# Automatic text processing and deeply annotated text corpora of Russian: interaction and mutual impact

*Leonid Iomdin*

A.A.Kharkevich Institute for Information Transmission Problems, Russian Academy of Sciences

*iomdin@iitp.ru*

# **Plan**

1. Annotated text corpora of Russian: what are they?

2. SynTagRus treebank: basic facts

3. SynTagRus and ETAP parser: living together
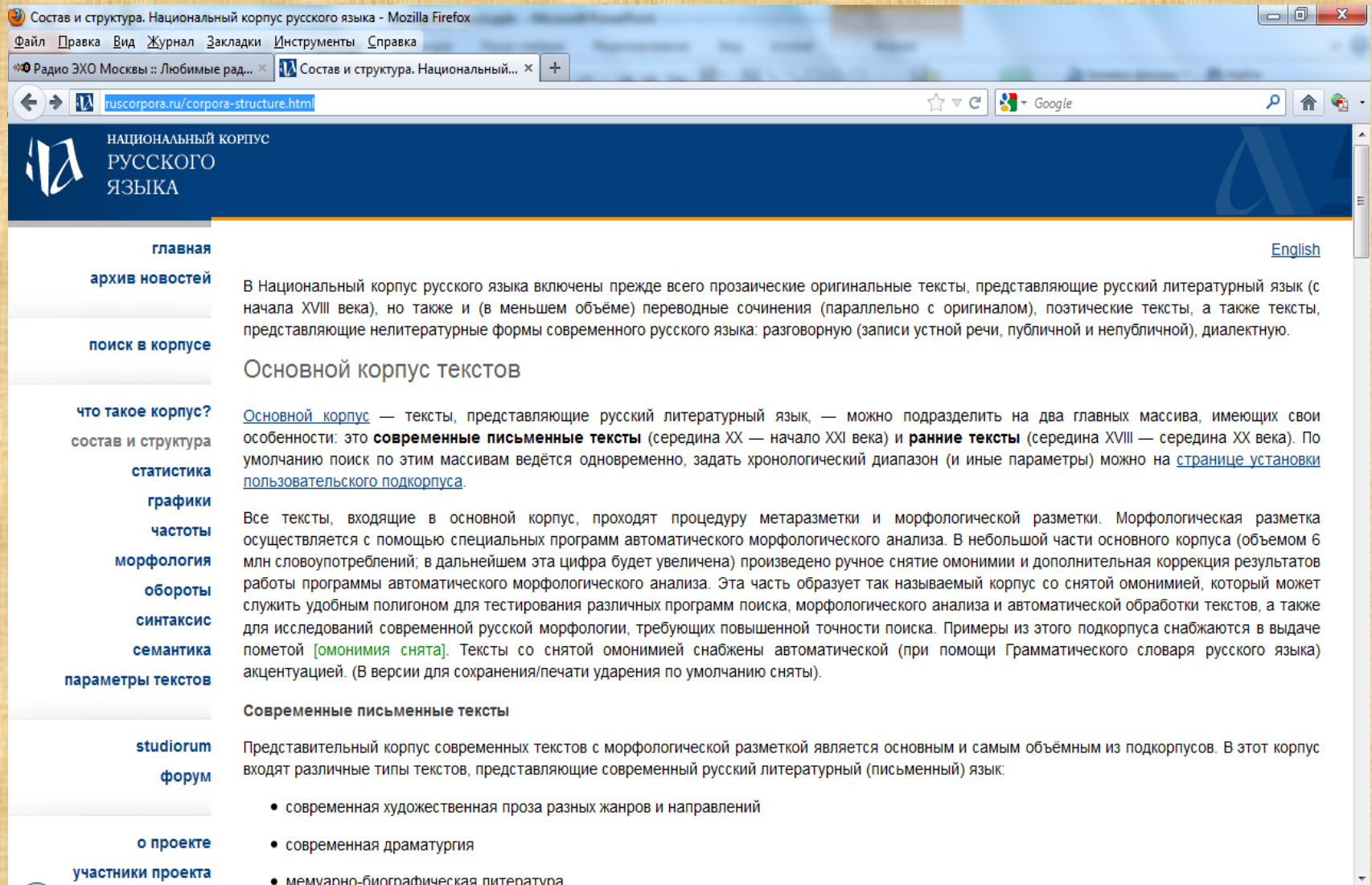
# Summary

The main focus will be on SynTagRus, a corpus of Russian texts annotated with dependency-type syntactic structures, lexical meanings, and lexical functions. Statistical data collected from the corpus are used to improve lexical and syntactic disambiguation of automatic parsing. Other uses of SynTagRus include construction of dependency parsers by machine learning techniques and regression testing of the rule-based parser.

