un scor de încredere pentru fiecare cuvânt din transcrierea produsă de RAV. Realizăm acest lucru în doi pași. Mai întâi, folosind probabilitățile aposteriori de la fiecare moment de timp  $p_k$ , extragem caracteristici pentru a reprezenta scorul de încredere al fiecărui token  $s_k^{(t)}$ . În cea de-a doua etapă, agregăm scorurile la nivel de token în scoruri de încredere la nivel de cuvânt  $s_j^{(w)}$ , pe baza token-urilor care aparțin fiecărui cuvânt. În continuare vom detalia acești doi pași; vezi și figura 2.4.a.



**Figura 2.4.a** Prezentare schematică generală a metodei propuse de estimare a scorurilor de încredere. Pe baza unui sistem de recunoaștere a vorbirii (RAV) de tip end-to-end, obținem probabilități  $p_k$  pentru fiecare token  $k$  condiționate de o rostire  $a$  și token-urile prezise anterior  $\hat{t}_{<k}$ . Pe baza acestor probabilități extragem scoruri de încredere la nivel de token  $s_k^{(t)}$ , pe care le agregăm apoi pentru a obține scoruri la nivel de cuvânt  $s_j^{(w)}$ . Dimensiunea vocabularului de token-uri este notată cu  $V$ , numărul de token-uri este notat cu  $T$  și numărul de cuvinte cu  $W$ .

**Extragerea caracteristicilor.** Pentru a măsura încrederea într-o predicție la nivel de token folosim două variante:

- Logaritmul probabilității (log-proba) celei mai probabile predicții date de clasificator, adică  $s_k^{(t)} = \log \max p$ . S-a observat empiric că acest tip de caracteristică oferă un baseline puternic pentru sarcinile de clasificare a greșelilor și detectare a eșantioanelor din afara distribuției [Hendrycks, 2016].
- Entropie negativă (neg-entropie) calculată peste vocabularul token-urilor la fiecare moment de timp, adică  $s_k^{(t)} = p^T \log p$ . O entropie mare înseamnă o incertitudine mare sau, invers, o entropie negativă mare implică o predicție încrezătoare. Deși entropia este de obicei utilizată ca măsură de încredere împreună cu tehnica de dropout [Gal, 2016], după cum vom vedea în curând, tehnica de dropout va putea fi folosită și peste probabilitățile originale.

**Agregarea.** Pentru a obține caracteristici la nivel de cuvânt din cele la nivel de simbol, experimentăm cu trei tipuri de funcții de agregare: suma, media, minimumul. Deoarece ambele caracteristici propuse sunt negative, însumarea a mai multe token-uri va duce la valori mai mici și, prin urmare, la scoruri de încredere mai mici; acest comportament poate fi de dorit deoarece cuvintele mai lungi sunt mai susceptibile de a fi eronate. De asemenea, atunci când însumăm logaritmul probabilităților, obținem un scor la nivel de cuvânt corespunzător probabilității log a întregii secvențe. Utilizarea agregării cu funcția de minim este justificată de faptul că am putea dori o încredere scăzută dacă cel puțin unul din token-uri are o încredere scăzută.

## Îmbunătățirea probabilităților la nivel de token

Propunem trei moduri de a face mai fiabile probabilitățile la nivel de token: scalarea temperaturii, tehnica de dropout și ansamblurile de modele. Presupunerea noastră este că îmbunătățind probabilitățile la nivel de token, vom obține și scoruri de confidență la nivel de cuvânt mai bune.

**Scalarea temperaturii** [Hinton, 2015; Guo, 2017] constă în împărțirea activărilor de tip logit (valorile de dinainte de stratul softmax) la un scalar  $\tau$  (cunoscut sub numele de temperatură). Valoarea lui  $\tau$  variază de la zero la infinit și controlează forma distribuției: când  $\tau \rightarrow 0$  obținem o distribuție uniformă, când  $\tau \rightarrow \infty$  obținem o distribuție Dirac localizată pe cea mai probabilă ieșire. Pe baza temperaturii  $\tau$  actualizăm probabilitățile la nivel de token la fiecare moment de timp, după cum urmează:

$$p'_k = \text{softmax}(\log(p_k) / \tau)$$

Apoi extragem caracteristici  $s_k^{(t)}$  peste probabilitățile actualizate, le agregăm în scorul la nivel de cuvânt  $s_j^{(w)}$  și, în cele din urmă, clasificăm cuvântul drept corect sau incorect:

$$P(\text{correct}) = \sigma(\alpha s^{(w)} + \beta)$$

Variabilele  $\alpha, \beta, \tau$  sunt parametri și sunt învățate prin optimizarea unei funcții de pierdere de entropie încrucișată (en., *cross-entropy loss*) pe un set de validare. Etichetele sunt setate la nivel de cuvânt prin alinierea la textul de bază cu transcrierea. Remarcăm faptul că parametrii  $\alpha$  și  $\beta$  nu modifică ordinea predicțiilor, dar ne permit să învățăm un model de estimare a încrederii calibrat.

**Dropout** [Srivastava, 2014] este o tehnică care maschează părți aleatorii ale activărilor într-o rețea, făcând rețeaua mai puțin predispusă la supra-antrenare (en., *overfitting*). În [Gal, 2016] s-a observat că tehnica de dropout induce o distribuție de probabilitate peste ponderile rețelei și poate fi în consecință utilizată pentru inferența Bayesiană aproximativă. Folosim această idee și calculăm probabilitățile la nivel de token obținute prin mai multe rulări folosind dropout:

$$p'_k = \frac{1}{N} \sum_n \hat{p}_k$$

unde  $\hat{p}$  specifică predicția dropout-ului. Probabilitățile actualizate sunt apoi utilizate pentru a extrage oricare dintre caracteristicile de incertitudine propuse (log-proba sau neg-entropie).

**Ansamblurile** [Lakshminarayanan, 2017] se bazează pe aceeași idee de mediere a predicțiilor din mai multe surse (ca în cazul dropout-ului), dar în acest caz ponderile provin din rețele antrenate independent (folosind inițializări diferite ale rețelei). În cazul nostru, calculăm media predicțiilor la nivel de token:

$$p'_k = \frac{1}{N} \sum_n p(t_k | \hat{t}_{<k}, a; \theta_n)$$

unde  $\{\theta_n\}_{n=1..N}$  specifică ansamblul de modele. Este important de subliniat că trebuie să avem același context pentru toate modelele din ansamblu, deci îl folosim pe cel dat de un model preantrenat.

Cele trei abordări prezentate pot fi combinate; de exemplu, mai întâi putem actualiza probabilitățile folosind scalarea temperaturii și apoi media mai multe predicții folosind tehnica de dropout.

### 2.4.3 Setup-ul experimental pentru experimente pe limba engleză

În această subsecțiune descriem bazele de date, sistemele de recunoaștere automată a vorbirii utilizate și metricile de evaluare utilizate pentru evaluarea scorurilor de încredere pe limba engleză

#### **Baze de date**

Pentru experimentele pe limba engleză am ales mai multe seturi de date, toate disponibile public și utilizate de comunitate mai ales pentru evaluarea sistemelor de recunoaștere a vorbirii.

**LibriSpeech** [Panayotov, 2015] este un corpus de aproximativ 1000 de ore de cărți audio citite. Datele au fost derivate din proiectul LibriVox și au fost atent segmentate și alinate. Folosim setul de date atât pentru antrenare, cât și pentru evaluare. Pentru antrenare folosim cele trei părți ale acestuia `clean100`, `clean360` și `other500`, în timp ce pentru validare și evaluare folosim părțile standard denumite `clean`, respectiv `other`.

**TED-LIUM2** [Rousseau, 2014] constă în discursuri și transcrierile acestora colectate de pe site-ul web TED. Utilizăm setul de date pentru evaluare și, în consecință, folosim numai subseturile predefinite `dev` și `test`.

**CommonVoice** [Ardila, 2020] este un set de date colaborativ cu transcrieri scurte care sunt citite de oameni din întreaga lume. Există mai multe versiuni ale setului de date și am folosit prima versiune. Folosim setul de date pentru evaluare și am definit subseturi de `dev` și `test` prin alegerea a 10% din rostiri în mod aleatoriu pentru fiecare dintre ele.

Tabelul 2.4.a prezintă dimensiunile părților de `test` pentru fiecare set de date de evaluare.

