

## **RAPORTARE ȘTIINȚIFICĂ**

### **FAZA DE EXECUȚIE NR. 3**

**CU TITLUL** Dezvoltarea modulului de limba romana, testarea si evaluarea metodelor proiectate in etapele anterioare

**Avizat,**

**Coordonator**

**Denumire:** Universitatea din Bucuresti

**Reprezentant Legal:** prof. dr. Mircea Dumitru

**Semnătură:**

**Ștampilă:**

**Agent Economic**

**Denumire:** Pluriva SRL

**Reprezentant Legal:**

**Semnătură:**

**Ștampilă:**

Director Proiect: prof.dr. Liviu P. Dinu

Semnătură:



Responsabil de proiect:

Semnătură:

## **Raport științific**

*privind implementarea proiectului in perioada ianuarie – septembrie 2018*

**Proiect:** *Sistem Inteligent de Generare Automată a Răspunsurilor (SIGAR).*

*SIGAR: Project 53BG/2016, funded by Romanian National Authority for Scientific Research and Innovation, CNCS/CCCDI – UEFISCDI, PNCDI III (Programul 2 - Creșterea competitivității economiei românești prin cercetare, dezvoltare și inovare. Transfer de cunoaștere la agentul economic “Bridge Grant”)*

**Durata Proiectului:** 1 octombrie 2016- 30 septembrie 2018

**Director:** Prof. univ. dr. Liviu P. Dinu

**Contractor:** Universitatea din Bucuresti

**Etapa:** 30 septembrie 2018 (unica)

### **Obiectivul principal**

Principalul obiectiv al proiectului constă în implementarea unei soluții bazate pe învățare automată, care va parsa și analiza baza de date de tichete în vederea construirii unor răspunsuri automate la întrebări sau probleme frecvent întâlnite.

**Obiectivul** proiectului pentru anul 2018 este:

**1. Dezvoltarea modului de limba romana, testarea si evaluarea metodelor proiectate in etapele anterioare.** Pentru îndeplinirea acestui obiectiv, au fost planificate următoarele activități:

Activitate 3.1: **Dezvoltarea de resurse specific proiectului pentru limba romana.**

În scopul realizării acestei activități următoarele acțiuni au fost planificate și realizate:

- analiza morfologica a datelor
- extragerea părților de vorbire pentru datele disponibile
- finalizarea adnotarilor pentru limba romana

stagii de  
in

Activitate 3.2: **Activitati diseminare.** Principalele acțiuni desfășurate în cadrul acestei activități au fost participarea la conferințe, efectuarea unor cercetare, diseminarea rezultatelor in intalniri formale si informale, organizarea unui seminar de cercetare, etc.

Activitate 3.3: **Evaluare si testare module pentru limba romana**

În scopul realizării acestei activități următoarele acțiuni au fost planificate și realizate:

- analiza datelor fara diacritice
- detectarea entitatilor in romana
- evaluare, testare, ajustare

Activitate 3.4: **Finalizarea modulelor, evaluarea, ajustarea rezultatelor.** Principalele acțiuni desfășurate în cadrul acestei activități au fost finalizarea resurselor pentru generarea raspunsurilor, testarea, evaluarea si ajustarea metodelor dezvoltate in etapele anterioare, intocmire rapoarte finale.

În cadrul acestei etape, principalul punct de lucru a fost determinat de optimizarea și testarea metodelor prezentate anterior pe limba română, în particular pe limba română utilizată în cadrul datelor Pluriva. O caracteristică aparte a acestor date constă în faptul că este adesea întâlnită o terminologie specifică, determinată de limbajul business și tehnic folosit de către utilizatorii sistemului, care nu este întâlnită în mod general în cadrul resurselor lingvistice publice disponibile pentru limba română. Mai mult decât atât, unelele deja existente pentru limba română neavând terminologia învățată în cadrul modelelor automate și nici acces la resurse similare, prezintă o eficiență foarte scăzută atunci când sunt aplicate în cazul de față. Procesul de curățare a datelor din Etapa 2, ne-a permis să construim câteva resurse specifice integrate și optimizate pentru procesarea datelor Pluriva. Printre acestea numărăm o metodă de adnotare a părților de vorbire care include și entități nume proprii, modele de word embeddings specifice limbii române utilizate în datele Pluriva și o unealtă de extragere a conținutului esențial din schimburile de mesaje dintre agenți și clienți.

