

# ГРАММАТИЧЕСКИЙ СЛОВАРЬ ДЛЯ АВТОМАТИЧЕСКОГО АНАЛИЗА ТЕКСТОВ XVIII–XIX ВЕКА: ПЕРВЫЕ РЕЗУЛЬТАТЫ<sup>1</sup>

**Поляков А. Е.** (pollex@mail.ru)

НПБ им. К. Д. Ушинского РАО, Москва, Россия

**Савчук С. О.** (savsvetlana@mail.ru),

**Сичинава Д. В.** (mitrius@gmail.com)

Институт русского языка им. В. В. Виноградова РАН,  
Москва, Россия

В статье излагаются основные принципы построения грамматического словаря и морфологического анализатора для текстов XVIII–XIX вв. веков с учётом орфографических, морфологических и лексических особенностей языка этого периода, выявленных на материале текстов Национального корпуса русского языка. Поскольку от данного анализатора требуется универсальность подхода и возможность работы с текстами разных типов и разными орфографическими режимами, то он должен состоять из нескольких модулей, применяемых к текстам различных типов в зависимости от степени проявления в них тех или иных орфографических и грамматических явлений. Его словарь построен на базе существующего грамматического словаря современного русского языка, а также словарей XIX в. и текстов Национального корпуса. Обсуждается несколько альтернативных возможностей реализации орфографических (предобработка, применение технологии параллельного корпуса, нормализация при разметке) и морфологических правил (нормализация при разметке, добавление нестандартных форм в парадигму). Проводится оценка первых результатов применения анализатора к текстам НКРЯ и предлагаются различные варианты улучшения результатов (введение новых правил, пополнение словаря и т. д.).

**Ключевые слова:** грамматический словарь, автоматический анализ текстов XVIII–XIX вв.

---

<sup>1</sup> Работа выполнена при поддержке Программы фундаментальных исследований Президиума РАН «Корпусная лингвистика».

# A GRAMMAR DICTIONARY FOR AUTOMATIC ANALYSIS OF THE XVIII–XIX<sup>TH</sup> CENTURY TEXTS: FIRST RESULTS

**Polyakov A. E.** (pollex@mail.ru)

Ushinsky's State Scientific Pedagogical Library RAE

**Savchuk S. O.** (savsvetlana@mail.ru),

**Sitchinava D. V.** (mitrius@gmail.com)

V. V. Vinogradov's Institute for the Russian Language RAS,  
Moscow, Russia

The paper presents the key principles of building a grammar dictionary and a morphological analyzer for XVIII–XIX<sup>th</sup> century Russian texts based on orthographical, morphological and lexical features exemplified by the Russian National Corpus (RNC). The analyzer should involve different modules applicable to different kinds of texts depending on their respective orthographical and grammatical phenomena. Several alternative ways of implementing orthographical and morphological rules are discussed (including pre-processing, online normalization etc.). Evaluation data of the first analysis results are presented.

**Key words:** grammar dictionary, automatic text analysis, texts of the XVIII–XIX<sup>th</sup> century

## 1. Постановка проблемы и способы ее решения

Последнее десятилетие характеризуется ростом интереса к сохранению письменного наследия — текстов предшествующих эпох, в частности, к оцифровке этих текстов и увеличению возможностей доступа и поиска (играет здесь роль и то, что старые тексты, созданные, как минимум, век назад, не охраняются авторским правом, находясь в общественном достоянии). Наблюдается рост числа исторических корпусов, электронных исторических библиотек; объёмы оцифрованных старых изданий, доступных в электронном виде, существенно выросли. Нужно упомянуть, в частности, систему Google Books, где возможен полнотекстовый поиск на разных языках, в том числе русском, по книгам и журналам, для раннего периода нередко представленным в полном виде и доступным для скачивания. Этот поиск уже становится неотъемлемым инструментом работы специалистов самого разного профиля, занимающихся историей тех или иных явлений. Уровень технологий вырос и позволяет изготавливать электронные издания в графических форматах на достаточно высоком уровне и в большем объёме.

Однако основная проблема создателей исторических корпусов и электронных библиотек с поиском по тексту по-прежнему остается открытой. Для морфологического анализа текстов предшествующих эпох (если такая задача вообще ставится) используется анализатор, рассчитанный на современный язык и частично (или даже полностью) на современную орфографию, потому что другого пока нет. Более того, от качества работы такого анализатора зависит уже качество распознавания отсканированного печатного текста в орфографии XVIII–XIX вв.; словоформы, не предсказываемые современным грамматическим словарём, могут быть распознаны неправильно (например, слово *пыль* с конечным ером — как *пыль*, а *красныя* — как *красная* или *красный*).

Морфологическая разметка в Национальном корпусе русского языка производится (полу)автоматически с помощью анализаторов Dialing (подкорпус со снятой омонимией) и Mystem (тексты с неснятой омонимией в составе разных подкорпусов). В основе обоих анализаторов лежит грамматический словарь современного русского языка [Зализняк 1977/2003]; словарь анализатора Mystem, по сравнению со словарём Зализняка, пополнен на материале часто встречающихся в Интернете и в поисковых запросах неологизмов 1990–2000-х годов и имён собственных. Результат такой разметки дает значительное количество ошибочных и/или гипотетических разборов. Как было показано в [Савчук, Сичинава 2009], процент погрешностей в разборе текстов XVIII в. находятся на уровне современных текстов, не нормированных (тексты электронной коммуникации) или в принципе не ориентированных на литературную норму (записей диалектной речи).

Повышение качества анализа текстов с большим количеством отклонений от стандарта предлагалось проводить двумя путями: 1) использование стандартного анализатора, отражающего грамматические и орфографические нормы современного литературного языка, с подключением дополнительных правил для отдельных словоформ и категорий слов и 2) использование нового морфологического анализатора, настроенного на определенные тексты. Первый путь до сих пор использовался для улучшения качества разметки НКРЯ, однако в большом пополняющемся корпусе список словоформ, отклоняющихся от стандарта и требующих применения дополнительных правил, неуклонно растет, что делает этот способ неэффективным.