**Tabelul 2.4.a** Dimensiunea seturilor de date (partea *test*) care sunt folosite pentru evaluarea metodelor pentru scorurile de încredere pe limba engleză.

Bază de date	Num. rostiri	Durăță (ore)
Libri clean	2.6K	5.4
Libri other	2.9K	5.3
TED	1.1K	2.6
CommonVoice	66K	72

### **Sistemul de recunoaștere automată a vorbirii**

Sistemul de recunoaștere automată a vorbirii (RAV) este implementat folosind librăria ESPnet [Watanabe, 2018]. Modelul implementează arhitectura de tip Transformer [Vaswani, 2017] și primește la intrare un banc de filtre Mel 80-dimensionale (extrase cu utilitarul Kaldi [Povey, 2011]) și produce la ieșire o secvență de token-uri. Vocabularul are dimensiunea de 5,000 de token-uri și este obținut prin segmentarea cuvintelor pe baza unui model de limbaj de tip unigram [Kudo, 2018]. Modelul este antrenat pe 960 de ore din baza de date LibriSpeech [Panayotov, 2015], iar datele sunt augmentate utilizând tehnicile SpecAugment (modificarea vitezei vorbirii, mascare în frecvență, mascare în timp) [Park, 2019]. Pentru decodare folosim un model de limbă, care este implementat tot ca Transformer și este antrenat pe 14,500 de cărți din domeniu public [Panayotov, 2015]. Vocabularul modelului de limbă constă din aceleași 5,000 de token-uri ca și modelul RAV.

Pentru experimentele cu ansambluri de modele, reantrenăm sistemul de RAV folosind aceeași arhitectură și date, dar inițializări ale ponderilor diferite. Repetăm procesul de patru ori obținând patru modele independente. Datorită constrângerilor de calcul, aceste modele au fost antrenate pentru un număr mai scurt de epoci decât sistemul principal (10 versus 120), dar am observat că curba funcției de pierdere a validării a început să se aplatizeze și că performanța sistemului de RAV este rezonabilă ( $5.5\% \pm 0.4$  WER pe Libri clean față de  $2.7\%$  obținut de modelul preantrenat).

### **Metrici de evaluare**

În mod ideal, ne dorim ca scorurile de încredere să fie corelate cu corectitudinea transcrierii, adică cuvintele corecte ar trebui să aibă un scor de încredere mare, în timp ce cuvintele incorecte, un scor scăzut. Pe baza lucrărilor anterioare [Hendrycks, 2016; Corbière, 2019; Malinin, 2020], folosim metrici care sunt utilizate în general pentru evaluarea clasificatorilor binari, dar variază pragul de discriminare între cele două clase. Mai precis, măsurăm aria de sub curba de precision-recall (area under precision-recall curve; AUPR) și aria de sub curba receiver operator curve (area under receiver operator curve; AUROC). Cu toate acestea, în funcție de ceea ce dorim să ne concentrăm (erori sau predicții corecte) obținem două variante: dacă suntem interesați de clasificări greșite, vom trata erorile ca pe o clasă pozitivă; pe de altă parte, dacă suntem interesați de clasificarea corectă, vom trata detecțiile reușite ca fiind clasa pozitivă. Prin urmare, pentru AUPR folosim două variante AUPR<sub>e</sub> (când erorile sunt tratate ca pozitive) și AUPR<sub>s</sub> (când succesele sunt tratate ca pozitive). Pentru AUROC se obține aceeași valoare pentru ambele opțiuni, deci nu este necesar să se facă această distincție.

#### **2.4.4 Rezultate experimentale pentru limba engleză**

Această secțiune prezintă rezultatele experimentale pe bazele de date în limba engleză. Începem cu o evaluare a caracteristicilor și a agregărilor acestora; după care raportăm rezultatele pentru variantele îmbunătățite: mai întâi, pentru tehnicile de scalare a temperaturii și dropout, apoi pentru ansamblurile de modele. Încheiem secțiunea cu o discuție privind alte legături cu starea artei și posibile direcții.

### **Caracteristici și metode de agregare**

Evaluăm caracteristicile pentru reprezentarea incertitudinii din token-uri și tehnicile de agregare propuse pe cele patru seturi de date descrise în secțiunea precedentă. Folosim modelul preantrenat pentru a obține predicții de text pentru toate fișierele audio din partea de testare al fiecărui set de date, și apoi estimăm

încrederea pe baza metodologiei propuse. Tabelul 2.4.b prezintă rezultatele pentru toate combinațiile de caracteristici și agregări.

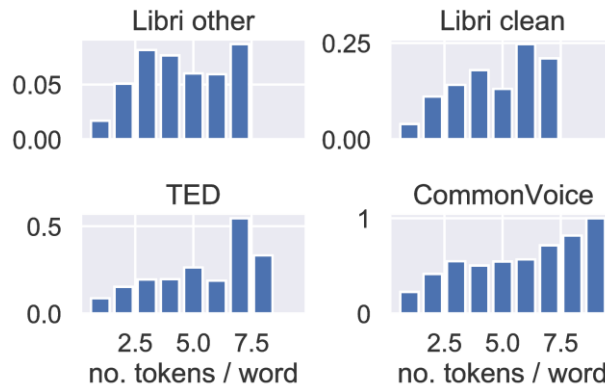
*Comparație a caracteristicilor.* Observăm că prin folosirea caracteristicile de tip probabilități log se obțin rezultate mai bune față de situația utilizării caracteristicilor pe bază de entropie în toate cazurile (indiferent de agregare sau setul de date). Singura excepție notabilă este setul de date CommonVoice în care rezultatele sunt comparabile.

*Comparație a metodelor de agregare.* În general, agregarea folosind suma funcționează mai bine cu caracteristicile de tip log-proba, în timp ce agregarea folosind minimul funcționează mai bine pentru caracteristicile de entropie. Suma s-ar putea să nu fie potrivită pentru caracteristicile de entropie, deoarece magnitudinea lor este mai mare decât în cazul caracteristicilor de tip log-proba, iar cuvintele lungi sunt penalizate prea mult prin lungime; dar, așa cum vom vedea mai departe, acest comportament poate fi atenuat prin metoda de scalare a temperaturii. Media este, în general, slabă pentru ambele caracteristici, sugerând că măsurile invariante în lungime sunt dăunătoare. Într-adevăr, o privire mai atentă la frecvența erorilor cu dimensiunea lungimii indică următorul fapt: cu cât un cuvânt este format din mai multe token-uri, cu atât este mai probabil că este incorect, vezi figura 2.4.b.

*Comparație între seturile de date.* Așa cum era de așteptat, modelul preantrenat are cea mai bună performanță în ceea ce privește datele din domeniu (2.7% WER pe Libri clean și 6.0% pe Libri other), performanța scăzând apoi brusc, pe măsură ce evaluăm datele din afara domeniului (13.3% pe TED și 28.6% pe CommonVoice). În fiecare dintre aceste situații, numărul de cuvinte care sunt clasificate corect se schimbă, trecând de la mai multe pe baza de date Libri la mai puține pe seturile de date TED și CommonVoice. Această observație explică de ce performanța pentru AUPRs scade în funcție de domeniul datelor, și, dimpotrivă, de ce se îmbunătățește performanța AUPRe. Din păcate, din exact acest motiv—performanța diferită a sistemului RAV de bază pe cele patru seturi de date—este imposibil de comparat metodele de încredere între seturile de date, deoarece acestea utilizează o altă adnotare [Ashukha, 2020].

**Tabelul 2.4.b** Rezultate pentru evaluarea scorurilor de încredere pentru toate combinațiile de caracteristici (caract.) și metode de agregare (agreg.) pe cele patru baze de date în engleză. Pentru toate cele trei metrici de evaluare utilizate (AUPRe, AUPRs, AUROC) valorile mai mari indică performanță mai bună. Raportăm eroarea la nivel de cuvânt pentru sistemul de RAV preantrenat pe fiecare din seturile de date; aceasta este listată în dreapta numelului bazei de date.