# THE Russian corpus: National Corpus of Russian, НКРЯ

**www.ruscorpora.ru**
(hosted by Yandex, to be changed soon, non-commercial partnership to be founded)

# National Corpus of Russian

Bratislava: Slovak corpus
anniversary conference

# National Corpus of Russian

1. **Main corpus (morphological annotation), over 300 million words. Subcorpus with resolved ambiguity: ("snjatnik"): 7 million words**
2. **Syntactic corpus (morphological and syntactic annotation, lexical meanings and lexical functions): 770,000 words, over 52,000 sentences)**

# National Corpus of Russian

3. **Newspaper corpus**
4. **Parallel aligned corpora** (English-Russian, Russian-English, German-Russian, Ukrainian-Russian, Russian-Ukrainian, Belarussian-Russian, Russian-Belarussian, multilingual)

# National Corpus of Russian

5. **Dialectal corpus**
6. **Poetic corpus**
7. **Learners' corpus**
8. **Oral speech corpus**
9. **Accentological corpus**
10. **Multimedia corpus**
11. **Church Slavonic corpus**

# SynTagRus

Syntactic Corpus of NRC and SynTagRus: the former is a subcorpus of the latter (is updated 3-4 times a a year, does not show lexical meanings or lexical functional annotation).

Bratislava: Slovak corpus anniversary conference

# SynTagRus

The corpus is created semi-automatically: first, every sentence is processed by the ETAP parser and then manually corrected by at least two linguist experts.

# SynTagRus

**Currently the treebank contains over 52,000 sentences belonging to texts of a variety of genres (contemporary fiction, popular science, newspaper and journal articles dated between 1960 and 2012, texts of online news Wikipedia articles etc.) and is steadily growing.**

# SynTagRus

**SynTagRus adopts a dependency-based annotation scheme, in a way parallel to the Prague Dependency Treebank but, in contrast, relying on the Meaning – Text theory by Igor Mel'čuk.**

# SynTagRus

**A sentence:**
*Naibol'šee vozmuščenie učastnikov mitinga vyzval prodolžajuščijsja rost cen na benzin, ustanavlivaemyx neftjanymi kompanijami*
'It was the continuing growth of petrol prices set by oil companies that caused the greatest indignation of the participants of the meeting'.

# SynTagRus

## and its representation:

# SynTagRus

**Nodes represent words (lemmas) assigned morphological and part-of-speech tags, whilst arcs are labeled with names of syntactic links. The tagging uses about 75 syntactic links, half of them proposed by Mel'čuk (1988).**

# SynTagRus

**Normally, one token corresponds to one node in the dependency tree. There are certain exceptions:**

- **composite words like _pjatidesjatiètažnyj_ 'fifty-storeyed', where one token corresponds to two or more nodes;**
- **multiword expressions like _po krajnej mere_ 'at least' where several tokens correspond to one node;**

# SynTagRus

§ **so-called phantom nodes for the representation of hard cases of ellipsis, or gapping, which do not correspond to any particular token in the sentence (cf. *Ja kupil rubašku, a on galstuk* 'I bought a shirt and he a tie'), which is expanded into *Ja kupil rubašku, a on kupilPHANTOM galstuk* 'I bought a shirt and he boughtPHANTOM a tie'**

# SynTagRus and ETAP

**Morphological Tagging of SynTagRus is based on a comprehensive morphological dictionary of Russian that counts about 130,000 entries (over 4 million word forms).**

**Recently, the dictionary was supplemented with full phonetic stress marking, which is used in a Russian text-to-speech synthesis system.**