Между тем база текстов НКРЯ существенно растёт. Расширился объём текстов XVIII в. (к февралю 2012 г. 3,8 млн слов) и текстов предшествующего периода — среднерусских текстов XIV–XVII веков (3 млн слов). Увеличение объема исторических текстов, а также их хронологического разнообразия делает более актуальным использование нового анализатора. Морфологический анализатор церковнославянского языка, разрабатываемый для церковнославянского подкорпуса НКРЯ [Поляков и др., 2012], для этой цели в текущем виде не подходит (его пробное применение к среднерусским текстам, без орфографической настройки, показало, что опознаётся лишь примерно половина словоформ). Он рассчитан на современную церковнославянскую орфографию (выработавшуюся с XVII в. и отличную как от предшествующей церковной, так и от последующей гражданской: с последовательным сложным распределением

пар омофоничных букв и т. п.), а также морфологический и лексический стандарт богослужбных текстов последних веков, когда церковнославянский язык ограничивается собственно функцией языка православной литургии. Между тем тексты XIV–XVII в., не говоря уже о XVIII–XIX вв., используют церковнославянские элементы лишь в комбинации с собственно русскими, и в ряде жанров собственно русские черты на всех уровнях существенно преобладают.

Возможно применить к текстам одновременно церковнославянский анализатор (ослабив в нём требования к орфографии) и современный анализатор, наоборот, добавив в него по крайней мере правила, упразднённые реформой 1918 года. Такая комбинация церковнославянского и современного анализаторов дает сравнительно неплохие результаты, но оставляет неразобранной некнижную и нецерковную лексику и грамматику, отсутствующую как в современном русском языке, так и в современном церковнославянском (например, компаратив типа *сильные* или множественное число среднего рода типа *злодействы*), не говоря уже о словообразовательной и фонетической вариативности и проблеме совмещения разборов обоих анализаторов для совпадающих словоформ. Поэтому разработка нового анализатора должна вестись с учетом специфики конкретных текстов, при этом предусматривая появление в новых текстах новых аналогичных словоформ, для которых должна быть предусмотрена возможность правильного разбора.

## 2. Специфика корпуса исторических текстов

Корпус исторических текстов отражает эпоху, когда корреляция между языком и жанром была гораздо сильнее и прямолинейнее, чем в современной языковой ситуации. Общеизвестно принятое русским классицизмом разделение литературного языка на три «штиля» (высокий, средний и низкий), определенных, в том числе, по насыщенности текста славянизмами; каждый из них был закреплён за своим жанром. Для среднерусского периода было характерно сосуществование текстов с более сильными церковнославянскими либо народными тенденциями; представление о «диглоссии» — жёстком распределении жанров между двумя языками — несколько упрощено, хотя и отражает определённые установки писавших. Если верно, что церковно-богословские тексты, такие, как проповедь, ориентировались (в период до 1740-х гг.) на чистый церковнославянский язык, а бытовые тексты, такие, как «грамотки» XVII–XVIII в. или личные дневники более позднего времени — на разговорный русский, то столь же верно и то, что одновременно определённые народные языковые черты проникали и в первые, а церковнославянские — и во вторые. Более равноправное комбинирование славянских и народных элементов было характерно для языка летописей, посланий (так называемый «гибридный язык»), позже — официальных документов, дипломатических и учебно-научных текстов. Наконец, необходимо отметить и такое явление, характерное и для среднерусской эпохи, и для языка Нового времени, как церковнославянские цитаты (из Библии, тогда еще не переведённой на русский, или из богослужбных текстов, до настоящего времени в общепринятой практике

церковнославянских) в составе русских текстов. Это явление присутствует и в современном языке как в виде фразеологизмов-славянизмов, в том числе библеизмов (*ничтоже сумняшеся, возвращается ветер на круги своя, коемуждо по делом его* и т.д.), так и собственно цитат, причём не только в специально богословских текстах. Явлением более или менее тесного взаимопроникновения двух систем, их «гибридизации», собственно говоря, и объясняется необходимость учета элементов лексики и грамматики фактически другого языка при анализе русских текстов.

Тексты разных периодов могут значительно отличаться по характеру языка, поскольку конец XVII—середина XIX в. — это самый интенсивный период формирования литературного языка нового типа. По количеству заимствованной лексики, по употребительности славянизмов, по орфографической нормированности достаточно существенно различаются между собой петровская эпоха, период классицизма, эпохи, связанные с именами Карамзина и Пушкина.

До недавнего времени Национальный корпус русского языка включал ранние тексты (XVIII — первая половина XX в.; о среднерусском периоде речь не идёт) только в современной орфографии, или, по крайней мере, с некоторыми отклонениями от современной нормы, но в целом ориентированные на орфографический режим 1918 г. или даже 1956 г. [Савчук, Сичинава, Гарипов 2006]. Это диктовалось, помимо традиции, и ориентацией на авторитетные научные издания советского и постсоветского времени. Вместе с тем давно обсуждается вопрос о включении в корпус значительного количества текстов в дореформенной орфографии [Соловьев, Ахтямов 2006; Савчук 2008], с сохранением упразднённых реформой 1918 года графем и орфографических правил, тем более что анализатор *Mystem* умеет учитывать большинство этих правил (использование *ѣ, ѓ, і, конечного њ, окончания -ыя, -ія, -аго, -яго*). Ценность оригинальной орфографии для филологического изучения текста сейчас не нуждается в особом аргументировании (см., в частности, работы М. И. Шапира). В настоящее время в экспериментальном порядке в основной корпус НКРЯ входит один текст целиком в дореформенной орфографии по изданию 1857 года — роман «Черная рада» Пантелеймона Кулиша, впервые распознанный и вычитанный для Корпуса с прижизненного издания; его изданий в новой русской орфографии, кажется, не существует, поскольку наиболее актуальной для читателя уже порядка ста лет считается авторская версия этого текста на украинском языке. Вообще для текстов, никогда не переиздававшихся в новой орфографии, сохранение в Корпусе дореформенного правописания кажется особо важной задачей; безвозвратная утрата этой информации после трудоёмкого сканирования и вычитки текста была бы совершенно нерациональна. Вместе с тем встаёт проблема обработки и совместного поиска по текстам разных типов. Например, для основного корпуса НКРЯ желательна возможность индексировать корпус так, чтобы при поиске точных форм можно было найти одновременно и дореформенное, и современное написание (например, чтобы по точному запросу словоформы *пъной* находилась бы и словоформа *пеной* в новоорфографических текстах).