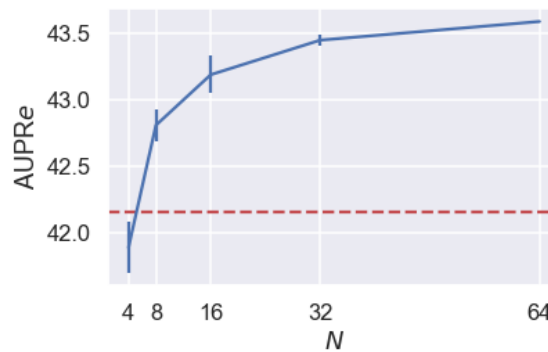
		Libri clean / 2.7%			Libri other / 6.0%		
caract.	agg.	AU-PR <sub>e</sub>	AU-PR <sub>s</sub>	AU-ROC	AU-PR <sub>e</sub>	AU-PR <sub>s</sub>	AU-ROC
log-proba	sum	21.55	99.21	82.41	29.99	98.10	81.75
log-proba	min	21.85	99.19	82.47	28.64	98.06	81.66
log-proba	avg	20.12	99.10	80.90	26.72	97.93	80.47
neg-entropy	sum	17.31	99.10	79.97	26.37	97.86	79.58
neg-entropy	min	19.94	99.09	80.55	26.75	97.82	79.64
neg-entropy	avg	17.55	98.95	77.72	24.26	97.59	77.46
		TED / 13.3%			CommonVoice / 28.6%		
caract.	agg.	AU-PR <sub>e</sub>	AU-PR <sub>s</sub>	AU-ROC	AU-PR <sub>e</sub>	AU-PR <sub>s</sub>	AU-ROC
log-proba	sum	39.97	95.88	79.95	48.98	77.71	64.84
log-proba	min	39.74	95.94	80.58	46.79	76.74	62.67
log-proba	avg	38.74	95.88	80.29	44.51	75.82	60.87
neg-entropy	sum	34.96	95.41	77.57	47.71	77.10	63.74
neg-entropy	min	37.55	95.56	79.01	45.51	76.00	61.21
neg-entropy	avg	36.28	95.42	78.29	42.64	74.83	58.75



**Figura 2.4.b** Raportul de erori ca o funcție de lungimea cuvântului. Raportul de erori este calculat ca numărul de cuvinte eronate împărțit la numărul total de cuvinte, iar lungimea cuvintelor este măsurată ca număr de token-uri.

### Scalarea temperaturii și tehnica de dropout

Evaluăm metodele de estimare a scorurilor de încredere după îmbunătățirea probabilităților token-urilor prin două din tehnicile descrise: scalarea temperaturii și tehnica de dropout. Folosim sistemul RAV preantrenat și raportăm rezultatele pe setul de testare TED. Parametrii pentru metoda de scalare a temperaturii sunt învățați pe partea dev a setului de date TED pentru fiecare setare de caracteristică și metodă de agregare. Când scalarea temperaturii este combinată cu tehnica de dropout, aplicăm mai întâi scalarea temperaturii (folosind aceeași temperatură) și apoi mediem probabilitățile obținute prin dropout. Pentru dropout mediem 64 de predicții independente. Tabelul 2.4.c prezintă rezultatele pentru toate combinațiile de caracteristici, agregări și tehnici de îmbunătățire. Rezultatele indică faptul că ambele metode propuse îmbunătățesc rezultatele la fel ca și combinația lor, ceea ce oferă în general cel mai bun rezultat. Observăm că tehnica de dropout aduce îmbunătățiri mai mari pentru caracteristicile log-proba, în timp ce caracteristica neg-entropie produce rezultate mai bune atunci când se utilizează scalarea temperaturii. Interesant, cele mai bune rezultate sunt acum obținute pentru neg-entropie cu agregare sumă (rândul 16). Figura 2.4.c arată că performanța dropout-ului se îmbunătățește cu numărul de rulări și se plafonează în jurul valorii alese de 64.



**Figura 2.4.c** Performanța AUPRe ca funcție de numărul de rulări de dropout  $N$  pe baza de date TED. Linia roșie orizontală indică rezultatul metodei fără dropout. Modelul folosește caracteristici de tip neg-entropie, agregarea pe bază de sumă și scalare a temperaturii.

### Ansambluri de modele

În cele ce urmează, prezentăm rezultate pentru estimarea scorurilor de încredere folosind ansambluri de modele și combinațiile acestora cu celelalte versiuni îmbunătățite (scalarea temperaturii și tehnica de dropout). Pentru fiecare dintre modelele reantrenate, care fac parte din ansamblu, folosim predicțiile modelului preantrenat pentru a selecta transcrierea pentru care dorim să estimăm scoruri de încredere; deci modelul reantrenat este folosit doar pentru scorul de încredere, prin extragerea caracteristicilor de încredere descrise anterior. Rezultatele sunt prezentate în tabelul 2.4.d. Pentru rândurile care nu folosesc ansamblu (rândurile 1, 2, 3 și 5 din tabel) evaluăm fiecare dintre cele patru modele individuale în mod independent și raportăm performanța medie. Modelul preantrenat (tabelul 2.4.c, rândul 13) are, în general, o performanță mai bună decât cele reantrenate (tabelul 2.4.d, rândul 1), sugerând că performanța predictivă a unui model se

**Tabelul 2.4.c.** Evaluarea scorurilor de încredere pe baza de date TED pentru combinații de caracteristici, metode de agregare și tehnici de îmbunătățire a probabilităților: scalarea temperaturii (ST) și tehnica de dropout (D). Punctul negru “•” indică dacă e folosită o anumită tehnică, iar linia “-” indică opusul.

nr.	caract.	agreg.	ST	D	AUPRe	AUPRs	AUROC
1	log-proba	sum	-	-	39.97	95.88	79.95
2			-	•	41.41	96.81	82.78
3			•	-	40.92	96.19	81.11
4			•	•	42.99	97.14	84.10
5	log-proba	min	-	-	39.74	95.94	80.58
6			-	•	42.08	96.94	83.76
7			•	-	39.84	95.98	80.74
8			•	•	42.17	97.00	83.93
9	log-proba	avg	-	-	38.74	95.88	80.29
10			-	•	41.19	96.95	83.73
11			•	-	38.97	95.99	80.66
12			•	•	41.32	97.06	84.08
13	neg-entropy	sum	-	-	34.96	95.41	77.57
14			-	•	33.14	96.22	79.45
15			•	-	42.16	96.91	83.50
16			•	•	43.59	97.62	85.51
17	neg-entropy	min	-	-	37.55	95.56	79.01
18			-	•	38.75	96.53	81.98
19			•	-	41.23	96.87	83.50
20			•	•	42.23	97.60	85.51
21	neg-entropy	avg	-	-	36.28	95.42	78.29
22			-	•	38.01	96.51	81.85
23			•	-	40.22	96.53	82.48
24			•	•	41.15	97.43	85.18

corelează cu performanța sa de estimare a încrederii. Printre cele trei metode de îmbunătățire propuse, observăm că scalarea temperaturii oferă cea mai mare creștere a performanței pentru toate cele trei valori (rândul 2). În mod surprinzător, metoda dropout-ului îmbunătățește numai performanța AUPR față de linia de bază (rândul 3). În combinațiile de două metode, scalarea temperaturii și ansamblul se completează reciproc și obțin performanțe mai bune.

**Tabelul 2.4.d** Evaluarea scorurilor de încredere pe baza de date TED pentru combinații ale metodelor îmbunătățite: scalarea temperaturii (ST), dropout (D) și ansambluri (A). Am folosit caracteristica neg-entropy și suma ca agregare. Punctul negru “•” indică dacă e folosită o anumită tehnică, iar linia “–” indică opusul.

nr.	ST	D	A	AUPRe	AUPRs	AUROC
1	–	–	–	28.58	95.30	75.79
2	•	–	–	32.00	96.32	79.47
3	–	•	–	27.49	95.51	75.67
4	–	–	•	30.89	96.26	78.89
5	•	•	–	31.10	96.40	79.06
6	•	–	•	34.57	96.95	81.64
7	–	•	•	28.94	96.26	77.93
8	•	•	•	33.00	96.84	80.82

În continuare, discutăm pe scurt diferite perspective legate de metodologia propusă.

**Mărirea setului de caracteristici.** Caracteristicile pe care le-am investigat au avantajul de a fi generale, pentru că valorifică probabilitățile la nivel de token, care sunt disponibile în majoritatea, dacă nu în toate, utilitățile pentru RAV-uri de tip end-to-end existente. Cu toate acestea, setul de caracteristici ar putea fi extins cu probabilități pentru sunetul de intrare sau textul generat, sau cu informații despre durata cuvântului, care pot fi extrase din ponderile de atenție.