## **Rezultatele activitatii in cadrul proiectului in perioada 1 ianuarie-30 septembrie 2018**

### **Activitatea științifică:**

În anul 2018 activitatea științifică a membrilor proiectului a constat în special pe analiza tipului de discurs utilizat în datele Pluriva, analiza morfologică și implementarea unui detector de părți de vorbire customizat pentru aceste date, capabil de asemenea să indice și prezența unor entități de tip nume propriu. În paralel au fost proiectați algoritmi care să identifice stări sau sentimente din cadrul textului, algoritmi care au un potențial practic în situația în care un client nemulțumit își exprimă această părere în interacțiunile cu agenții firmei.

În ultima parte a proiectului am dorit să ne concentrăm pe cercetarea abordărilor care s-ar dovedi fezabile pentru procesarea textelor în limba română și mai exact pentru procesarea tipurilor de texte pe care agentul economic le stochează. Conform analizelor efectuate în cadrul activităților din Etapa 2, am constatat că aceste texte nu doar că prezintă un limbaj de specialitate, dar că acoperă și un lexic care nu este standard și care nu poate fi analizat utilizând resursele disponibile pentru limba română (DEX online, Romanian Treebank, POS taggers sau Named Entity Recognizers) din cauza lipsei adnotărilor pre-existente pentru acest tip de date.

Un part of speech (POS) tagger sau detector de părți de vorbire reprezintă o unealtă esențială pentru orice tip de analiză text. În cazul nostru am analizat uneltele și resursele deja existente pentru limba română și am concluzionat că acestea nu doar că prezintă o performanță relativ scăzută pentru limba română, dar că pot fi aproape inutile pentru datele agentului economic care conțin un lexic specializat ce nu face parte decât pe jumătate din vocabularul limbii române așa cum este el normat în DEX. Așadar, am fost nevoiți să găsim soluții alternative pentru adnotarea automată a acronimelor, cuvintelor împrumutate din engleză sau a celor ne-normate. Am ales să utilizăm corpusul *Romanian UD treebank (RoRefTrees)* (Barbu Mititelu et al., 2016) adnotat deja cu părți de vorbire și arbori sintactici. Corpusul are o distribuție neuniformă de genuri: literatură - 1818 propoziții, legislativ - 1606 propoziții, medical - 1210 propoziții, traduceri din FrameNet - 1092 propoziții, texte academice - 950 propoziții, știri - 933 propoziții, științific - 362 propoziții, wikipedia - 251 propoziții, altele - 1301 propoziții. În total conține 218 mii de cuvinte adnotate cu părți de vorbire. Utilizăm ultima versiune a acestui corpus pentru a extrage seturi de reguli de adnotare a părților de vorbire. Ca mecanism de învățare utilizăm Single Classification Ripple Down Rules (RDR) care s-au dovedit a fi eficiente (*Dat Quoc Nguyen et al., 2016*) pentru acest tip de problema, prin care înățăm anumite structuri condiționale de tip arbore de decizie care conduc la determinarea unei părți de vorbire în funcție de contextul în care se află un cuvânt. Avantajul acestei metode constă în simplitate și capacitatea de a fi descriptivă și a modifica regulile învățate din corpus prin supervizare umană.

Câteva exemple de reguli învățate:

- cuvântAnterior1 == "va" : conclusion = "VERB"
- cuvânt == "până" și următorulCuvânt == "ce" : conclusion = "ADP"
- suffixUltimele3Litere == "șpe" : conclusion = "NUM"
- cuvântAnterior1 == "să" : conclusion = "VERB"
- cuvânt == "este" and următorulPOS == "CCONJ" : conclusion = "AUX"
- cuvântAnterior2 == "liceul" : conclusion = "PROP"

În plus, regulile învățate din corpus au fost verificate augmentate manual iar cuvintele care lipsesc în corpusul de antrenare au fost analizate separat astfel:

**Pas 1.** determină dacă un cuvânt lipsește din dicționarul extras din datele de antrenare sau pentru care nu se poate deduce partea de vorbire folosind una din regulile anterioare (de exemplu acronime: fac (factură), Ms (mulțumesc), task, sample, etc.)