### 3. Принцип работы анализатора. Составные элементы анализатора

Можно сформулировать требования, которым должен отвечать анализатор для обработки исторических текстов:

А) *Универсальность*: анализатор должен уметь работать с текстами, обладающими различными характеристиками:

- 1) обрабатывать тексты в новой и дореформенной орфографиях, учитывая информацию, которая несёт каждая из них (например, не просто приравнивать омофоничные буквы, но и отличать *всь* от *всѣ*);
- 2) обрабатывать тексты разных жанров от религиозных, написанных на церковнославянском или приближенном к нему, до бытовых писем в свободной орфографии (далекой от нормализации). В перспективе следует предусмотреть использование анализатора для обработки современных текстов, содержащих отступления от литературной нормы иного рода: текстов электронной коммуникации, записей речи на региональных вариантах русского языка (бытующих в русской диаспоре или иноэтнической среде) и пр.

Б) *Открытость* — способность пополняться и видоизменяться, настраиваться на разные типы текстов, возможность «обучаться» на основании пополненного словаря и т. п.<sup>2</sup>

Несмотря на требование универсальности, одновременно могут сосуществовать и модификации анализатора для текстов разных периодов и/или жанров, использующие ряд правил, специфических именно для этих текстов. Например, к текстам XX–XXI вв. едва ли нужно применять правило, согласно которому частицы *б(ы)*, *ж(е)* и *ли/ль* могут писаться слитно с предыдущим словом (что особо часто встречается в XVIII в.); это приведёт исключительно к паразитическим разборам типа *мысль* = *мы* + *ль*, *стали* = *ста* + *ли*, *ниже* = *ни* + *же* и др. (в текстах XVIII в. разборы такого рода можно отсеивать при помощи специального правила). Аналогично, такие специфические для письменности раннего XVIII в., а также XVII и предшествующих веков правила, как пропуск мягкого знака (*толко* = *только*) и отсутствие в графической системе буквы *й* (*таино* = *тайно*) приведёт к излишним неправдоподобным разборам в современных текстах (типа *банка* = *банька*, *заика* = *зайка* и т. п.). Список периодов и жанров, требующих особых модификаций анализатора, нужно будет установить опытным путём после анализа текстов, входящих в Национальный корпус русского языка. В частности, можно предположить, что различные модификации потребуются для текстов бытовых грамоток XVII в., для текстов разных периодов,

---

<sup>2</sup> В отношении несовременных текстов (представляющих собой по определению закрытый, хотя и очень большой класс) о принципе открытости анализатора можно говорить с известной долей условности, однако это обстоятельство не играет практической роли до тех пор, пока все или почти все исторические тексты не будут включены в корпус. В настоящее время мы еще очень далеки от этого (достаточно сказать, что не только оцифровано, но и вообще издано лишь незначительное меньшинство сохранившихся от XVII–XVIII вв. текстов, в том числе и бытовых текстов в нестандартной орфографии).

ориентированных на церковнославянский язык, для собственно русских текстов XVIII в., первой половины XIX в. и нескольких дальнейших периодов.

Исходя из этих требований, перед разработчиками стоят два типа задач.

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

2. *Орфографические*: обеспечение лемматизации форм, имеющих отклонения от стандартных написаний. К решению второй (орфографической) задачи имеется ряд возможных подходов.

1) Предобработка (preprocessing) — нормализация текста, своего рода «перевод» его на стандартный язык, предваряющая морфологический анализ. Нормализация текста широко применяется в устных и диалектных корпусах разных языков и в подкорпусе электронной коммуникации НКРЯ [см. об этом Гришина, Савчук 2009] в ручном режиме. При подготовке диалектных текстов для диалектного подкорпуса НКРЯ предобработка текстов, записанных в фонетической транскрипции, осуществлялась в полуавтоматическом режиме по технологии, разработанной И. Б. Качинской и Т. А. Архангельским при участии одного из соавторов данной статьи. Сначала с исходным текстом в фонетической транскрипции (так называемый текст-1) работал автоматический модуль-детранскриптор, переводящий транскрипцию в условную орфографическую запись (текст-2), а потом эта запись редактировалась и вручную переводилась (в части грамматики и фонетики) на литературный язык; этот перевод (текст-3) в дальнейшем автоматически анализировался при помощи Mystem, после чего в нём вручную снималась омонимия и проставлялась специфическая для диалектного языка разметка [Качинская 2010]. Для исторических корпусов английского языка в автоматическом режиме применяется модуль VARD [Baron, Raison 2009], который, по оценке его авторов, дает хорошие результаты при обработке широкого спектра текстов с ненормативной орфографией. Для орфографической нормализации французских текстов эпохи Возрождения разработан инструмент VariaLog, который в принципе может быть использован и для других языков [Lay 2012]. Недостатки метода предобработки заключаются в значительной затрате ресурсов для подготовки нормализованной версии текста и зачастую неоднозначности такой нормализации при объёме корпуса в несколько миллионов слов<sup>3</sup>.

