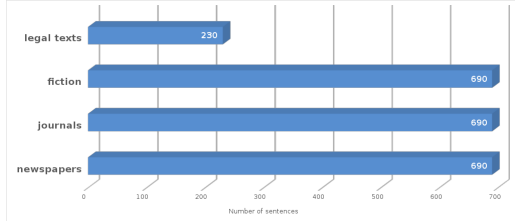


The Representation of MWEs in the Lithuanian Dependency Treebank

Jolanta Kovalevskaitė, Erika Rimkutė, Loïc Boizou

Lithuanian Dependency Treebank (LDT) *Alksnis*

- LDT is a part of Clarin-LT infrastructure.
- The set period of working on LDT covers 2015-2016.
- The goal is to prepare 2300 sentences annotated according to the dependency grammar.
- The corpus itself consists of several text types:



LDT Annotations

Morphology

- According to MULTEXT-East format.
- Each part of speech is annotated using an individual set of morphological categories (from 2 to 14). Examples:
 - turi* ('he/she has'), lemma *turėti* ('to have'), annotation Vgmp3s--n--ni- (verb, general, main form, present tense, 3rd person, singular, -gender, -voice, not negative, -definiteness, -case, not reflexive, indicative mood, -degree);
 - vertinimų* ('evaluations'), lemma *vertinimas* ('evaluation'), annotation Ncmppgn- (noun, common, masculine, plural, genitive, -name).

Syntax

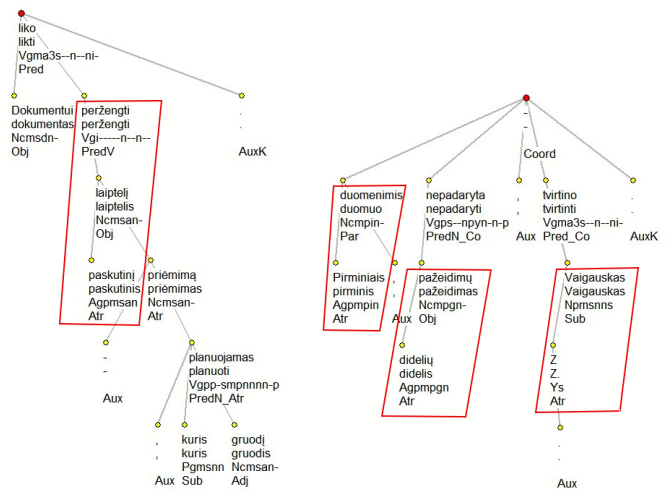
- Dependency model adapted from the Prague Dependency Treebank analytical layer.
- There are 18 syntactic functions (some of them combined, e.g. Pred_Atr, Pred_Co).
 - Pred (simple predicate)
 - PredN (compound nominal predicate)
 - PredV (compound verbal predicate)
 - Sub (subject)
 - Obj (object)
 - Adj (adjunct)
 - Atr (attribute)
 - Coord (coordination)
 - AuxC (subordinate conjunction)
 - AuxP (preposition)
 - AuxZ (emphasizing word)
 - Aux, AuxK (other auxiliary)
- The syntactic analysis is produced by a rule-based parser (Boizou and Zamblera 2014).
- Both morphological and syntactic annotations are performed automatically, but revised manually by linguists.

Types of MWEs in Lithuanian

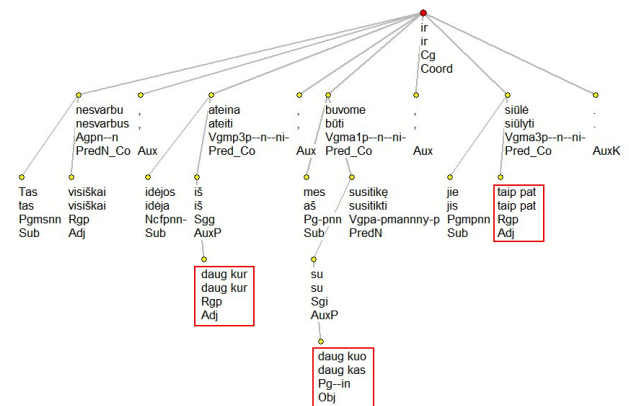
MWE type	Definition	Examples
Collocations	arbitrary, analyzable and flexible MWEs	nominal collocation: <i>atsakingas sprendimas</i> ('responsible decision') <i>pirminiais duomenimis</i> ('preliminarily preliminary:INS.PL data:INS.PL') verbal collocation: <i>priimti sprendimą</i> ('to make a decision')
Idioms	semantically non-compositional MWEs	nominal idiom: <i>Achilo kulnas</i> ('Achilles heel') <i>vargais negalais</i> ('with difficulty') verbal idiom: <i>kasti karo kirvį</i> ('to dig the hatchet')
Proverbs	syntactically complete expressions with no slots to be filled	<i>kas ne su mumis — tas prieš mus</i> ('he who is not with us is against us')
Named entities	Geographical names, names of persons, companies, institutions...	<i>Valdas Adamkus</i> <i>Kauno rajonas</i> <i>Via Baltica</i>
MWEs of grammatical nature	consist of two or more words (composed of inflective or uninflected parts of speech) and form semantically and syntactically unified, non-compositional unit that performs one syntactic function	multi-word adverbs: <i>taip pat</i> ('also, too') <i>iš anksto</i> ('in advance') multi-word pronouns: <i>kai kurie</i> ('some'), <i>nė vienas</i> ('none') multi-word particles: <i>vargu ar</i> ('hardly') multi-word prepositions: <i>iki pat</i> ('to, until') multi-word conjunctions: <i>vis dėlto</i> ('however, nevertheless')

The Representation of MWEs Types in LDT

In LDT, we annotate all words of different MWE types separately, except for those of grammatical nature. The formal flexibility of the Lithuanian MWEs (mostly of collocations, idioms and named entities) and the rather free word order are important reasons to treat each word of these MWEs as a single syntactic node with its proper morphological and syntactic annotation.



MWEs of grammatical nature are treated as single lexical units already on the morphological level, and appear as single nodes in tree structures.



Future Works

- The solutions for the representation of each type of MWEs are still under development.
- Possible options:
 - to apply more than one principle of annotation with respect to a particular MWE type;
 - to have a special label *MWE* for all MWEs.

References

- Boizou L., Kovalevskaitė J., Rimkutė E. Automatic Lemmatisation of Lithuanian MWEs. In *Proceedings of 20th Nordic Conference of Computational Linguistics NODALIDA 2015*. NEALT proceedings, vol. 23. Linköping, ACL anthology, 2015. <http://aclweb.org/anthology/W/W15/W15-1808.pdf>
- Boizou L., Zamblera F. Syntactic Engine for the Lithuanian Language. In *Proceedings of the Sixth International Conference Baltic HLT 2014: Human Language Technologies — The Baltic Perspective*. Amsterdam, Berlin, Tokyo, Washington, DC, IOS Press, 2014, 69–75. <http://ebooks.iospress.nl/publication/38006>
- Kovalevskaitė J., Boizou L., Rimkutė E. 2016: Lietuvių kalbos žodžių junginių morfologinių ir sintaksinių ypatybių sąsajos. *Darbai ir dienos* 64, 115–133.