**Învățarea caracteristicilor.** Inspirați de studiile precedente [Corbière, 2019; Chen, 2019], am experimentat, de asemenea, cu învățarea unei rețele de estimare a încrederii pe baza caracteristicilor extrase din modelul end-to-end (mai exact, activări logit și pre-logit). Cu toate acestea, experimentele noastre nu au reușit să arate îmbunătățiri față de rezultatele prezentate.

**Tratarea cuvintelor pierdute.** Pentru a genera etichete pentru scorurile de încredere, aliniem textul de referință la textul prezis și marcăm cuvintele corecte din textul prezis ca pozitive și substituțiile și inserțiile ca negative. Această abordare este tipică în literatura de evaluare a scorurilor de încredere, dar ratează erorile făcute prin ștergerea (pierderea) cuvintelor din referință. Mai multe lucrări au abordat această problemă [Seigel, 2014; Ragni, 2018]; lăsăm pentru viitor să extindem abordarea curentă pentru această sarcină.

#### 2.4.5 Setup-ul experimental pentru experimente pe limba română

Pentru experimentele pe limba română folosim seturile de date de evaluare prezentate în secțiunea 2.1.1: RSC-eval (5.5 ore de vorbire citită) și SSC-eval1+2 (3.5+1.5 ore de vorbire spontană).

Sistemul de transcriere de vorbire folosit pentru aceste experimente este cel introdus în secțiunea 2.2.1. Sistemul RAV este implementat folosind librăria ESPnet și se bazează pe o arhitectură de tip Transformer. Vocabularul sistemului cuprinde 1,000 de subcuvinte și a fost extras din transcrierile materialelor audio cu care a fost antrenat întregul sistem: cele 517 ore de vorbire din seturile de date RSC-train, SSC-train1+2 și SSC-train3+4-trans-v4. Pentru decodare se folosește un model de limbă adițional implementat tot ca Transformer și antrenat pe cele ~352M de cuvinte din seturile de date news2017 și talkshows. Performanța sistemului de RAV folosit, exprimată sub forma erorii la nivel de cuvânt (WER), este de 3.4% WER pe RSC-eval, 15.3% WER pe SSC-eval1, respectiv 23.5% WER pe SSC-eval2.

Din cauza complexității modelului CTC-Transformer și timpului mare de inferență pe configurația hardware disponibilă în acest proiect, nu s-au mai putut antrena modele adiționale cu inițializări diferite ale ponderilor, pentru ansamblurile de modele.

Metricile de evaluare introduse în secțiunea anterioară (*i.e.* AUROC, AUPRs, AUPRs) sunt universale, independente de limbă și, în consecință, vor fi folosite și pentru experimentele pe limba română.

## 2.4.6 Rezultate experimentale pentru limba română

Pentru limba română am reluat următoarele experimente realizate și pentru sistemul de estimare a încrederii pentru limba engleză:

- evaluarea diverselor caracteristici ce pot fi utilizate pentru estimarea încrederii (Tabelul 2.4.e);
- evaluarea diverselor metode de agregare a caracteristicilor la nivel de subcuvânt pentru a obține scoruri de încredere la nivel de cuvânt (Tabelul 2.4.e);
- evaluarea tehnicii dropout pentru îmbunătățirea scorurilor de încredere de mai sus Tabelul 2.4.f).

### Caracteristici și metode de agregare

În această secțiune evaluăm caracteristicile pentru reprezentarea incertitudinii din subcuvinte și tehnicile de agregare propuse pe cele trei seturi de date prezentate în secțiunea 2.1.1: RSC-eval (5.5 ore de vorbire citită) și SSC-eval1+2 (3.5+1.5 ore de vorbire spontană). Tabelul 2.4.e prezintă rezultatele pentru toate combinațiile de caracteristici și agregări.

**Tabelul 2.4.e** Rezultate pentru evaluarea scorurilor de încredere pentru toate combinațiile de caracteristici (caract.) și metode de agregare (agg.), pe cele trei seturi de evaluare în limba română. Pentru toate cele trei metrice de evaluare utilizate (AU-PR<sub>e</sub>, AU-PR<sub>s</sub>, AU-ROC) valorile mai mari indică performanță mai bună.

caract.	agg.	RSC-eval / 1.8%			SSC-eval1 / 11.0%			SSC-eval2 / 14.0%		
		AU-PR <sub>e</sub>	AU-PR <sub>s</sub>	AU-ROC	AU-PR <sub>e</sub>	AU-PR <sub>s</sub>	AU-ROC	AU-PR <sub>e</sub>	AU-PR <sub>s</sub>	AU-ROC
log-proba	sum	15.81	98.77	74.50	27.19	96.14	74.57	15.79	91.90	61.52
	min	18.30	98.91	78.66	28.32	96.34	77.24	13.86	91.99	60.46
	avg	16.22	98.74	75.75	26.69	96.04	75.38	14.04	92.09	60.60
neg-entropy	sum	15.82	98.73	75.65	26.44	96.03	75.31	14.30	92.17	61.09
	min	16.28	98.86	77.79	27.10	96.13	76.14	12.75	91.47	57.94
	avg	13.55	98.65	74.02	24.97	95.75	73.62	13.37	91.68	58.88

*Comparație a caracteristicilor.* În ceea ce privește tipul de caracteristici (log-proba vs. neg-entropy), rezultatele pe seturile de date în limba română indică exact aceleași concluzii ca și în cazul limbii engleză:

- pe seturile de date pe care sistemul RAV are o acuratețe relativ bună (în cazul nostru RSC-eval și SSC-eval1) caracteristicile de tip probabilități log obțin rezultate mai bune decât caracteristicile pe bază de entropie indiferent de agregare
- pe seturile de date mai dificile (SSC-eval2 pentru română, respectiv CommonVoice pentru engleză) caracteristicile bazate pe entropie agregate cu metoda *sum* prezintă rezultate similare cu, uneori chiar mai bune decât caracteristicile de tip probabilități log.

*Comparație a metodelor de agregare.* În ceea ce privește metodele de agregare, concluziile sunt relativ diferite pentru limba română:

- pe limba română agregarea minimul funcționează mai bine indiferent de caracteristici, în timp ce pe limba engleză caracteristicile de tip probabilități log erau agregate mai bine folosind suma;
- agregarea prin medie este, în general, slabă pentru ambele caracteristici, sugerând că măsurile invariante în lungime sunt dăunătoare;
- excepții de la concluziile de mai sus apar din nou pentru seturile de date mai dificile (SSC-eval2 pentru română, respectiv CommonVoice pentru engleză) pentru care agregarea prin sumă produce cele mai bune rezultate, indiferent de caracteristicile utilizate.

*Comparație între seturile de date.* Așa cum era de așteptat, modelul preantrenat are cea mai bună performanță în ceea ce privește datele din domeniu (1.8% WER pe RSC-eval și 11.0% pe SSC-eval1) și performanță mai slabă pe date din afara domeniului (14.0% pe SSC-eval2). În fiecare dintre aceste setări, numărul de cuvinte care sunt clasificate corect se schimbă, trecând de la mai multe pe setul RSC la mai



puține pe seturile de date SSC. Această observație explică de ce performanța pentru AUPRs este mare pe RSC și scade pe SSC-eval1 și 2, și, dimpotrivă, de ce performanța pentru AUPRe este mai slabă pentru RSC și mai mică pentru SSC-eval1. Din păcate, din exact acest motiv—performanța diferită a sistemului RAV de bază pe cele patru seturi de date—este imposibil de comparat metodele de încredere între seturile de date. Performanța pentru AUPRe pe setul de date SSC-eval2 face din nou notă discordantă: dat fiind faptul că sistemul RAV are performanțe slabe de transcriere pe acest set de date ne-am fi așteptat ca acele cuvinte transcrise greșit să poată fi foarte ușor de identificat (un AUPRe de valoare mare), însă acest lucru nu se întâmplă.

### ***Tehnica de dropout***

În această secțiune evaluăm metodele de estimare a scorurilor de încredere după îmbunătățirea probabilităților subcuvintelor prin tehnica de dropout. Această tehnică presupune transcrierea datelor de intrare de N=4, 8, 16, 32, respectiv 64 de ori folosind dropout. Mediem apoi cele N predicții independente. Tabelul 2.4.f prezintă rezultatele pentru toate combinațiile de caracteristici și agregări și dropout cu N=64.