---

<sup>3</sup> Вот как выглядел бы фрагмент текста после предварительной обработки (в данном случае нормализация потребовалась для половины словоформ):

```
<distinct form="Што">что</distinct>  
<distinct form="касаецца">касаецца</distinct>  
<distinct form="да">до</distinct> брата князь Александра  
<distinct form="Михаиловича">Михайловича</distinct> я вам  
<distinct form="ево">его</distinct>  
<distinct form="ваяж">воаж</distinct> в  
<distinct form="первам">первом</distinct>  
<distinct form="писме">письме</distinct>  
<distinct form="обстоятельна">обстоятельно</distinct>  
<distinct form="аписывал">описывал</distinct> с  
<distinct form="челавекам">человеком</distinct> Тургенева
```

- 2) Применение технологий параллельного корпуса — одновременное использование как оригинального текста, так и перевода орфографии на современные нормы, выровненных по предложениям или даже словоформам [Meuer 2009]. Национальный корпус русского языка уже поддерживает параллельную технологию — выравнивание текстов по предложениям, которая используется, прежде всего, в двуязычных и многоязычном параллельных подкорпусах НКРЯ (ср. [Добровольский, Кретов, Шаров, 2005], [Sitchinava 2012]). Фактически эта технология представляет собой расширение предыдущей: нормализованный текст становится доступен пользователю. Существенным недостатком такого подхода, с нашей точки зрения, является то, этот нормализованный текст по крайней мере в значительной части случаев не будет опираться на существующие критические или массовые издания (как это бывает в случае включения в корпус старых текстов в пореформенной орфографии), а будет представлять собой продукт деятельности разработчиков корпуса. Тем самым последние фактически берут на себя также функции текстологов, готовящих для широких кругов пользователей претендующее на научность издание старого текста в новой орфографии. Как представляется, подобный подход для нужд исторических корпусов НКРЯ нецелесообразен.
- 3) Учет различных уровней вариативности в грамматическом словаре анализатора. Именно данный подход лежит в основе предлагаемого в настоящей публикации. В грамматический словарь анализатора, выстроенный на базе лексикографических источников, дополнительно включаются словоформы, не распознаваемые или распознаваемые неправильно и неполно существующими анализаторами. Одновременно работают орфографические правила (они могут быть также включены в механизм индексации корпуса), позволяющие опознать в цепочке, отсутствующей в словаре, то или иное словарное слово (разбор слов *с* *і* как слов *с* *и*, разбор слов на *-тца* как слов на *-ться*); это алгоритмическая нормализация, не фиксируемая в виде какого бы то ни было промежуточного текста. При этом могут быть включены отдельным списком блокирующие правила для случаев частотных «паразитических» разборов (ср. обсуждаемые выше случаи типа *ниже* = *ни* + *же*, *стали* = *ста* + *ли*). Применение данного подхода к анализу конкретных текстов разных типов и периодов покажет, в какой степени он позволит сократить количество неправильных разборов, и, возможно, обнаружит участки, в принципе не поддающиеся автоматическому анализу и требующие ручной нормализации.

---

<distinct form="пакоинава">покойного</distinct> Ивана  
<distinct form="Сергеича">Сергеевича</distinct>, и  
<distinct form="пажалавал">пожаловал</distinct>  
<distinct form="кушел">кушал</distinct> у меня на  
<distinct form="другои">другой</distinct> день  
<distinct form="возвароту">возвороту</distinct>  
<distinct form="сваево">своего</distinct>  
<distinct form="ис">из</distinct> Теплых Станов, и с тех пор не видал;  
[Вас. Бор. Голицын Влад. Бор. Голицыну 8 апреля 1771 г.].



## 4. Грамматический словарь

### 4.1. Общие принципы

Общие принципы и алгоритм работы анализатора, области его применения были изложены в [Поляков 2012]. Грамматический словарь определяется как список лексем языка с приписанной информацией о словоизменении. Каждая лексема в словаре содержит, как минимум, следующую информацию: 1) основа с указанием чередований; 2) постоянные признаки лексемы (часть речи, род, одушевленность, переходность, и т.д.); 3) код словоизменительного типа (парадигмы). Вот пример записи для некоторых глаголов с чередованиями в основе:

но(с ш)+ить	V,ipf,tr	V4
б(и ь е)+ть	V,ipf,tr	V11
пе(к ч )+ь	V,ipf,tr	V8
ж(г ж ег е)+чь	V,ipf,tr	V8*g

Грамматический словарь анализатора складывается из нескольких модулей, соответствующих различным периодам истории русского языка. При анализе конкретного текста выбирается модуль, соответствующий типу и периоду создания текста.

- 1) Современный модуль строится на основе Грамматического словаря Зализняка [Зализняк 1977/2003], который фиксирует лексический и грамматический стандарт конца XX века. Тем не менее, этот модуль позволяет анализировать значительную часть словоформ, встречающихся в текстах XVIII–XIX века, если они не имеют существенных орфографических и грамматических отличий от современной нормы.
- 2) Модуль XVIII–XIX века строится на основе анализа корпуса реальных текстов, а также исторических словарей русского языка, включая:
  - Словарь Академии Российской (1789–1794);
  - Словарь церковнославянского и русского языка (ЦСРЯ) (1847);
  - Полный русский орфографический словарь (1898);
  - Словарь русского языка XVIII века.
- 3) Модуль XVII века и более ранних периодов строится в основном на основе анализа корпуса реальных текстов и лишь частично на основе исторического словаря (Словарь русского языка XIV–XVII века).
- 4) Модуль для церковнославянского языка сейчас существует в рамках отдельного церковнославянского корпуса, а возможность его применения для исторического корпуса пока неочевидна.

