

ГРАММАТИЧЕСКИЙ СЛОВАРЬ

ДЛЯ АВТОМАТИЧЕСКОГО АНАЛИЗА ТЕКСТОВ XVIII—XIX в.:

ПЕРВЫЕ РЕЗУЛЬТАТЫ

А.Е. Поляков, НПБ им. К.Д. Ушинского РАО

С.О. Савчук, Д.В. Сичинава, Институт русского языка

им. В.В. Виноградова РАН

Мотивация



Отсутствие морфологического анализатора для исторических корпусов и электронных библиотек



Отсутствие анализатора для распознавания отсканированного печатного текста в орфографии XVIII—XIX вв.



Неудовлетворительные результаты поиска и низкое качество распознавания

Цели

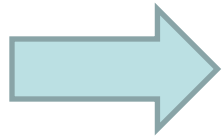
- Добиться удовлетворительных результатов при автоматическом анализе текстов следующих категорий.
- I. Тексты к. XVIII - н. XX в. в дореволюционной орфографии из оригинальных источников, не издававшиеся в XX в.
- *Среда. Маія 5 дня 1787 года. Усть не имѣя говорю: глась мой слушаютъ въ отдаленности и многіе оному повинуются. Иногда размѣряю и теченіе времени: иногда созываю людей на работу.*
- Основная проблема – правила графики и орфографии

- II. Тексты XVIII в. в современной и оригинальной орфографии.
- *Младыи* челоѡѡкъ всегда имѡетѡ съ *благочестными* и *добродѡтелными людми* *обходїтїся*, отѡ которыхѡ бы онѡ добру *научїтца* могѡ. Также и съ такіми *людми* которые и честное имя и непорочное *жїтїе* имѡютѡ. А отѡ такіхѡ, которые легкомысленно и *слочестно* живутѡ, бѡгати *яко* бы отѡ яду или лютого мору. [Юности честное зеркало, 1717]
- Основная проблема – пополнение словаря + исторические парадигмы + правила графики и орфографии

- III. Тексты XIV-XVII в. и рукописные тексты XVIII в.
- *И не точію тѣхъ иже градъ враждебно, яко зміеве, своими зубы держащихъ, носихъ множество и своевѣрныхъ враждующихъ ми страшахся, иже присѣдять о насъ тайно въ ловителѣхъ ко Еллиномъ* [Иван Тимофеев. Временник (1610-1617)]
- *и пажалои ты приежьцаи х празнику Никалаю чюдворцу а ка мнѣ въ госте а я млсте твоеи вседошна рать писавы Івашка Сидарав челомъ бьетъ* [И. Сидоров И.А. Снарскому (1701)]
- Основная проблема – пополнение словаря + исторические парадигмы + отсутствие норм орфографии

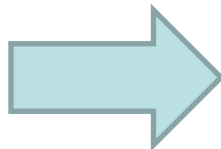
Задачи

- 1. *Грамматические*: разработка грамматических парадигм для лексем, отсутствующих в современном словаре.



- Грамматический словарь

- 2. *Орфографические*: обеспечение лемматизации форм, имеющих отклонения от стандартных написаний.



- Правила

Грамматический словарь

список лексем с приписанной информацией о словоизменении.

1) основа с указанием чередований

2) постоянные признаки лексемы (часть речи, род, одушевленность, переходность)

3) код словоизменительного типа (парадигмы)

б(и|ь|е)+ть V,ipf,tr V11

пе(к|ч|)+ь V,ipf,tr V8

но(с|ш)+ить V,ipf,tr V4

Модули грамматического словаря

Современный модуль

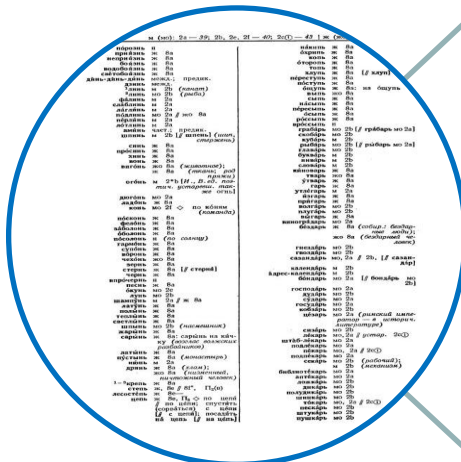
- Грамматический словарь Зализняка

Модуль XVIII—XIX века

- Исторические словари
- Корпус текстов

Модуль XVII века

- Корпус текстов
- Словарь русского языка XIV—XVII в



Формирование словника

Первое направление:
расширение словника
за счет включения в
него материалов
исторических словарей

- – Словарь Академии Российской (1789–1794);
- – Словарь церковнославянского и русского языка (ЦСРЯ) (1847);
- – Полный русский орфографический словарь (1898);
- – Словарь русского языка XVIII века.

