

Teme de licență 2016
Dan Cristea

1. Parametrizarea Portalului COROLA

în cotuletă cu lector dr. Alex Moruz (alex.moruz@gmail.com), drd. Andrei Scutelnicu (andreiscutelnicu@gmail.com)

COROLA este un acronim de la *CORpus of ROmanian LAnguage*. Numele complet al proiectului este Corpus Computațional Reprezentativ al Limbii Române Contemporane, un proiect prioritar al Academiei Române, în curs de realizare la Institutul de Inteligență Artificială București și Institutul de Informatică Teoretică din Iași. Portalul proiectului poate fi accesat la adresa: <http://78.96.45.245/>. Actualmente Portalul găzduiește tehnologii care realizează următoarele operațiuni:

- curățirea textelor primite de la colaboratorii noștri (edituri, autori etc.) de formate de tipărire - proces semiautomat;
- completarea metadatelor ce trebuie să însoțească documentele (autor, editură, an apariție etc.) - proces semiautomat;
- lansarea unui lanț de prelucrări lingvistice care urmăresc adnotarea automată cu informații de segmentare și de natură morfologică: granițe de propoziții și cuvinte, părți de vorbire și leme. Ulterior, după perfecționarea instrumentelor corespunzătoare, se urmărește adăugarea și a altor niveluri de analiză: sintactică, semantică etc.

Activitățile manuale se derulează prin intermediul unei interfețe web.

Proiectul de licență urmărește parametrizarea activităților Portalului pentru a separa procedurile aplicabile oricărei limbi de cele specifice unei limbi anume. În felul acesta, adaptarea lui pentru o anumită limbă ar implica particularizări și accesul la resurse specifice (lanțuri specifice de prelucrare etc.).

În esență se va dezvolta un *framework* de creare a corpusurilor, care poate interesa foarte mulți cercetători preocupați de construirea de corpusuri pentru limbile care nu dispun încă de astfel de resurse.

2. Monitorizarea activităților în Portalul COROLA

Se urmărește monitorizarea dinamicii corpusului COROLA. În particular, complexitatea deosebită a acestui corpus, datorată sutelor de mii ori milioane de fișiere de text și voce care vor trebui să-l compună și a adnotărilor adăugate lor, parțial automat, parțial manual, precum și dinamica achiziționării lui, care se va desfășura pe parcursul a mai multor ani, impune o monitorizare atentă a evoluției lui. IIT a elaborat deja (Cristea et al., 2014) un graf de dependențe între schemele de adnotare care să facă posibilă asignarea unui unic identificator (corespunzător nodului din graf ce caracterizează adnotarea pe care o include) fiecărui fișier aflat într-o anumită fază de adnotare. Simultan, se urmărește ca fișierele să aibă asociate în

partea lor de metadate, între altele, și schema de codificare, tipul adnotărilor (manuală sau automată), instrumentele de adnotare folosite, versiunile lor etc. Existența grafului de dependențe între diferitele scheme de adnotare permite regăsirea și complementarea informației din metadate referitoare la adnotarea unui anumit text. Va exista o legătură directă între metadate, graful de dependențe și schemele de adnotare. Această viziune va face posibilă:

- 1) identificarea adnotărilor contribuite de experți manual față de echivalentul lor inclus automat de tehnologie;
- 2) urmărirea lanțului de adnotări adăugate unui fișier, pe măsura rulării acestuia și continuarea în cazul întreruperilor;
- 3) calcularea minimului de elemente *stand-off* necesar de refăcut în cazul înlocuirii unui modul cu altul perfecționat;
- 4) elaborarea unor indicatori globali și unor statistici care să oglindească stadiul evoluției corpusului,
- 5) definirea restricțiilor asupra interogărilor KWIC (*key word in context*) posibil de adresat fiecărui fișier în funcție de adnotările pe care le include etc.

3. Formularea interogărilor GGS în limbaj natural

În cotutelă cu drd. Radu Simionescu (radsimu@gmail.com)

GGs-3 (*Graphical Grammar Studio*) este un sistem cadru de dezvoltare a grafurilor de tip rețele tranzitive recursive (*recursive transitive networks* - RTN), prin care se pot descrie reguli simbolice de prelucrare a limbajului natural. Puterea lor de prelucrare depășește cu mult cea a expresiilor regulate, pentru că limbajul GGS-urilor acceptă definirea de variabile pentru memorarea temporară de valori, parsarea în avans sau retroactivă fără consumarea tokenilor de intrare, definirea de macro-uri etc. GGS-3 a fost aplicat până acum cu succes la definirea de reguli pentru corectarea erorilor adnotărilor la parte de vorbire lăsate de un POS-tagger statistic, la antrenarea de parsere sintactice și la generarea de gramatici sintactice (Simionescu, 2014, 2015, 2016).