Для адекватного анализа текстов XVIII–XIX века, часть которых представлена в дореформенной орфографии, необходимо добавить в грамматическую модель анализатора следующие формы:

- 1) формы, характерные для всех периодов:

- деепричастия совершенного вида от основы презенса (*прийдя, увидя, взгромоздзясь*), которые вполне употребительны в современном языке, а тем более в языке XVIII–XIX века;
  - вариант частицы *-ся* после гласных (*валюся, валилася*), который употребляется в современном языке (в некоторых идиолектах), а также в языке XVIII–XIX века;
  - сравнительная степень на *-ей* (*сильней*) и с префиксом *по-* (*посильнее, посильней*);
- 2) формы, характерные для XIX века и более ранних периодов:
- адъективные флексии (*-аго/-яго, -ья/-ія*);
  - особые формы местоимений (*ея, онъ, онь, онднъ, онднхъ*);
  - творительный падеж 3-го склонения на *-ію* (*милостію, помощію*);
- 3) формы, характерные для XVIII века и более ранних периодов:
- усеченные формы прилагательных (*красна/о/ы/у*), которые формально совпадают с краткими формами, но имеют другое грамматическое значение;
  - сравнительная степень на *-яе* (*сильняе, скоряе*);
  - глагольные флексии *-ти* и *-ши* (*ходиши, ходити*), которые, скорее всего, должны трактоваться как церковнославянизмы (см. ниже).
- 4) церковнославянские формы:
- формы имперфекта (*творяше, творяхомъ, творяху*) и аориста (*творихъ, творихомъ, твориша*), которые нередко бывают омонимичны (*делахъ, делахомъ*);
  - частотные формы косвенных падежей существительных (*градомъ, градъхъ, троцы, троцьхъ*);
  - частотные лексемы (*иже, яко, понеже, вельми*);
- и т. д.

Анализатор, помимо словаря, сопровождается отдельными правилами для анализа текстов в разных орфографических режимах — алгоритмической нормализацией, одновременной с приписываемой морфологической разметкой. Они работают по следующему принципу: если словоформа не предсказывается на основании грамматического словаря (а может получить только гипотетический разбор), то предпринимаются попытки, исходя из её буквенного состава, регулярной замены в ней тех или иных букв и анализа получившейся нормализованной словоформы. Если эта словоформа получает негипотетический анализ (и, возможно, этот анализ удовлетворяет тем или иным грамматическим ограничениям), то он подставляется в качестве её разбора. Таковы графические (типа  $\emptyset \Rightarrow \phi$ , см. ниже, 4.3.А) и орфографические (типа *цы*  $\Rightarrow$  *ци*, см. ниже, 4.3.Б) правила, применимые ко всем словоформам, удовлетворяющим данным критериям.

Вместе с тем должен быть задействован также ряд индивидуальных правил, вводящих конкретные орфографические варианты для конкретных лемм; например, такова индивидуальная вариативность типа *естьли~если* (неизменяемое), *потчевать~подчивать*, *ветчина~вядчина* (проводится по всей парадигме).

## 4.2. Источники для формирования словника грамматического словаря

Источниками для формирования словника грамматического словаря являются, с одной стороны, существующие исторические словари русского языка (см. выше, п. 4.1), с другой — результаты анализа корпуса реальных текстов. Лексемы из этих двух источников будут добавляться к основному модулю, составленному на основе Грамматического словаря Зализняка. Этот путь представляется нам самым коротким для достижения практических целей — адекватного анализа текстов исторического корпуса. Почему нецелесообразно ограничиваться только словарями? Прежде всего потому, что значительную часть словника исторических словарей составляют книжные и устаревшие слова (включая церковнославянизмы), которые не очень часто встречаются в реальных текстах, особенно в художественных и бытовых. В то же время широко представленные в текстах собственные имена в словарях не описаны. Таким образом, создание электронного грамматического словаря исключительно на базе лексикографических источников, значительная часть которого не будет «работать» при анализе текстов, кажется нам неэффективным. Напротив, словник, полученный на основе текущего состояния корпуса и пополняемый при добавлении новых текстов, как раз и будет отражать реальное употребление.

Корпусной словник создается на базе частотного списка словоформ, извлеченных из текстов, которым при автоматическом анализе приписываются леммы и грамматические характеристики. Большая часть словоформ успешно разбирается современным анализатором; меньшую часть (около четверти) составляют словоформы, которые отсутствуют в современном грамматическом словаре. При автоматическом анализе они получают гипотетические разборы, в дальнейшем им вручную приписываются правильные леммы и грамматическая информация. По составу это а) церковнославянская лексика (актуальная для корпуса), частично перекрывается с лексикой из ЦСРЯ б) варианты, в) собственные имена и отыменные прилагательные [Савчук 2012].

Пополнение словника грамматического словаря из обоих источников будет осуществляться поэтапно, по мере подготовки электронных версий словарей, одной стороны, и развития корпуса исторических текстов, с другой.

## 4.3. Методы анализа вариативности.

Для правильного анализа и сведения к единой лемме языковых вариантов на различных уровнях анализатор может использовать правила следующего типа.

### А) Графические правила

Графические правила основаны на приравнивании графем, встретившихся в тексте, графемам, входящим в порождаемые словарём формы, например,  $\theta \Rightarrow \phi$ ;  $\xi \Rightarrow \text{кс}$ . В текстах, где распределение омофоничных графем

было неустойчивым и зависело от конкретной орфографической школы (прежде всего это среднерусские тексты и тексты раннего XVIII в., особенно бытовые), может быть задействован модуль, рассматривающий такие пары букв как равнозначные варианты и при нормализации заменяющий один на другой глобально. Вместе с тем для текстов, где достаточно устойчиво установилась смысловозначительная функция ряда омофонов (например, в русской «гротовской» орфографии *миръ* vs. *міръ*), может использоваться модуль, учитывающий эти различия при постановке леммы и морфологическом разборе.