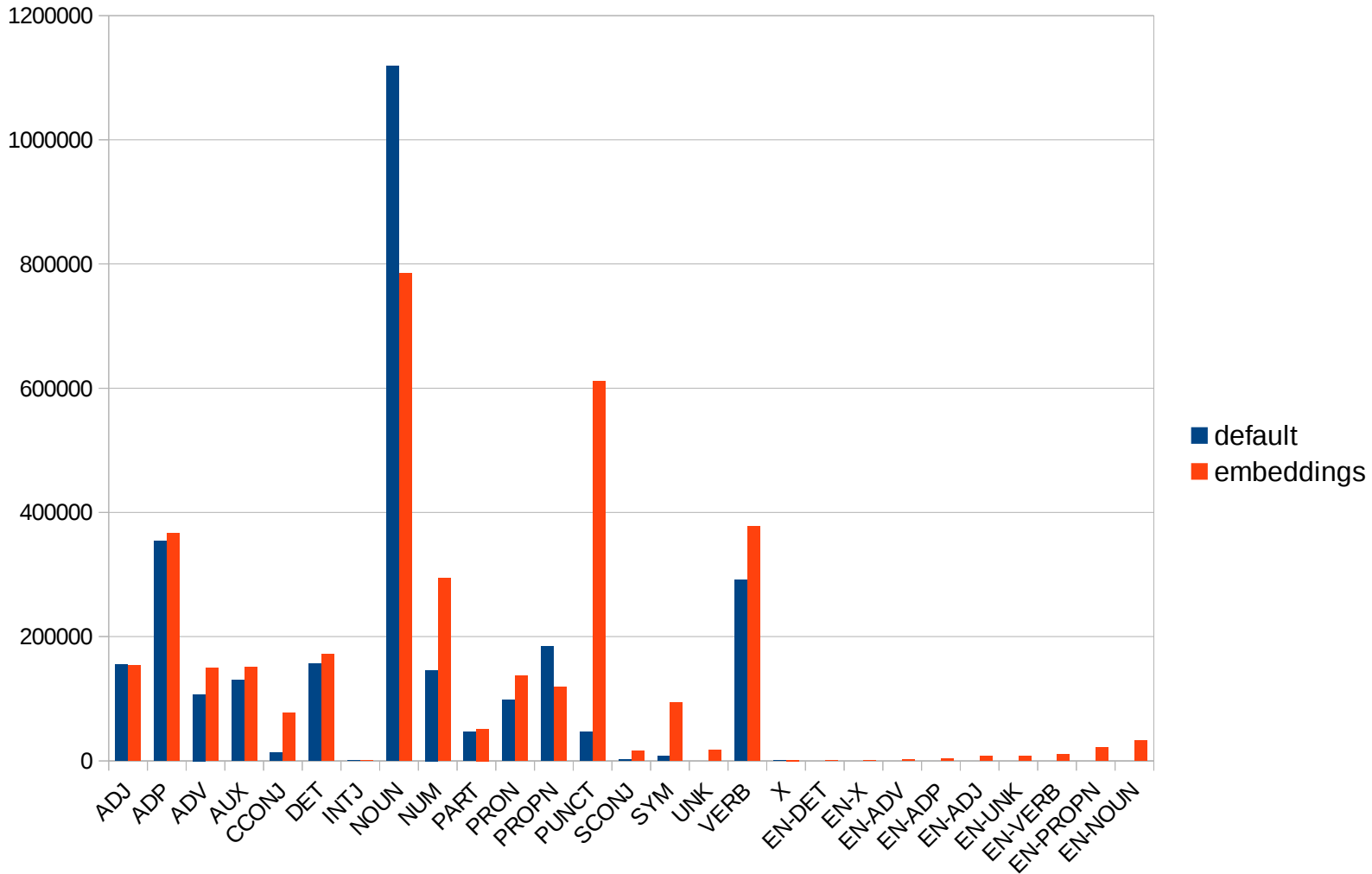
**Pas 2.** verifică dacă există word embeddings antrenate care includ acel cuvânt