**Tabelul 2.4.f** Rezultate pentru evaluarea tehnicii dropout (N=64) pentru îmbunătățirea scorurilor de încredere.

		RSC-eval / 1.8%			SSC-eval1 / 11.0%			SSC-eval2 / 14.0%		
caract.	agg.	AU-PR <sub>e</sub>	AU-PR <sub>s</sub>	AU-ROC	AU-PR <sub>e</sub>	AU-PR <sub>s</sub>	AU-ROC	AU-PR <sub>e</sub>	AU-PR <sub>s</sub>	AU-ROC
log-proba	sum	17.76	98.97	77.14	28.96	96.85	77.31	16.03	90.08	58.47
	min	18.90	99.17	82.02	29.01	97.10	79.90	13.58	90.73	56.77
	avg	15.73	99.13	81.02	27.75	96.95	79.22	13.02	90.78	56.22
neg-entropy	sum	9.53	98.75	71.19	20.40	96.20	71.68	15.18	90.45	56.80
	min	14.43	99.11	80.63	27.01	96.90	78.86	11.85	90.22	53.74
	avg	11.83	99.02	78.71	25.50	96.70	77.84	11.94	90.36	53.85

Comparând rezultatele din Tabelul 2.4.e cu rezultatele din Tabelul 2.4.f observăm îmbunătățiri pentru seturile de date RSC-eval și SSC-eval1, însă o scădere de performanță pentru setul de date SSC-eval2. Nici îmbunătățirile obținute pe RSC-eval și SSC-eval1 nu sunt foarte mari (AU-ROC crește de la 78.66 la 82.02 pe RSC-eval și de la 77.24 la 79.90 pe SSC-eval1).

Toate metricile de performanță se îmbunătățesc pe măsură ce sunt mediate mai multe rezultate de transcriere (pe măsură ce N crește) și, la fel ca și pentru limba engleză, creșterile de performanță se plafonează pentru N=64.

Un ultim considerent demn de menționat este acela că pentru a utiliza tehnica dropout în estimarea încrederii este nevoie de un timp de calcul net superior situației inițiale. Practic, pentru N=64 trebuie realizate 64 de transcrieri, deci timpul de calcul va fi de aproximativ 64 de ori mai mare. În concluzie, avantajele minime obținute prin utilizarea acestei tehnici pentru îmbunătățirea probabilităților subcuvintelor sunt contrabalansate de un dezavantaj semnificativ ce ține de timpul de calcul.

### **2.4.7 Utilizarea scorurilor de încredere propuse pentru generarea de date adnotate**

În activitatea 2.12 din etapa 2019 am propus utilizarea scorurilor de încredere la nivel de cuvânt puse la dispoziție de un sistem de RAV pentru a selecta cuvintele presupus transcrise corect de acesta. Selecția se face impunând un prag minim pentru scorul de încredere la nivel de cuvânt, prag sub care cuvântul în cauză este considerat ca fiind transcris greșit. Secvențele de cuvinte presupus transcrise corect, împreună cu datele audio corespunzătoare, ar urma să formeze un nou set de date adnotat. Pentru ca aplicarea acestei metode să genereze seturi de date adnotate suficient de mari, este necesar ca setul de date neadnotat utilizat la intrarea sistemului să aibă o dimensiune considerabilă. De exemplu, în acest proiect am utilizat seturile de date SSC-train3-raw și SSC-train4-raw care însumează 913 ore de vorbire. Metoda descrisă mai sus se poate aplica numai în cazul în care transcrierea acestui set de date neadnotat cu ajutorul sistemului de RAV se poate

realiza în mod eficient. Din păcate experimentele noastre au arătat că sistemul RAV ESPnet transcrie o oră de vorbire într-un interval de 2 până la 7 ore, în funcție de cât de similară este acustica respectivei ore de vorbire raportat la setul de date de antrenare. În acest context, transcrierea setului de date de 913 ore ar fi durat între 76 de zile și 266 de zile. În consecință, din considerente ce țin de infrastructura de calcul, nu am putut utiliza noile metode de estimare a scorurilor de încredere pentru generare de date audio adnotate.

## 2.5 Activitatea 3.14 - Diseminare

Diseminarea rezultatelor proiectului a fost realizată prin intermediul website-ului proiectului (<https://tadarav.speed.pub.ro>) și prin publicarea mai multor articole științifice după cum urmează:

1. 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.
2. 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.
3. 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.
4. 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. S, pp. S77–S89, 2020, ISSN 1453-8245, ISI IF 0.661.
5. 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.
6. 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.

Dintre articolele listate mai sus, al patrulea este deja indexat în Web of Science (Thompson Reuters - ISI), al doilea este deja indexat IEEE Xplore și în curs de indexare în Web of Science, primul și al treilea au apărut în volumul conferințelor respective și sunt în curs de indexare în Web of Science, iar al cincilea și al șaselea vor apărea în volumul conferinței respective și vor fi indexare în Web of Science.

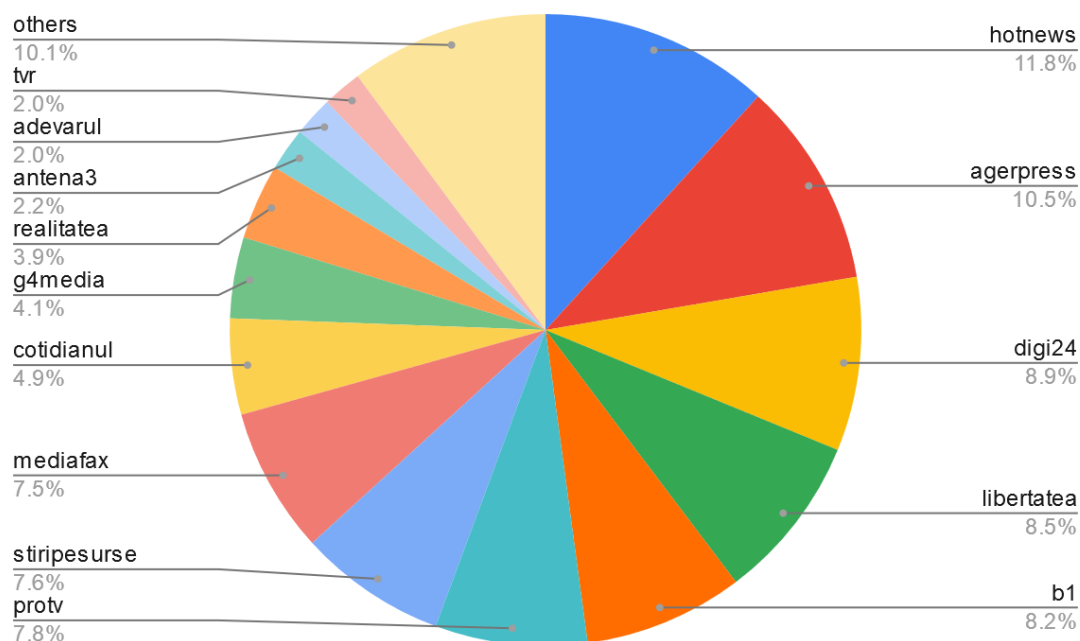
## 2.6 Crearea unui sistem de RAV îmbunătățit

### 2.6.1 Actualizarea modelelor de limbă pentru transcriere de vorbire

În această secțiune sunt prezentate sumar demersurile realizate de colectivul Speed în vederea actualizării modelelor de limbă utilizate în sistemele de transcriere de vorbire. Etapele parcurse și descrise în cele ce urmează sunt: (i) achiziția unui nou corpus de text pentru actualizarea modelelor de limbă; (ii) actualizarea aplicației de preprocesare a textului de antrenare pentru modelele de limbă; (iii) crearea unor noi modele de limbă; (iv) îmbunătățiri aduse sistemelor RAV prin utilizarea noilor modele de limbă. Aceste aspecte sunt prezentate în detaliu în [Manolache, 2020].

Un nou corpus de text, format din articole scrise, a fost achiziționat de pe mai multe website-uri de știri românești folosind o aplicație Java. Această aplicație a funcționat în mod continuu între 26.06.2018 - 22.05.2020, verificând periodic feed-urile RSS a 24 de website-uri de știri și descărcând orice articol întâlnit. Fiecare sursă de știri este extrasă pe un fir de execuție separat, făcând astfel procesul de extracție mai rapid. Procedura de extracție constă în descărcarea conținutului HTML, analiza acestuia în vederea determinării tag-urilor unde ar trebui să se regăsească textul și apoi preluarea textului propriu-zis.

