

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

1. Platforma de prelucrare a textelor TEPROLIN

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

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

Păiș et al., (2020) descriu un studiu de caz în care TEPROLIN rulează pe mai multe fire de execuție în RELATE și adnotează corpusul legislativ din proiectul MARCELL. În total, au fost

adnotați aprox. 456 de milioane de tokeni, *ceea ce demonstrează că testarea și validarea platformei s-au încheiat cu succes.*

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

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

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

- a. `python3 -m venv /calea/către/mediul/dedicat/teprolin`
- b. `pip3 install -r requirements.txt`

2.1 Download și acces public

Versiunea publică (și finală) a platformei TEPROLIN se află pe GitHub, la adresa <https://github.com/racai-ai/TEPROLIN>. O versiune anterioară mai veche, care conține software proprietar, se află la <https://gitlab.com/raduion/teprolin>. Pentru a avea acces la această versiune, trebuie să aveți cont pe GitLab și să solicitați accesul autorului platformei.

3. Referințe

Păiș, V., Tufiș, D. și Ion, R. (2020) A Processing Platform Relating Data and Tools for Romanian Language. In Proceedings of the 12th Language Resources and Evaluation Conference (LREC), European Language Resources Association, Marseille, France, pages 81—88.