#### Б) Орфографические правила

Орфографические правила связаны с нормализацией определённых орфограмм: иными словами, определённые буквы заменяются на другие лишь в контексте некоторых третьих, при этом обе буквы присутствуют в алфавите грамматического словаря. Таково, например, правило *цы => ци*, при соблюдении которого стандартизированный разбор получают написания типа *цыновка*, *цыдулка*, *цыгарка*, *цыгейка*, нормативные до 1956 г., причём некоторые из них встречаются в текстах и после этой даты, или правило *тца => ться/тся*, при котором восстанавливаются оба (финитный и инфинитивный) разбора для глагольных словоформ, где конечное сочетание фонетически совпало и в таком виде до XVIII в. отражалось на письме в текстах, не считавшихся грубо неграмотными. Особую роль такие правила играют в бытовой письменности XVII–XVIII вв., где полностью разрешёнными приёмами (как и, например, в поздних берестяных грамотах) является упрощение двойных согласных, передача на письме оглушения и озвончения, пропуск знака для ь и т. п.

#### 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=НАДЕЯТЬСЯ gr=praes,3p,pl;  
словоформа *надеютца* получает разбор lex=НАДЕЯТЬСЯ  
gr=praes,3p,pl=distort.

Аналогично словоформа *носитца*, для которой возможны две нормализации, получает два разбора — инфинитивный (*носится*) и финитный (*носится*), а словоформа *отца* только разборы от слова ОТЕЦ (ввиду отсутствия словоформ *\*от(ь)ся*).

#### В) Морфологические правила

Морфологические правила, в отличие от орфографических, устроены с учётом информации о грамматическом разборе словарной словоформы. Например, вводится правило, которое можно записать как *-яе, comр,аnom <=*

-*ee*, *сопр.* Это значит, что если в тексте встретилась не получающая словарного анализа словоформа, кончающаяся на *-яе*, например, *сильняе*, а замена *-яе* на *-ее* в этой словоформе даёт словарную форму (*сильнее*), и притом эта форма есть словоформа сравнительной степени, то словоформа типа *сильняе* получает такой же грамматический разбор, что и словоформа типа *сильнее* плюс помету *апот* («аномальная морфологическая форма»). Правило это может быть сформулировано и более широко, так, чтобы учитывать односложные окончания (*сильняй*) и префиксальную сравнительную степень (*посильняе*). Аналогичные правила могут вводить (с добавлением пометы *апот*) ненормативные варианты постфикса *-ся* вместо *-сь* (*оставалосся*) или наоборот (*запершегосся*).

Альтернативный способ получения адекватного морфологического анализа несловарных форм состоит в том, чтобы включить нужные формы в состав парадигм. В процессе совершенствования анализатора предполагается опробовать оба этих способа.

Пример применения морфологического правила:

в тексте представлена словоформа *имееши*;  
из-за наличия конечного *ши* проверяется правило *ши=шь*;  
словоформа *имеешь* словарем порождается и разбирается как  
*lex=ИМЕТЬ gr=praes,2p,sg*;  
из-за наличия разбора *2p,sg* словоформа *имееши* получает разбор  
*lex=ИМЕТЬ gr=praes,2p,sg=апот*.

При этом, например, несловарная словоформа *воши* не получает разбора, идентичного разбору словоформы *вошь* (поскольку она не глагольная).

Г) *Списочный способ* — задание вариативности списками, на ограниченных классах единиц. Таково, например, орфографическое правило, согласно которому ряд конкретных корней может (или даже практически обязан) в среднерусских и церковнославянских текстах сокращаться (записываться под титлом), ср. такие словоформы, как *Г(о)с(по)дь*, *м(е)с(я)ц*, *гл(агол)ет* и т. д. К списочным правилам относятся и блокирующие правила для частотных паразитических разборов типа *стали = ста ли*.

Указанные правила имеют определенный порядок применения; так, в первую очередь применяются графические правила, в том числе списочные; затем — орфографические, морфологические и списочные правила, блокирующие паразитические разборы.

В процессе формирования словаря и оптимизации работы анализатора предусматривается несколько циклов обработки текстов. Каждая новая версия анализатора будет проходить проверку на корпусе: после анализа результатов в словарь и правила анализатора будут вноситься пополнения и коррективы, после чего этот цикл повторяется уже с новой версией. Далее мы изложим оценку результатов первой версии анализатора на корпусе текстов XVIII–XIX веков.

## 5. Оценка результатов

### 5.1. Состав несловарных словоформ

Первый вариант анализатора на основе первой версии словаря был опробован на экспериментальном корпусе объемом около 4 млн словоупотреблений<sup>4</sup>, который включает 256 тыс. различных словоформ. Результаты представлены в таблице.

<b>разобрано</b>	185 221	72,4%
<b>гипотезы</b>	63 904	25,0%
<b>не разобрано</b>	6 780	2,6%

Наибольший интерес для дальнейшей разработки словаря представляет анализ словоформ, которые не были опознаны как слова русского языка или получили гипотетические разборы, и оценка предложенных гипотез. Список непознанных форм в настоящее время полностью проанализирован, всем формам вручную приписаны леммы, и они дополнили список разобранных форм. Перечислим наиболее массовые случаи.