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 2.6.a.



**Figura 2.6.a** Sursele de știri componente ale corpusului news2020. Procentele sunt relative la dimensiunea întregului corpus (255M cuvinte)

Aplicația de procesare a limbajului natural ce are scopul de normalizare a seturilor de date text ce urmează a fi folosite pentru antrenarea modelelor de limbă a fost supusă mai multor actualizări. Cea mai notabilă actualizare este legată de procesarea cuvintelor ce conțin cratimă. Această actualizare a fost necesară deoarece s-au întâmpinat mai multe cazuri în care sistemul RAV genera cuvinte care nu conțineau cratimă precum: ”ratat o”, ”lăsându se”, ”TVA ul”. Cauza principală a acestei probleme a fost aplicația de procesare a limbajului natural folosită pentru preprocesarea textului folosit mai departe pentru antrenarea modelelor de limbă folosite în sistemele RAV. Vechea aplicație de procesare a limbajului natural căuta cuvintele cu cratimă într-un lexicon dedicat. Dacă cuvântul era găsit în lexicon, acesta rămânea neschimbat, în caz contrar, cratima era înlocuită cu spațiu (e.g. “lovit o”, “lăsându se”, “FMI ul”). Lexiconul dedicat pentru cuvintele cu cratimă nu conținea toate cuvintele posibile cu cratimă deoarece acest lucru ar fi însemnat ca vocabularul pentru modelul de limbă să crească exponențial. De exemplu pentru acronimul SMURD, mai multe forme ar fi trebuit inserate în vocabular: SMURD-ul, SMURD-ului. Alt exemplu mai problematic îl reprezintă diferitele forme ale verbelor (e.g. gerunziu, imperativ): “lăsându-mă”, “lăsându-te”, “lăsându-l”, “nelăsându-ne”. Din cauza acestei probleme, lexiconul a fost limitat la cele mai frecvente cuvinte cu cratimă, iar restul cuvintelor întâlnite în text aveau cratima înlocuită cu spațiu. Acest fapt a dus la crearea unui model de limbă ce conține cuvinte incomplete: “ul”, “ului”, “lăsându” ce urmau să apară la în ieșirea sistemului ASR ca și erori. Soluția la această problemă a constat în modificarea procedurii de procesare a cuvintelor cu cratimă. Lexiconul a fost păstrat, dar numai pentru cuvinte compuse, precum: Târgu-Jiu. Pe lângă lexicon, au fost create 2 liste adiționale, una de prefixe cu cratimă (e.g. te-, ne-, istorico-), respectiv una de sufixe cu cratimă (e.g. -se, -vă, -ul). Având în vedere aceste resurse, procesarea cuvintelor cu cratimă întâlnite în text se face conform Algoritmului 2.6.a.

Astfel, după etapa de procesare de text, textul folosit pentru antrenarea modelului de limbă conținea cuvinte precum: “să”, “-mi”, “să”, “-i”, “-m-”, “a”, “ratat”, “-o”, “FMI”, “-ul”. Noul sistem ASR ce folosește modelul de limbă generează acum text de forma: “**TVA -ul** a crescut cu 3 la sută”. În final, un simplu procesor de text va unii cratimele de cuvintele adiacente: “**TVA-ul** a crescut cu 3 la sută”.

Cu privire la modelele de limbă noi, nu au fost create tipuri noi de modele de limbă, ci s-au 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 (352M de cuvinte), pe când noile modele de limbă au fost antrenate pe corpusul news002 plus corpusul news2020 (352M + 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

<b>if</b> (cuvântul cu cratimă este întâlnit în lexiconul pentru cuvinte cu cratimă)	Exemple
<b>then</b> se păstrează forma inițială	<i>Coca-Cola</i>
<b>else</b> cuvântul este împărțit	
<b>if</b> (una din părți se regăsește în lexiconul de cuvinte comune) && (cealaltă parte se regăsește în lista de sufixe/prefixe)	TVA <i>-ul</i>
<b>then</b> cuvântul și sufixul/prefixul sunt separate cu spațiu	<i>TVA-ul =&gt; TVA -ul</i>
<b>else</b> cratima este înlocuită cu un spațiu, iar forma inițială este trecută într-o listă de cuvinte nerecunoscute pentru o analiză manuală ulterioară	<i>două-trei</i> <i>bugetari-particulari</i>

*Algoritm 2.6.a Procedura actualizată de procesare a cuvintelor cu cratimă*

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.

Pentru etapa de decodare am experimentat cu diferite ordine n-gram (2-gram și 3-gram), dimensiuni ale vocabularului (200k, 250k and 300k words) și n-gram pruning (no pruning, 1e-7, 3e-7). Cele mai bune rezultatele obținute după evaluarea pe seturile de transcrieri pentru evaluare pot fi observate în Tabelul 2.6.a.

Din punct de vedere al PPL, cel mai bun model de limbă, LM-2020-3g-large-200k, prezintă o scădere relativă între 40% și 45% pe seturile de evaluare, comparativ cu modelul de limbă de bază (Tabel 2.6.a linia 1). Cu privire la rezultatele OOV, se poate observa o îmbunătățire în cazul seturilor de evaluare de vorbire spontană (SSC-eval1 și SSC-eval2). În cazul setului de evaluare de vorbire citită, RSC-eval, numai modelele cu un vocabular de 300k cuvinte au generat un rezultat mai bun.

Similar cu modelele de limbă pentru decodare prezentate în Tabelul 2.6.a, în Tabelul 2.6.b sunt evaluate modelele de limbă destinate pentru etapa de reevaluare lingvistică, anume 4-gram și RNN.

Analizând rezultatele din Tabelul 2.6.b, putem observa pentru modelul LM-2020-4g-large-200k o scădere relativă de 7% și 40% PPL pe seturile RSC-eval și SSC-eval2 comparativ cu modelul de bază cu rescoring (Tabelul 2.6.b, linia 1). Pentru setul de evaluare SSC-eval1, modelul LM-2020-5g-RNN a obținut rezultate mai bune, cu o scădere relativă de 14% PPL. Din punct de vedere al rezultatelor OOV, modelul LM-2020-4g-large-300k a generat cele mai bune rezultate pentru seturile de evaluare de vorbire spontană, pe când pentru setul RSC-eval cel mai bun rezultat a fost obținut la modelul de bază și LM-2020-5g-RNN. Procentajul de 0.01% OOV pentru modelul de bază și LM-2020-5g-RNN a fost obținut datorită faptului că aceste modele nu am avut o dimensiune limită pentru vocabular, față de modelele 4-gram care au avut limită.

**Tabelul 2.6.a** Modele de limbă pentru decodare îmbunătățite (LM-2020-\*) vs model de limbă pentru decodare de bază (LM-2017). Numele modelelor de limbă indică ordinul n-gram (2g, 3g), dacă s-a aplicat sau nu n-gram pruning (small = 3e-7, large = no pruning) și dimensiunea vocabularului (200k, 300k words)

Model de limbă	PPL [%]			OOV [%]		
	RSC-eval	SSC-eval1	SSC-eval2	RSC-eval	SSC-eval1	SSC-eval2
LM-2017-2g-large-200k	316	219	356	0.13	1.96	0.91
LM-2020-2g-small-300k	393	280	411	<b>0.09</b>	<b>0.37</b>	<b>0.28</b>
<b>LM-2020-3g-large-200k</b>	<b>178</b>	<b>131</b>	<b>210</b>	0.19	0.51	0.37

**Tabelul 2.6.b** Modele de limbă pentru rescoring îmbunătățite (LM-2020-\*) vs model de limbă pentru rescoring de bază (LM-2019). Numele modelelor de limbă indică ordinul n-gram (4g, RNN), n-gram pruning (large = no pruning) și dimensiunea vocabularului (200k, 250k, 300k words)

Model de limbă	PPL [%]			OOV [%]		
	RSC-eval	SSC-eval1	SSC-eval2	RSC-eval	SSC-eval1	SSC-eval2
LM-2019-RNN	176	132	325	0.01	0.82	0.66
LM-2020-RNN	246	<b>113</b>	206	<b>0.01</b>	0.43	0.35
LM-2020-4g-large-200k	<b>162</b>	118	<b>196</b>	0.19	0.51	0.37
LM-2020-4g-large-250k	163	119	197	0.14	0.44	0.34
LM-2020-4g-large-300k	165	120	198	0.09	<b>0.37</b>	<b>0.28</b>

Noile sisteme RAV au fost obținute folosind modelele de limbă prezentate anterior. 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 2.6.c sunt prezentate rezultatele WER ale noilor sisteme RAV cu și fără rescoring.