**Pas 3.** dacă nu există, partea de vorbire a cuvântului este UNK (necunoscut)

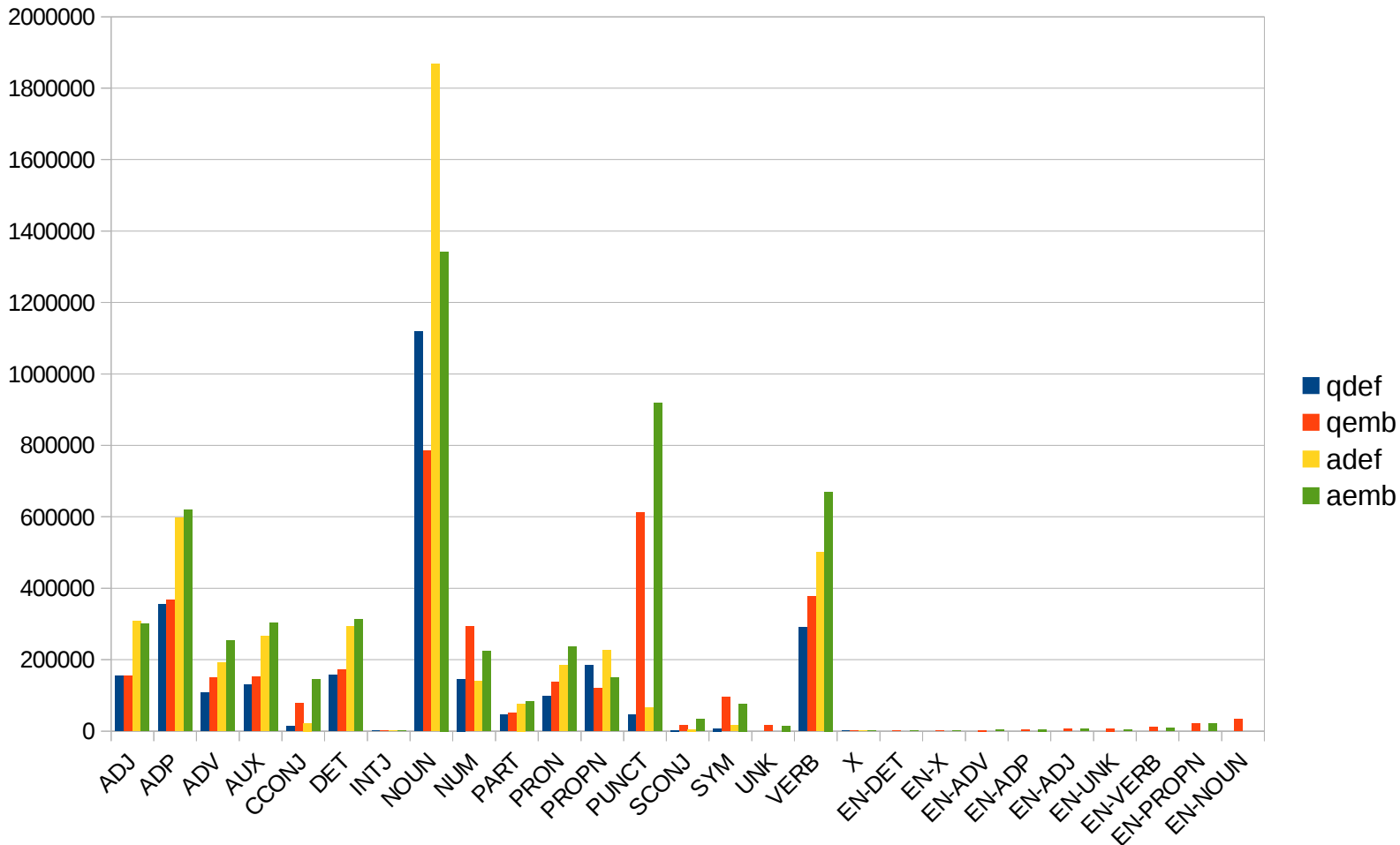
**Pas 4.** altfel, caută cele mai similare 10 cuvinte în funcție de distanța euclidiană dintre reprezentările vectoriale ale acestora și a cuvântului dat

**Pas 5.** atribuie cuvântului partea de vorbire a celui mai apropiat vecin în raport cu distanța euclidiană, care fie se găsește în dicționar sau pentru care se poate aplica una din regulile anterioare