Не распознаны или неправильно опознаны сочетания знаменательных частей речи с частицами *-то(-та, -ат), же(ж), ли(ль), бы(б), -де, -ка*, в написании которых в разные периоды наблюдались колебания. Согласно современным правилам, частицы *-то, -де, -ка* пишутся через дефис, *бы, ли, же* — раздельно. Раздельное написание этих частиц рекомендовалось уже в XIX в [Грот, 1873], однако и в XX в. вплоть до реформы 1956 г. активно использовались дефисные написания. В изданиях XVIII в. в написаниях частиц не было последовательности, можно встретить дефисные, слитные и раздельные написания: *них-же, месте-же, нихже, мыже, таковаже, пили-б, ожидали-б, где-б, пилиб, ожидалиб, еслиб*. В НКРЯ сочетания с раздельным и дефисным написанием частиц анализируются как две леммы, при этом автоматический анализ в большинстве случаев правильный. Эти же решения следует использовать и в исторической части словаря. Основную трудность как в НКРЯ, так и в конструируемом анализаторе представляет правильная идентификация частиц в составе сочетаний при слитном написании (*ежелиж, былаб* и под.).

Другую большую группу словоформ, не получивших разборов, составляют архаические формы склонения и спряжения. Среди них заметное

<sup>4</sup> Экспериментальный корпус составлен на основе текстов XVIII — первой трети XIX вв., входящих в состав НКРЯ. По жанровому составу корпус XVIII в. разнообразен: доля художественных текстов и публицистики — по 24%, церковно-богословские тексты составляют 19%, научные тексты ф — 17%, официальные документы — 11%, бытовые тексты (письма, дневники) — 5%. Приблизительно в тех же пропорциях представлены и тексты XIX в. Хронологически тексты экспериментального корпуса распределяются следующим образом: 1700–1730 — 6%, 1731–1780 — 43%, 1781–1799 — 30%, 1800–1830 — 21%. Подробнее о составе корпуса XVIII в. см. [Савчук, Сичинава 2009].

место занимают: существительные, прилагательные в форме тв. п. на *-ою* (263 формы): *Гришкою, Кабардою, прежестокою*; причастия в форме им. п. мн. ч. на *-ии* (101) / *-ьи* (18): *входящи, нарицающи, приеждаемы, украшенны*; причастия на *-яй* (128): *возвышай, вступаай*; краткие причастия на *-ущ / ащ* (27): *блистающ, властвующ*; формы имперфекта: *живаше, знаяше* (10); *бяху, стояху, мняху* (22). Часть глагольных форм была лемматизирована, но все предложенные гипотезы оказались ошибочными: форма 2 л. наст. в. на *-ши* (278): *дееши, жаждеши*; имперфекта (122): *поучаше, презираше, моляшеся; бежаху, зваху, побиваху, нарицахуся, удивляхуся, являхуся*; аориста: *искусиша, победиша* (39); *несохом, победихом* (50) и др. К архаичным глагольным формам примыкают также диалектно-просторечные формы 3 л. ед. ч. на *-ут/-ют* для глаголов 2 спряжения (40): *купют, просят, проводят, посмотрют, готовятся, находятся*. Все они займут свое место в парадигмах исторического модуля словаря.

Третью многочисленную группу составляют орфографические варианты: *безщчетну, возмеш, баталиах, полицыи, прокломащи, поосчрении* (написания с *щ* вместо *ц* совершенно регулярны, например, у Татищева, который в своем орфографическом трактате утверждал об избыточности буквы *щ* и свою фамилию писал как *Татисчев*) и др. Здесь анализ отдельных групп вариантов приведет к формулированию частных формальных правил анализа соответствующих орфограмм (см. п. 4.3).

Четвертую группу составляют многочисленные собственные имена — топонимы, имена, фамилии и отчества лиц, литературных героев и мифологических персонажей, причем многие из них присутствуют в текстах в нескольких вариантах: *Шлиссельбург* (утвердившийся впоследствии вариант), *Шлюссембурх, Шлютенбурх, Шлютелбург, Шлютельбург, Слюсинбург, Слютелбург, Слютельбург, Валпарейсо, Валпарейзо, Вальпарейзо* (в современной передаче — *Вальпараисо*, город в Чили), *Елизавета, Елисавета, Елисавет, Елисаветф, Елисавет, Ньютон, Нейтон, Невтон, Дон-Кихот, Дон-Кишот, Донкишот, Микель-Анджело, Мишель Анжело* и др. Примыкают к этой группе производные от собственных имен — прилагательные и существительные: *европскии, коперниканскии, ефесския, Антошка, Бомонтша* и т. д.

Часть неопознанных форм (около 1%) объясняется ошибками набора и сканирования, в результате проверки таких «псевдоформ» по текстам вносятся исправления.

## 5.2. Оценка гипотетических разборов.

Среди словоформ, получивших гипотетические разборы, можно выделить зоны с высоким уровнем предсказуемости (25–30%, то есть одна гипотетическая лемма из трех или четырех предложенных оказывается правильной), зоны средней предсказуемости (предлагается от пяти гипотез, из них одна правильная) и зоны с нулевой предсказуемостью (ни одна из предложенных гипотез не является правильной).

Примеры высокой предсказуемости обнаружили, в частности, существительные с основой на -к- (около 4% словоформ, получивших варианты разборов), спрягаемые формы глаголов, за исключением архаичных (более 6% словоформ), формы прилагательных (около 35% всех словоформ с гипотетическими разборами):

**доимка** доимка?=N,f,inan=sg,nom=N33\* | доимок?=N,m,inan=sg,gen=N13\*  
| доимка?=N,f,anim=sg,nom=N33\*

**напоют** напоать?=V,ipf=ind,pres,pl,3,act=V1 | напоать?=V,ipf,intr=ind,pres,pl,3,act=V1 | напоать?=V,pf=ind,fut,pl,3,act=V1

**комфузные** комфузный?=A=pl,nom/acc=A1\* | комфузной?=A=pl,nom/acc=A1b | комфузная?=N,f,inan=pl,nom/acc/acc=A1 | комфузный?=A=pl,nom/acc=A1 | комфузные?=N,pl,inan=pl,nom/acc/acc=A1