În Tabelul 2.6.c 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 2.6.c linia 1). Însă, pe setul de evaluare RSC-eval sistemul RAV a înregistrat o creștere de 16% WER.

## 2.6.2 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 2.6.d 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. Acesta a fost antrenat folosind setul de antrenare *baseline*: RSC-train1 + SSC-train1 + SSC-train2.

**Tabelul 2.6.c** 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	<b>1.8</b>	11.0	14.0
LM-2020-RNN	LM-2020-3g-large-200k	<b>2.7</b>	<b>11.0</b>	<b>14.4</b>	2.1	<b>10.0</b>	<b>12.9</b>
	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

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 aliniată 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. În secțiunea 2.3.2 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. 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 2.6.e 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 2.6.d.** 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 2.6.e.** 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 2.6.f 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. 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 2.6.f.** 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

### 2.6.3 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 2.6.1) și model acustic (conform secțiunii 2.6.2), cel mai performant sistem de RAV Speed are actualmente performanțele listate în Tabelul 2.6.g.

*Tabelul 2.6.g. 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

## 2.7 Bibliografie

- [Ardila, 2020] Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M. Tyers, and Gregor Weber, “Common Voice: A massively-multilingual speech corpus,” in International Conference on Language Resources and Evaluation, 2020.
- [Ashukha, 2020] Arsenii Ashukha, Alexander Lyzhov, Dmitry Molchanov, and Dmitry Vetrov, “Pitfalls of in-domain uncertainty estimation and ensembling in deep learning,” in International Conference on Learning Representations, 2020.
- [Barbu, 2018] V. Barbu-Mititelu, D. Tufiş, E. Irimia, “*The Reference Corpus of the Contemporary Romanian Language (CoRoLa)*,” in Proceedings of LREC 2018, Japan, p.1178-1185.
- [Chen, 2019] Tongfei Chen, Jiri Navrátil, Vijay Iyengar, and Karthikeyan Shanmugam, “Confidence scoring using whitebox meta-models with linear classifier probes,” in International Conference on Artificial Intelligence and Statistics, 2019, pp. 1467–1475.
- [Corbière, 2019] Charles Corbière, Nicolas Thome, Avner Bar-Hen, Matthieu Cord, and Patrick Pérez, “Addressing failure prediction by learning model confidence,” in Advances in Neural Information Processing Systems, 2019, pp. 2902–2913.
- [Cortina, 2016] Isaías Sánchez Cortina, Jesús Andrés-Ferrer, Alberto Sanchis, and Alfons Juan, “Speaker-adapted confidence measures for speech recognition of video lectures,” *Computer Speech & Language*, vol. 37, pp. 11–23, 2016.
- [Del-Agua, 2018] M. A. Del-Agua, A. Gimenez, A. Sanchis, J. Civera, and A. Juan, “Speaker-adapted confidence measures for ASR using deep bidirectional recurrent neural networks,” *Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 7, pp. 1198–1206, 2018.
- [Errattahi, 2018] Rahhal Errattahi, Salil Deena, Asmaa El Hannani, Hassan Ouahmane, and Thomas Hain, “Improving ASR error detection with RNNLM adaptation,” in IEEE Spoken Language Technology Workshop, 2018, pp. 190–196.
- [Gal, 2016] Yarin Gal and Zoubin Ghahramani, “Dropout as a Bayesian approximation: Representing model uncertainty in deep learning,” in International Conference on Machine Learning, 2016, pp. 1050–1059.
- [Georgescu, 2017] A.-L. Georgescu, H. Cucu, C. Burileanu, “*Speed’s DNN Approach to Romanian Speech Recognition*,” in Proc. 9th Conference on Speech Technology and Human-Computer Dialogue (SpeD), 8p, 2017.
- [Georgescu, 2018] A.-L. Georgescu, H. Cucu, “Automatic annotation of speech corpora using complementary GMM and DNN acoustic models,” In Proc. TSP 2018, pp. 1-4.
- [Georgescu, 2019a] A.-L. Georgescu, C. Manolache, G. Pop, D. Oneață, H. Cucu, D. Burileanu, C. Burileanu, “*Proiect component TADARAV: Raport științific și tehnic în extenso 2019*”.



- [Georgescu, 2019b] A.-L. Georgescu, H. Cucu, C. Burileanu, “*Kaldi-based DNN architectures for speech recognition in Romanian*,” in the Proceedings of the 10th Conference on Speech Technology and Human-Computer Dialogue (SpeD), 2019, Timișoara, Romania.
- [Georgescu, 2020] 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, 2021] 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.
- [Guo, 2017] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger, “On calibration of modern neural networks,” in International Conference on Machine Learning, 2017, pp. 1321–1330.
- [Hadian, 2018] Hossein Hadian, Hossein Sameti, Daniel Povey, and Sanjeev Khudanpur, “End-to-end speech recognition using lattice-free MMI,” in Interspeech, 2018, pp. 12–16.
- [Hazen, 2002] Timothy J Hazen, Stephanie Seneff, and Joseph Polifroni, “Recognition confidence scoring and its use in speech understanding systems,” *Computer Speech & Language*, vol. 16, no. 1, pp. 49–67, 2002.
- [Hendrycks, 2016] Dan Hendrycks and Kevin Gimpel, “A baseline for detecting misclassified and out-of-distribution examples in neural networks,” in International Conference on Learning Representations, 2016.
- [Hinton, 2015] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean, “Distilling the knowledge in a neural network,” arXiv preprint arXiv:1503.02531, 2015.
- [Hori, 2019] Takaaki Hori, Shigeki Karita, Tomoki Hayashi, Shinji Watanabe, “*Advanced Methods for Neural End-to-End Speech Processing – Unification, Integration, and Implementation*,” Interspeech 2019.
- [Huggins-Daines, 2006] David Huggins-Daines, Mohit Kumar, Arthur Chan, Alan W. Black, Mosur Ravishankar, and Alexander I. Rudnicky. “*Pocketsphinx: A free, real-time continuous speech recognition system for hand-held devices*.” in Proc. 2006 IEEE International Conference on Acoustics Speech and Signal Processing Proceedings, vol. 1, pp. I-I. 2006.
- [Kalgaonkar, 2015] Kaustubh Kalgaonkar, Chaojun Liu, Yifan Gong, and Kaisheng Yao, “Estimating confidence scores on ASR results using recurrent neural networks,” in IEEE International Conference on Acoustics, Speech and Signal Processing, 2015, pp. 4999–5003.
- [Karita, 2019] Shigeki Karita, Nanxin Chen, Tomoki Hayashi, Takaaki Hori, Hirofumi Inaguma, Ziyang Jiang, Masao Someki, Nelson Enrique Yalta Soplín, Ryuichi Yamamoto, Xiaofei Wang, Shinji Watanabe, Takenori Yoshimura, and Wangyou Zhang, “A comparative study on transformer vs RNN in speech applications,” in Workshop on Automatic Speech Recognition and Understanding, 2019, pp. 449–456.
- [Kemp, 1997] Thomas Kemp and Thomas Schaaf, “Estimating confidence using word lattices,” in Eurospeech, 1997.
- [Kim, 2017] Suyoun Kim, Takaaki Hori, and Shinji Watanabe, “Joint CTC-attention based end-to-end speech recognition using multi-task learning,” in Proc. ICASSP, pp. 4835-4839, 2017.
- [Kudo, 2018] Taku Kudo, “Subword regularization: Improving neural network translation models with multiple subword candidates,” in Association for Computational Linguistics, 2018, pp. 66–75.
- [Lakshminarayanan, 2017] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell, “Simple and scalable predictive uncertainty estimation using deep ensembles,” in Advances in Neural Information Processing Systems, 2017, pp. 6402–6413.
- [Li, 2019] Qiuqia Li, PM Ness, Anton Ragni, and Mark JF Gales, “Bi-directional lattice recurrent neural networks for confidence estimation,” in IEEE International Conference on Acoustics, Speech and Signal Processing, 2019, pp. 6755–6759.
- [Lüscher, 2019] Christoph Lüscher, Eugen Beck, Kazuki Irie, Markus Kitzka, Wilfried Michel, Albert Zeyer, Ralf Schlüter, and Hermann Ney, “RWTH ASR systems for LibriSpeech: Hybrid vs attention,” in Interspeech, 2019, pp. 231–235.