Și nu în ultimul rând, am dorit să identificăm dacă din cuvintele adnotate, avem de fapt cuvinte provenite din engleză pentru care nu am putea identifica partea de vorbire cu foarte multă acuratețe sau care ar putea rezulta ca o excepție de la regulile de adnotare. Așadar, am ales să adnotăm părțile de vorbire cu un prefix "EN-", dacă acel cuvânt face parte din synsets existente în WordNet. În plus, am implementat funcții care să trateze automat lipsa diacriticelor din corpusul nostru, astfel că sunt executate căutări ale cuvintelor cu și fără diacritice. În mod similar, am mai introdus câteva reguli de adnotare care să identifice cuvinte cu formă de titlu sau de nume propriu. În graficul de mai jos sunt reprezentate distribuțiile părților de vorbire adnotate pe corpusul de întrebări. Modelul de vectori de cuvinte folosit este de asemenea antrenat pe același corpus, care a fost pre-procesat și curățat înainte. Modelul de bază (în albastru, default) are tendința de a identifica mai multe substantive (NOUN), dar o analiză manuală mai atentă asupra adnotărilor generate a arătat că aceasta este tendința algoritmului de a adnota substantive când partea de vorbire este necunoscută, generând multe greșeli pentru cuvinte care nu apar în dicționarul învățat din corpus. Pe corpusul nostru, acuratețea predicțiilor este de 65% cu sistemul de bază, pe când acuratețea pe datele de antrenare este de 95.75%. Folosind structurile vectoriale de reprezentări ale cuvintelor, am putut introduce în adnotator informații suplimentare pe baza cărora deciziile sunt mult mai bune. După cum se observă și în diagrama de mai jos, adnotatorul cu embeddings identifică un număr semnificativ mai mare de cuvinte din fiecare categorie, mai puțin categoria substantivelor. Are loc o redistribuire a părților de vorbire ceea ce face ca sistemul să fie mai robust și cu un nivel de acuratețe mai mare (85.3%) conform unor evaluări manuale ale rezultatelor adnotate.



Distribuțiile comparative ale părților de vorbire pentru corpusul de întrebări (qdef și qemb) și pentru cel de răspunsuri (adef și aemb) se regăsesc în următoarea diagramă:



Numărul de substantive greșit adnotate este cu mult mai mare în corpusul de răspunsuri (galben și verde), comparativ cu cel de întrebări, ceea ce face și redistribuirea lor mai amplă în cadrul celorlalte părți de vorbire. În ambele corpora observăm și o redistribuire a numelor proprii, numărul total al acestora scăzând după utilizarea reprezentărilor vectoriale.

Un alt aspect demn de observat este faptul ca un numar de aproximativ 86000 de părți de vorbire din corpusul de întrebări sunt parte din limba engleză, iar un număr comparabil de 82000 de părți de vorbire din corpusul de răspunsuri au ca rezultat prefixarea cu EN-. Ceea ce ne indică faptul că acele cuvinte sunt greu de adnotat fără eroare și că deși ele sunt adnotate cu o parte de vorbire (DET, ADV, VERB etc.), cel mai probabil este rezultatul algoritmului de antrenare. În acest fel, putem filtra cuvintele străine sau le putem pune în containere separate pentru procesare.

Metodele de detectare a entităților necesită un corpus vast adnotat manual cu anumite tipuri de entități și cu părțile de vorbire ale cuvintelor care sunt vecine cu acele entități. În cazul nostru, neavând la dispoziție o resursă de asemenea magnitudine, am ales să ne folosim de tag-ul PROPN pentru a identifica numele proprii. În această categorie includem nume de firmă, nume de persoane, nume de oraș etc. Toate aceste nume se identifică cu ușurință în corpusul agentului economic deoarece aceste cuvinte sunt în mod majoritar scrise cu literă mare la început. În mod paradoxal, deși schimburile de replici dintre agenți și clienți sunt caracterizate de un număr semnificativ de cuvinte nenormate și lipsa diacriticelor, în cazul numelor proprii este păstrată

norma de scriere, mai exact 25.7% din cuvintele înregistrate în modelele antrenate sunt scrise cu literă mare, iar aproximativ 8% sunt cuvinte aflate la începutul unei propoziții, rămânând un procent de 17.7% de entități nume proprii pentru care putem deduce tag-ul doar pe baza ortografiei. Acestea au fost adăugate în dicționarul de adnotare folosit de către POS tagger.