Второе направление:
пополнение словника
словами, извлеченными
из корпуса текстов XVIII-
XIX в., которые не
учтены в лингвистических
источниках

- Несловарные слова:
- Имена собственные и образованные от них прилагательные.
- Лексемы, отсутствующие в словаре (*ажитёр, анатомить, бомбаст*)
- Морфологические архаизмы (*гресех, острови, знаеши, бяше, бяху* и пр.).
- Варианты различного рода.

Правила

А) Графические правила

- основаны на приравнивании графем, встретившихся в тексте, графемам, входящим в порождаемые словарём формы, например, $\theta \Rightarrow \phi$; $\xi \Rightarrow \kappa\sigma$.

Б) Орфографические правила

- связаны с нормализацией определённых орфограмм. Например, $\text{цы} \Rightarrow \text{ци}$: *цыновка, цыдулка, цыгарка, цыгейка*

В) Морфологические правила

- устроены с учётом информации о грамматическом разборе словарной словоформы. Например, $\text{-яе, comr, anom} \leq \text{-ее, comr}$.

Г) Списочный способ

- задание вариативности списками, на ограниченных классах единиц. Например, ряд конкретных корней в среднерусских и церковнославянских текстах записываться под титлом, ср. такие словоформы, как *Г(о)с(по)дь, м(е)с(я)ц, гл(агол)ет* и т. д.

Пример работы орфографического правила

в тексте
представлена
словоформа
надеютца;

из-за наличия
конечного *тца*
проверяется
правило
тца=ться/тся

словоформа
надеются
словарем не
порождается;

словоформа
надеются
словарем
порождается и
разбирается
как
lex=НАДЕЯТЬ
СЯ
gr=praes,3p,pl;

словоформа
надеютца
получает
разбор
lex=НАДЕЯТЬ
СЯ
gr=praes,3p,pl
=dis_{ort.}

Порядок применения правил

графические правила, в том числе списочные



орфографические



морфологические



списочные правила, блокирующие паразитические разборы

Циклы обработки текстов

В процессе формирования словаря и оптимизации работы анализатора предусматривается несколько циклов обработки текстов.



Оценка результатов

Экспериментальный корпус –
4 млн с/у, ок. 256 тыс. разных с/форм

Состав текстов	
художественные	24%
публицистика	24%
церковно-богословские	19%
научные	17%
официальные	11%
бытовые	5%
Хронологическое распределение	
1700-1730	6%
1731-1780	43%
1781-1799	30%
1800-1830	21%

Не распознаны или неправильно опознаны

сочетания знаменательных частей речи с **частицами** *-то(-та, -ат), же(ж), ли(ль), бы(б), -де, -ка: таковаже, пили-б, пилиб, еслиб.*

архаические формы склонения и спряжения: *возвышаяй, живяше, знаяше, стояху*

орфографические варианты: *безщчетну, возмеш, баталиах*

собственные имена — топонимы, имена лиц, литературных героев и мифологических персонажей: *Шлюссембурх, Елисаветфь, Невтон, Мишель Анжело*

разобрано	185221	72,4%
гипотезы	63904	25,0%
не разобрано	6780	2,6%

Оценка правильных разборов



**Высокий показатель
однозначных разборов (73%)**

В корпусе мало текстов начала XVIII в.
Основной частью текстов представлена в
современной орфографии

Увеличение доли оригинальных
текстов XVII и 1-ой трети XVIII в.
приведет к росту количества
ошибочных разборов

Ближайшие планы

внедрение всех изменений в словарь;

коррекция экспериментального корпуса с учетом выявленных ошибок;

тестирование новой версии словаря на экспериментальном корпусе;

тестирование словаря на отдельных текстах более раннего периода, относящихся к разным жанрам.

- Работа выполнена при поддержке Программы фундаментальных исследований Президиума РАН «Корпусная лингвистика».
- А.Е. Поляков pollex@mail.ru, НПБ им. К.Д. Ушинского РАО
- С.О. Савчук savsvetlana@mail.ru, Д.В. Сичинава mitrius@gmail.com Институт русского языка им. В.В. Виноградова РАН