- [Malinin, 2020] Andrey Malinin and Mark Gales, “Uncertainty in structured prediction,” arXiv preprint arXiv:2002.07650, 2020.
- [Manolache, 2019] C. Manolache, “*Natural language processing using artificial intelligence: keyword spotting and speech transcription intelligibility*“, MSc. Thesis, June 2019 (scientific coordinator: Conf. H. Cucu).
- [Manolache, 2020] 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.
- [Ogawa, 2017] Atsunori Ogawa and Takaaki Hori, “Error detection and accuracy estimation in automatic speech recognition using deep bidirectional recurrent neural networks,” *Speech Communication*, vol. 89, pp. 70–83, 2017.
- [Ovadia, 2019] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek, “Can you trust your model’s uncertainty? Evaluating predictive uncertainty under dataset shift,” in *Advances in Neural Information Processing Systems*, 2019, pp. 13991–14002.
- [Panayotov, 2015] V. Panayotov, G. Chen, D. Povey and S. Khudanpur, “*Librispeech: An ASR corpus based on public domain audio books*,” in Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 5206-5210, 2015.
- [Park, 2019] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le, “SpecAugment: A simple data augmentation method for automatic speech recognition,” in *Interspeech*, 2019, pp. 2613–2617.
- [Paszke, 2019] Adam Paszke et al., “*PyTorch: An Imperative Style, High-Performance Deep Learning*,” *Advances in Neural Information Processing Systems* 32, pp. 8024-8035, 2019.
- [Povey, 2011] Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, Jan Silovsky, Georg Stemmer, and Karel Vesely, “The Kaldi speech recognition toolkit,” in *Workshop on Automatic Speech Recognition and Understanding*, 2011.
- [Ragni, 2018] Anton Ragni, Qiujia Li, Mark JF Gales, and Yongqiang Wang, “Confidence estimation and deletion prediction using bidirectional recurrent neural networks,” in *IEEE Spoken Language Technology Workshop*, 2018, pp. 204–211.
- [Rousseau, 2014] Anthony Rousseau, Paul Deléglise, and Yannick Estève, “Enhancing the TED-LIUM corpus with selected data for language modeling and more TED talks,” in *International Conference on Language Resources and Evaluation*, 2014, pp. 3935–3939.
- [Seigel, 2013] Mathew Stephen Seigel, *Confidence estimation for automatic speech recognition hypotheses*, Ph.D. thesis, University of Cambridge, 2013.
- [Seigel, 2014] Matthew Stephen Seigel and Philip C Woodland, “Detecting deletions in ASR output,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2014, pp. 2302–2306.
- [Sennrich, 2016] Rico Sennrich, Barry Haddow, Alexandra Birch, “Neural Machine Translation of Rare Words with Subword Units,” In Proc. 54th Annual Meeting of the Association for Computational Linguistics, pp 1715-1725, 2016.
- [Sperber, 2017] Matthias Sperber, Graham Neubig, Jan Niehues, and Alex Waibel, “Neural lattice-to-sequence models for uncertain inputs,” in *Empirical Methods in Natural Language Processing*, 2017, pp. 1380–1389.
- [Srivastava, 2014] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov, “Dropout: A simple way to prevent neural networks from overfitting,” *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 1 2014.

- [Swarup, 2019] Prakhar Swarup, Roland Maas, Sri Garimella, Sri Harish Mallidi, and Björn Hoffmeister, “Improving ASR confidence scores for Alexa using acoustic and hypothesis embeddings,” in *Interspeech*, 2019, pp. 2175–2179.
- [Tokui, 2019] Seiya Tokui, Ryosuke Okuta, Takuya Akiba, Yusuke Niitani, Toru Ogawa, Shunta Saito, Shuji Suzuki, Kota Uenishi, Brian Vogel, and Hiroyuki Yamazaki Vincent, “*Chainer: A deep learning framework for accelerating the research cycle*,” In Proc. of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pp. 2002–2011, 2019.
- [Tüske, 2019] Zoltán Tüske, Kartik Audhkhasi, and George Saon, “Advancing sequence-to-sequence based speech recognition,” in *Interspeech*, 2019, pp. 3780–3784.
- [Vaswani, 2017] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.
- [Vesely, 2013] Karel Vesely, Mirko Hannemann, and Lukas Burget, “Semi-supervised training of deep neural networks,” in *Workshop on Automatic Speech Recognition and Understanding*, 2013, pp. 267–272.
- [Vyas, 2019] Apoorv Vyas, Pranay Dighe, Sibong Tong, and Hervé Bourlard, “Analyzing uncertainties in speech recognition using dropout,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2019, pp. 6730–6734.
- [Watanabe, 2017] Shinji Watanabe, Takaaki Hori, Suyoun Kim, John R. Hershey and Tomoki Hayashi, “*Hybrid CTC/Attention Architecture for End-to-End Speech Recognition*,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 8, pp. 1240–1253, 2017.
- [Watanabe, 2018] Shinji Watanabe, Takaaki Hori, Shigeki Karita, Tomoki Hayashi, Jiro Nishitoba, Yuya Unno, Nelson Enrique Yalta Soplín, Jahn Heymann, Matthew Wiesner, Nanxin Chen, Adithya Renduchintala, and Tsubasa Ochiai, “ESPnet: End-to-End Speech Processing Toolkit,” in *Proc. Interspeech*, pp. 2207–2211, 2018.
- [Weintraub, 1997] Mitch Weintraub, Françoise Beaufays, Zeév Rivlin, Yochai Konig, and Andreas Stolcke, “Neural-network based measures of confidence for word recognition,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 1997, vol. 2, pp. 887–890.
- [Yu, 2010] Dong Yu, Balakrishnan Varadarajan, Li Deng, and Alex Acero, “Active learning and semi-supervised learning for speech recognition: A unified framework using the global entropy reduction maximization criterion,” *Computer Speech & Language*, vol. 24, no. 3, pp. 433–444, 2010.

### **3 Structura ofertei de servicii de cercetare și tehnologice**

Laboratorul de cercetare *Speech and Dialogue* (Speed) din cadrul Universității Politehnica din București (UPB), reprezentantul UPB în proiectul TADARAV, oferă pe platforma ERRIS serviciile de cercetare și tehnologice enumerate în Tabelul 3.

*Tabelul 3. Servicii de cercetare și tehnologice oferite de Laboratorul de cercetare Speech and Dialogue*

<b>Serviciu</b>	<b>Detalii</b>
Serviciu și aplicație web de transcriere de documente ce conțin vorbire în limba română	<a href="https://transcriptions.speed.pub.ro">https://transcriptions.speed.pub.ro</a>
Serviciu și aplicație web de identificare de cuvinte cheie în documente ce conțin vorbire în limba română	<a href="https://keywords.speed.pub.ro">https://keywords.speed.pub.ro</a>
Serviciu și aplicație web de restaurare de diacritice în limba română	<a href="https://diacritics.speed.pub.ro">https://diacritics.speed.pub.ro</a>
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

Laboratorul de cercetare *Speech and Dialogue* (SpeeD) este prezent pe platforma ERRIS la adresa <https://erris.gov.ro/SpeeD---UPB>.

#### **4 Locuri de muncă susținute prin program**

Echipa de cercetare a Universității Politehnica din București pentru proiectul component TADARAV este prezentată în Tabelul 4.

*Tabelul 4. Echipa de cercetare UPB*

<b>Nr.</b>	<b>Nume</b>	<b>Calitatea</b>	<b>Poziția</b>	<b>Normă</b>
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	<b>Membru cercetător nou</b>	<b>Întreagă</b>
6	Gheorghe POP	ACS	<b>Membru cercetător nou</b>	<b>Întreagă</b>
7	Cristian MANOLACHE	ACS	<b>Membru cercetător nou</b>	<b>Întreagă</b>

#### **5 Valorificarea și îmbunătățirea competențelor și resurselor existente la nivelul consorțiului**

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