Pentru finalizarea resurselor, am implementat o metodă de curățare și extragere a conținutului esențial din date. Această metodă poate fi folosită pentru pre-procesare și trimitere a conținutului esențial către agregator sau poate fi folosită pentru extragerea unui sumar care să fie livrat agentului care citește conținutul ticketului. De asemenea poate fi folosită pentru îmbunătățirea categorisirilor de date pe subiecte și pentru indexare și căutarea mai rapidă a elementelor duplicate.

Pentru a antrena un model de extragere de conținut esențial, am construit un corpus paralel din perechi de documente aliniate linie cu linie. Dintr-un document care conține o secvență de întrebări – răspunsuri, am eliminat liniile care nu conțin informații esențiale legate de subiectul abordat. Aceste informații pot cuprinde propoziții întregi, elemente de structura a mailurilor (header, formulări), tabele html, elemente de conversație. Corpusul inițial cuprinde 2000 de documente selectate aleatoriu cumulând în total peste 200,000 de cuvinte, iar în urma înlăturării elementelor care nu fac parte din conținutul esențial, rămânem cu un corpus filtrat care conține 73,000 de cuvinte. Utilizăm cele două componente aliniate pentru a eticheta liniile într-un format binar – linii care trebuie păstrate sau care trebuie scoase. Având liniile etichetate putem antrena un clasificator care să distingă între liniile care sunt irelevante și cele cu conținut esențial. Pentru asta am extras trăsături la nivel lexical folosind tf-idf și am antrenat un model de regresie logistică pe acele perechi de trăsături – etichete. Un dezavantaj major al acestei metode constă în faptul că nu ia în considerare propozițiile vecine în momentul adnotării. Pentru a trata această problemă, am utilizat probabilitățile liniilor vecine de a fi parte din conținut astfel:

- dacă clasificatorul afișează 1, atunci linia este cu conținut esențial
- dacă clasificatorul afișează -1 și probabilitatea cumulată a propoziției curente cu cea anterioară este mai mare de 0.2, atunci linia este cu conținut esențial
- dacă propoziția anterioară este -1 și probabilitatea cumulată a propoziției curente cu cea anterioară este mai mare 0.2, atunci linia anterioară este cu conținut esențial
- altfel, linia curentă nu este cu conținut esențial

Pragul de 0.2 a fost stabilit în urma mai multor procese de cross-validare, iar cumularea probabilităților liniilor de a fi sau nu parte din conținut esențial îmbunătățește considerabil calitatea rezultatelor.

În urma implementării, am decis să extragem conținutul esențial din întregul corpus și evaluarea rezultatelor după un criteriu manual și subiectiv. Astfel că, am inspectat un sample de 1000 de texte din care am extras conținut relevant pentru a observa ce tipuri de greșeli face algoritmul, mai exact am observat că modelul era antrenat într-un mod mult prea restrictiv, eliminând propoziții cu caracter esențial în situații în care acest lucru nu era necesar. Motivele pentru acest comportament rezultă din felul în care corpusul a fost adnotat și din faptul că trăsăturile cu tf-idf extrase nu au capacitatea de a generaliza pentru date cu varietate mare lexicală. Altfel spus, datele de antrenare și datele pentru care am executat filtrare, au lexic diferit fapt care duce la overfitting. Pentru a compensa asta, am antrenat structuri de reprezentare vectorială a propozițiilor folosind centroidul de word embeddings. Antrenarea clasificatorului sau regresorului s-a executat pe centroizii de reprezentări vectoriale iar rezultatele s-au dovedit



promițătoare (F1 0.714, 0.73 acuratețe). Analizând manual adnotările făcute de clasificator, am putut observa că acestea erau mult prea permissive, adnotând drept conținut fraze care ar fi făcut în mod normal parte din clasa celor eliminate. Așadar am ales să agregăm modelul care execută predicțiile pe baza reprezentărilor tf-idf cu cel care utilizează reprezentările vectoriale într-un model combinat. Tabelul următor conține un rezumat al rezultatelor obținute în identificarea și extragerea conținutului esențial.