**ETAP-3 morphological analyzer uses the dictionary to produce morphological annotation of words belonging to the corpus, which includes the lemma, POS tags, and, depending on POS, a set of morphological features.**

# SynTagRus and ETAP

**The current version of SynTagRus contains partial lexical functional annotation. For collocations that could be presented with the apparatus of lexical functions, the tagging includes information on values and attributes of such lexical functions.**

Bratislava: Slovak corpus anniversary conference

# SynTagRus and ETAP

# SynTagRus and ETAP

# SynTagRus and ETAP

**The current version of SynTagRus displays word senses as they are presented in the combinatorial dictionary of ETAP.**

# SynTagRus and ETAP

*Пока же исследователи по-разному толкуют "первоисточник", а в качестве доказательства своей правоты приводят отдельные археологические находки последних лет.*

'**So far, the researchers interpret differently the primary source, and as a proof of their being right they present isolated archeological findings of the recent years**'

# SynTagRus and ETAP

# SynTagRus and ETAP

# SynTagRus and ETAP



**Dialog**

**CD Names**

ТОЛКОВАТЬ1

| ТОЛКОВАТЬ1 |
| ТОЛКОВАТЬ2 |

OK

Cancel

**CD text**

```
24576 20:16:10 11-03-2009        ТОЛКОВАТЬ1
        COMMENT:"ИСТОЛКОВЫВАТЬ, ИНТЕРПРЕТИРОВАТЬ"
        EXAMPLE:"СЛОВАРЬ ТОЛКУЕТ VIS-A-VIS КАК 'ЛИЦОМ К ЛИЦУ'"
        POR:V    ЛИ
        SYNT:НЕСОВ!,ТРАНЗИТ,СТР-ЕМ,СТР-СЯ,СТР,СОГЛАКТ-3-2,БЕНЕФ
        DES:'ДЕЙСТВИЕ','ФАКТ','АБСТРАКТ'
        D1.1:'ЛИЦО','МЕХАНИЗМ','СИСТЕМА','ИНФОРМАЦИЯ'
        D2.1:ВИН,'ИНФОРМАЦИЯ'
        D3.1:КАК1
        _S0:ТОЛКОВАНИЕ
TRAF:1-КОМПЛ.11
TRAF:1-КОМПЛ.51
TRAF:2-КОМПЛ.15
TRAF:2-КОМПЛ.16
TRAF:1-КОМПЛ.28
TRAF:1-КОМПЛ.27
TRAF:НЕАКТ-КОМПЛ.12
        ***************************
ZONE:EN
        TRANS:INTERPRET
TRAF:EXPANS.64
LR:ТОЛКОВАНИЕ
```
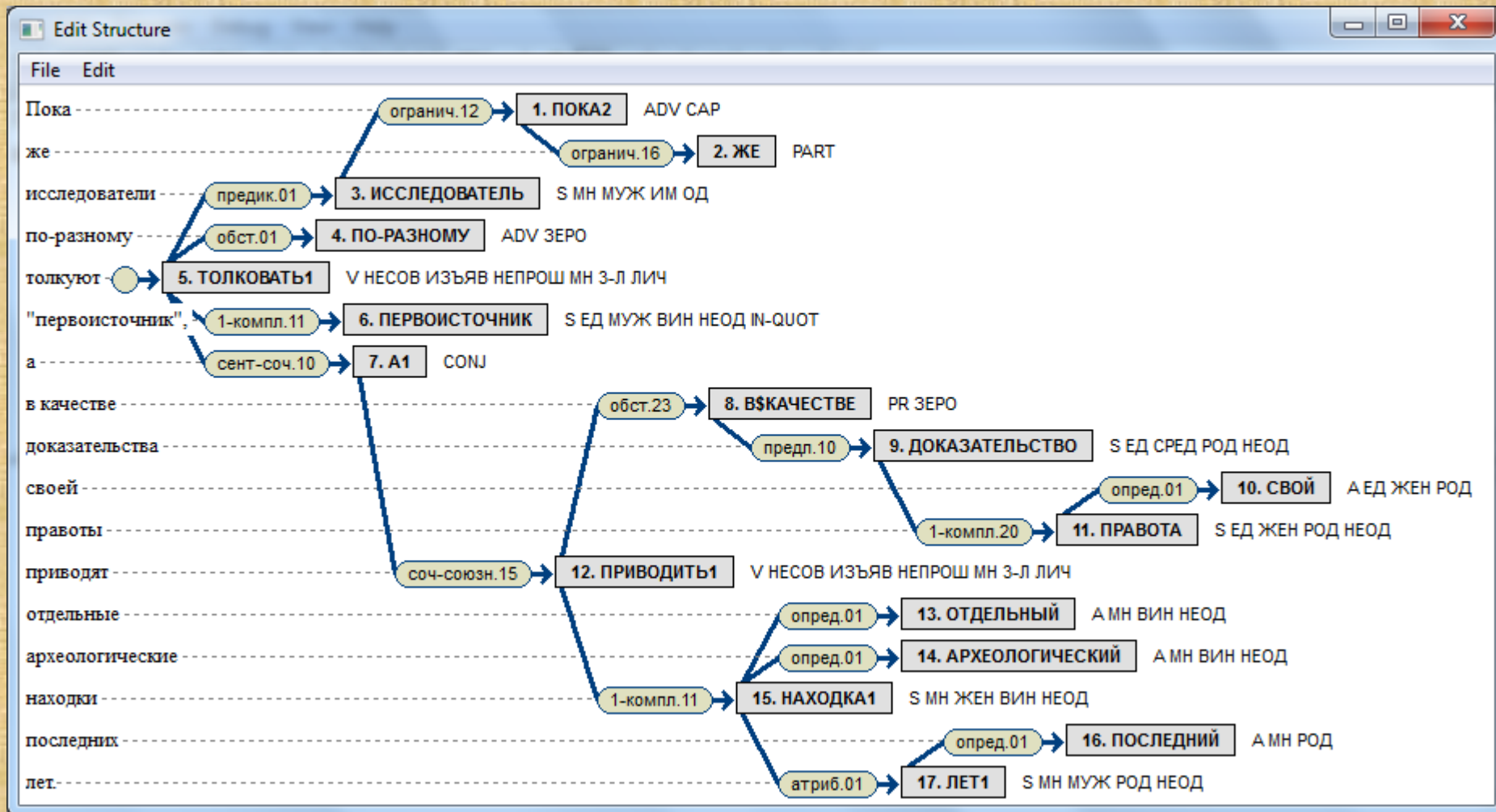
# SynTagRus and ETAP