Datorită flexibilității și puterii lui de exprimare, se speră ca GGS-3 să poată fi utilizat la formularea de constrângeri pentru regăsirea de contexte de apariție ale cuvintelor în corpusul COROLA (v. mai sus).

Proiectul urmărește realizarea unei interfețe în limbaj natural care să preia exprimări simple de condiții și să genereze secvențe GGS-3, care, adresate corpusului COROLA (sau oricărui altui corpus care respectă același format), să genereze liste de contexte (co-ocurențe) care respectă condițiile.

Exemple de astfel de interogări KWIC (*key word in context*):

- *Vreau ocurențele lemei 'mișca' (vb)*. => se va genera o listă de ocurențe centrate pe formele verbului 'a mișca', cu contexte stânga-dreapta care includ, implicit, doar fraza în care apare verbul;

- 'sufla' urmat de 'în' (contexte dreapta în lungime de 5 cuvinte) => se va genera o lista de ocurențe care încep cu o formă a verbului 'a sufla' și continuă de prepoziția 'în' și cel mult încă 3 cuvinte;
- *exemple ale verbului 'a duce' în care apare un complement circumstanțial instrumental* => se va genera o listă de contexte frazale în care verbul 'a duce' este urmat de un complement circumstanțial instrumental (ex: "Ieri, Ion a fost atât de amabil să mă ducă cu mașina la gară.").

4. Inferarea gramaticii unei limbi

Proiect propus de drd. Radu Simionescu (radsimu@gmail.com) - luați legătura direct cu dânsul.

5. Recunoașterea relațiilor semantice în texte

Ariton Andrei - aryton_andrey19@yahoo.com, andrei.ariton@info.uaic.ro - pt relații de rudenie
 Andreia Băiceanu - andrea_baiceanu@yahoo.com, mihaela.baiceanu@info.uaic.ro - pt relații afective și sociale

Tema urmărește antrenarea unui program care să recunoască diferite tipuri de relații semantice care sunt exprimate într-un text liber între mențiuni de entități. O entitate poate fi: **persoană** sau **parte fizică a unei persoane, instituție, locație geografică** etc. Exemple de mențiuni de entități de tip **persoană**: *Ion, el, bărbatul cu pălărie de soare, sa*; **părți fizice de persoane**: *mâna sa dreaptă, ochii, piciorului* etc.; **instituții**: *Universitatea "Alexandru Ioan Cuza", Primărie, PNL* etc.; locații geografice: *Iași, Bulevardul Copou, Palas Mall* etc.

Iată și câteva exemple de relații:

- **referențialitate** de tip **part-of**: X part-of Y dacă X este o parte fizică a persoanei Y

1:[*Ion*] și-a acoperit 2:[*ochii*] cu 3:[*mâna*]. => [2] part-of [1], [3] part-of [1]

- diferite tipuri de rudenie (**kinship**): X **daughter** Y dacă X este fiica lui Y, X **parent** Y dacă X este părinte pentru Y etc.

1:[*fiica lui*] 2:[*Zamolxis*] => [1] daughter-of [2]

Când 1:[*ți*]-ai vizitat 2:[*părinții*] ultima oară? => [2] parent [1]

- relații spațiale:

Când urci 1:[*Bulevardul Copou*], *pe stânga ai să întâlnești clădirea* 2:[*Universității*]. => stabiliți voi ce relație se poate afirma între entitățile [1] și [2]?...

În antrenare se vor folosi două corpusuri care conțin notații ale entităților și relațiilor, construite manual: QuoVadis și MappingBooks, ambele amintite la curs. Sarcina voastră este de a crea patternuri de recunoaștere a relațiilor, de a găsi trăsăturile cele mai relevante pentru recunoașterea relațiilor și de a antrena pachete de statistică (Weka etc.) în acest scop.

6. Deducerea abilităților de manevră în jocuri interactive

proiect propus de studentul Petru Manea