Method	$F_1$	Accuracy
tf-idf classifier	0.746	0.890
emb classifier	0.714	0.873
tf-idf context proba	0.775	0.897
emb context proba	0.738	0.878
combined	0.774	0.893

Drept concluzii, putem vedea că utilizarea probabilităților contextuale ale liniilor îmbunătățește acuratețea și metrica F1 de evaluare, dar că această metrică nu este neapărat suficientă în evaluarea exactă a frazelor care pot fi eliminate sau păstrate. În urma procesului de sampling și inspectare manuală a datelor adnotate, am putut observa diferențe semnificative în momentul folosirii celor două tipuri de clasificatoare într-o formă combinată, astfel că forma combinată oferă un echilibru mai bun între clasificatorul prea restrictiv și cel prea permisiv. Considerăm că acest rezultat este deosebit de important pentru agentul economic, permițându-i acestuia să extragă conținutul esențial care ulterior poate fi folosit pentru adnotare, căutare sau procesare.

## Activitati de diseminare

În aceasta etapa, membrii proiectului au reușit să publice 9 lucrări (3 lucrari la conferinte categoria A, 2 la conferinte categoria B, 1 la conferinte tip C, 2 lucrari la workshop-uri asociate unor conferinte tip A si o lucrare la granita dintre științele umaniste și cele reale) la unele din cele mai importante conferințe dedicate procesării limbajului natural și lingvisticii computaționale (COLING, EMNLP, CICLING, LREC, workshopuri asociate EMNLP, etc), si o serie de alte articole au fost trimise spre publicare.

Articolul *Lexical Analysis and Content Extraction from Customer-Agent Interactions*, autori: Sergiu Nisioi, Anca Bucur and Liviu P. Dinu, în curs de publicare (W-NUT la EMNLP 2018) în care sunt prezentate problemele cu care ne confruntăm în cadrul analizei de text pe date cu conținut specific, cum sunt cele de la Pluriva și o metodă de eliminare a redundanțelor din texte cu scopul de a păstra elementele esențiale din schimbul de replici dintre agenți și clienți.

Articolul *Exploring Optimism and Pessimism in Twitter Using Deep Learning*, în curs de publicare la EMNLP 2018 este printre primele articole din literatura care abordează într-o maniera computationala detectarea starii de optimism/pesimism din texte. Daca pentru analiza opiniilor si detectarea sentimentelor sunt peste 7000 de articole publicate, acest rezultat se numara printre primele articole (a fost considerat chiar primul articol din domeniu de catre unul dintre referentii anonimi) si consideram ca va fi extrem de util pentru o mare varietate de utilizatori. Un rezultat extrem de important a fost cel publicat in articolul *Ab Initio: Automatic*

*Latin Proto-word Reconstruction* (COLING 2018) in care reusim sa prezicem protocuvantul latin din care provin cuvintele in limbile romanice moderne cu o acuratete superioara rezultatelor raportate anterior, devenind astfel state of art. Tehnicile computationale utilizate sunt folosite in premiera, si metoda folosita are in vedere inclusiv predictia cuvintelor din care au provenit cuvintele romanesti moderne.

Un articol la intersecția dintre științele umaniste și cele reale, care abordează aspecte filozofice, etice și teoretice care țin de impactul tehnologiei în contemporaneitate, este *Posthumanism, Technology, and Monstrous Life Forms*, autor: Anca Bucur, publicat la ISEA 2018. Aceste aspecte devin relevante din punct de vedere social și economic, în momentul în care tehnologia capătă agentivitate și un potențial de înlocuire a forței de muncă umană. Un alt articol publicat este *A Detailed Evaluation of Neural Sequence-to-Sequence Models for In-domain and Cross-domain Text Simplification*, autori: S Stajner, S Nisioi, LREC 2018, care expune problemele și avantajele utilizării unei rețele neuronale de tipul celei implementate pentru generarea răspunsurilor, în cadrul unui task care are ca scop simplificarea conținutului și extragerea esențialului din textele originale. În baza rezultatelor obținute aici, am putut determina metoda optimă, conform opiniei noastre, de a esențialul din datele avute la dispoziție.

În articolul *Full Inflection Learning Using Deep Neural Networks* (CICLING 2018, in curs de publicare) au fost folosite rețelele neuronale în scopul modelării formelor flexionare ale substantivelor și verbelor din română (pornind de la forma de dictionar), și am obținut cele mai bune rezultate raportate până acum.