*Толкуют о сооружении местного Сити, почти как в Москве, и, разумеется, с высоченным небоскребом напротив Смольного на невском правобережье.*

**'They talk about the construction of a local City, almost like in Moscow and, naturally, with a very high skyscraper opposite Smolny on the Neva river bank.**

# SynTagRus and ETAP

Bratislava: Slovak corpus anniversary conference

# SynTagRus and ETAP
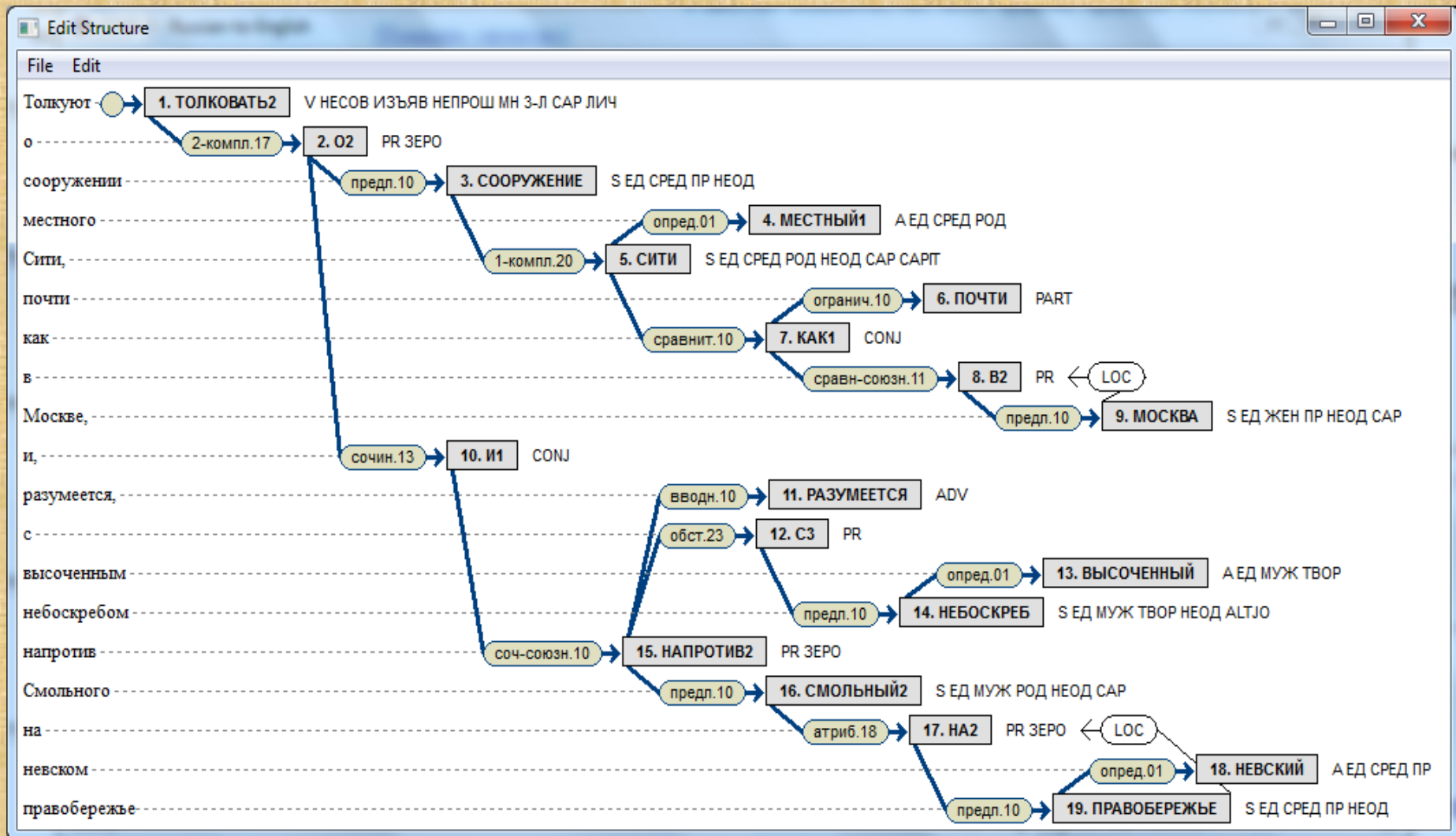
# SynTagRus and ETAP

Syntagrus is not only a linguistic resource but also a computational resource which is used

- to collect various statistical data;
- to create training sets for machine learning;
- to develop automatic parsers (Nivre-Boguslavsky-Iomdin 2008)

# SynTagRus beyond ETAP

- Over 30 free licenses provided to universities and academic institutions;
- 2 commercial licenses provided to big IT companies

# Three main uses of SynTagRus within ETAP

1. It provides the statistics of the different syntactic constructions, lexical co-occurrences, patterns of ambiguities etc., which is used at several points of the algorithm if the statistical component is activated.

Bratislava: Slovak corpus anniversary conference

# Three main uses of SynTagRus within ETAP

2. It serves as an efficient and accurate evaluation resource, which is used to evaluate the performance of ETAP parser and in this way find and resolve some of the system's bottlenecks.

# **Three main uses of SynTagRus within ETAP**

3. It is used for regression testing of ETAP.

# Regression Testing

- Periodically, ETAP is run on the whole corpus. Sentences that receive parses exactly equivalent to those stored in the corpus (between 30 and 35 percent of the bulk of the corpus) are selected as basis for regression testing.
- ETAP is then run on this test set to see if changes introduced in the dictionary, rules, or software affected the state of the test set.
- Regression testing has proven helpful in ensuring the stability of the parser and eventually improving it. Regression testing helps improve the SynTagRus itself: sometimes the discrepancies in parses detected by regression test runs point to erroneous annotation in the corpus, which is then corrected.