De asemenea, membrii proiectului au participat la stagii de cercetare, au susținut mai multe conferințe în țară și străinătate, au organizat a doua conferința internațională RAAI 2018 (Recent Advances in Artificial Intelligence), a fost continuată în cadrul proiectului seminarul internațional bilunar de lingvistică matematică și computațională „Solomon Marcus”, pe parcursul căruia a continuat seria de prezentări începută în 2016. În cadrul prezentărilor participă cu regularitate atât membrii proiectului, cât și masteranzi, doctoranzi sau postdoctoranzi din Departamentul de Informatică, din alte departamente ale Universității din București, din cadrul Facultății de Automatică a Universității Politehnica din București, membri ai partenerului Pluriva, dar și alți colegi din industrie interesați de tematica abordată. Considerăm că prin această activitate aducem un beneficiu proiectului prin activitatea de diseminare, și, în același timp realizăm noi punți de comunicare între mediul academic și industrie.

A fost întreținută pagina web a proiectului (<http://nlp.unibuc.ro/projects/sigar.html>), astfel încât toate rezultatele așteptate la finalul acestei etape (pagina web, articol trimis spre publicare, raport) au fost realizate.

### **Activități administrative și de organizare**

Membrii echipei s-au întâlnit periodic pentru a prezenta și discuta aspecte legate de organizare, de repartizarea temelor de cercetare, de colaborare între membrii echipei, de stabilirea conferințelor și jurnalelor în care se trimit spre publicare rezultatele cercetării și de convenirea asupra unui calendar comun cu respectarea termenelor de atingere a obiectivelor proiectului și de livrare a rezultatelor. Între membrii proiectului și partenerul economic a fost o continuă comunicare.

Au fost îndeplinite astfel toate activitățile administrative necesare bunei desfășurări a proiectului.

## Concluzii

- Obiectivele etapei au fost substanțial realizate.

### Publicații și conferințe susținute de către membrii proiectului în anul 2018:

1. Cornelia Caragea, Liviu P. Dinu, Bogdan Dumitru, **2018**. *Exploring Optimism and Pessimism in Twitter Using Deep Learning*. In Proc. **EMNLP 2018**, Brussels, Belgium, 2018 (to appear)
2. Sergiu Nisioi, Anca Bucur, Liviu P Dinu, **2018**. *Lexical Analysis and Content Extraction from Customer-Agent Interactions*. Proceedings of the 2018 EMNLP Workshop W-NUT: The 4th Workshop on Noisy User-generated Text, pages 132–136 Brussels, Belgium, 2018.
3. Alina Ciobanu, Liviu P. Dinu, **2018**. *Ab Initio: Latin Proto-word Reconstruction*. In Proc. **COLING 2018 (main conference)**, p. 1604-1614, Santa Fe, USA, 2018 .
4. Alina Ciobanu, Liviu P. Dinu, **2018**. *Simulating Language Evolution: A Tool for Historical Linguistics*. In Proc. **COLING 2018 (demo section)**, pp68-72, Santa Fe, USA, 2018
5. Anca Bucur, **2018**. Posthumanism, Technology, and Monstrous Life Forms. In Proceedings of ISEA 2018 - 24th International Symposium on Electronic Art, p 313-316, Durban, South Africa, 2018.
6. Bogdan Dumitru, Alina Ciobanu, Liviu P Dinu, **2018**. *ALB at SemEval-2018 Task 10: A System for Capturing Discriminative Attributes*. In Proc **SemEval 2018 (SemEval@NAACL-HLT 2018)**, task 10, p. 963-967, New Orleans, USA, 2018
7. Maria Sulea, Bogdan Dumitru, Liviu P. Dinu. *Full Inflection Learning Using Deep Neural Networks*, **CICLing 2018**, Hanoi, Vietnam, 2018
8. Liviu P Dinu, Ana Uban, **2018**. *Analyzing Stylistic Variation across Different Political Regimes*, **CICLing 2018**, Hanoi, Vietnam, 2018
9. [Sanja Stajner](#), Sergiu Nisioi, **2018**. *A Detailed Evaluation of Neural Sequence-to-Sequence Models for In-domain and Cross-domain Text Simplification*. Proceedings of 11th edition of the Language Resources and Evaluation Conference (**LREC 2018**), p. 3026-3033, 7-12 May 2018, Miyazaki (Japan) .

Director Proiect,  
Prof. dr. Liviu P. Dinu

*Lincoln*