Из архаических форм в зоне высокой предсказуемости находятся формы прилагательных мн. ч им. п. на -ия /-ья, -аго:

**одинакия** одинакий?=A=pl,nom/acc=A3 | одинакий?=A=pl,nom/acc=A3\*  
| одинакая?=N,f,inan=pl,nom/acc/acc=A3 | одинакой?=A=pl,nom/acc=A3b | одинакое?=N,n,inan=pl,nom/acc/acc=A3

Эти формы правильно опознаются как прилагательные и отграничиваются от форм существительных с омонимичными финалями -ия:

**министра** **министра**?=N,f,inan=sg,nom=N37 | министра?=N,topn,f,inan=sg,nom=N37 | министерий?=N,persn,m,anim=sg,gen/acc=N17 | министерий?=N,m,inan=sg,gen=N17 | министра?=N,persn,f,anim=sg,nom=N37

Пример средней степени предсказуемости:

**неведь** неведь?=N,f,inan=sg,nom/acc=N41 | неведь?=N,m,anim=sg,nom=N12 | неведь?=N,topn,m,anim/inan=sg,nom=N12 | **неведь**?=CONJ/PART | неведь?=PREP

Наибольшие сложности представляет идентификация архаических глагольных форм. Все они попадают в зону нулевой предсказуемости. Причина объяснена выше — отсутствие в словнике в достаточном объеме церковнославянских глаголов и парадигм спряжения.

Форма **вознесоста** (2 л. двойст. числа аориста глагола *вознести*) получает гипотетические разборы, ни один из которых не является верным:



вознесост?=N,m,anim=sg,gen/acc=N11 | вознесост?=N,m,inan=sg,gen=N11  
| вознесоста?=N,f,inan=sg,nom=N31 | вознесост?=N,persn,m,anim=  
sg,gen/acc=N11 | вознесост?=N,topn,m,inan=sg,gen=N11

Форме **видиши** (2 л. ед. ч. наст. вр. глагола *видети*) приписаны ошибочные разборы<sup>5</sup>:

видиша?=N,f,inan=pl,nom/acc!sg,gen=N34 |  
видисать?=V,pf=imp,sg,2,act=V6t | видисать?=V,ipf=imp,sg,2,act=V  
6t | видиша?=N,persn,m,anim=pl,nom!sg,gen=N34

Нулевую предсказуемость имеют формы им.п. мн.ч. прилагательных на *-ии* (современное *-ие*).

великии – великия?=N,topn,f,inan=pl,nom/  
acc!sg,d!sg,gen!sg,loc=N37

глухии – глухия?=N,f,inan=pl,nom/acc!sg,d!sg,gen!sg,loc=N37 |  
глухия?=N,topn,f,inan=pl,nom/acc!sg,d!sg,gen!sg,loc=N37

В отличие от прилагательных существительные на *-ия* опознаются с высокой степенью правильности:

**девизии** – девизия?=N,f,inan=pl,nom/acc!sg,d!sg,gen!sg,loc=N37 |  
девизия?=N,topn,f,inan=pl,nom/acc!sg,d!sg,gen!sg,loc=N37 |  
девизие?=N,n,inan=sg,loc=N27

**энергии** – энергия?=N,f,inan=pl,nom/acc!sg,d!sg,gen!sg,loc=N37 |  
энергия?=N,topn,f,inan=pl,nom/acc!sg,d!sg,gen!sg,loc=N37

Если учесть, что прилагательных в этой зоне в 2 раза больше, чем существительных, в разбор словоформ следует включить грамму прилагательного, применив специальные фильтры. В частности, для форм, оканчивающихся на *-ии* (их около 700), справедливо следующее наблюдение:

<sup>5</sup> Архаические глагольные формы — проблема и для анализатора Mystem, который в отдельных случаях «взрывается» ложными гипотезами. В частности, для формы доставляеши в НКРЯ предлагается 8 лемм и 37 гипотетических вариантов анализа.

Лемма доставляеши – сущ, фам, одуш, м, 0 / сущ, неод, с, ед, 0 / сущ, имя, одуш, м, 0  
Лемма доставляеший – прил, мн, кратк

Лемма доставляешишь – глаг, нп, нсв, повел, действ, 2л, ед / глаг, нп, св, повел,  
действ, 2л, ед / глаг, перех, нсв, повел, действ, 2л, ед / глаг, перех, св, повел,  
действ, 2л, ед

Лемма доставляеш – сущ, фам, одуш, м, мн, им / сущ, неод, м, мн, вин, геогр / сущ,  
неод, м, мн, им, геогр / сущ, имя, одуш, м, мн, им /

Лемма доставляеша – сущ, фам, одуш, м, мн, им / сущ, фам, одуш, м, ед, род / сущ, фам,  
одуш, м-ж, мн, им / сущ, фам, одуш, м-ж, ед, род

и т.д.

- часть основы на гласную + *-нии* => существительное (*Гавании, Казании, Исмении, докончании, догорении*);
- часть основы на согласную + *-нии* => прилагательное (*главнии, вчерашнии, доволнии*);
- часть основы на н + *-нии* => прилагательное или причастие (*украшеннии, смиреннии*).

В зоне нулевой предсказуемости находится также подавляющее большинство (более 300) форм тв.п. существительных и прилагательных на *-ою* (*ескадрою, Козмою, кавалерственною, манетною*), притяжательные прилагательные на *-ин* (*богинин, венерин, минервин* и др.).

Список словоформ с гипотетическими разборами находится в стадии обработки: словоформам с ошибочными разборами вручную приписываются леммы и грамматические признаки, для форм со средней предсказуемостью рассматриваются способы сокращения предлагаемых гипотез.

### 5.3. Сокращение количества гипотез

Сокращения предлагаемых гипотез можно добиться путем использования правил, которые будут применяться к определенным классам словоформ, образующим закрытые или пополняемые списки. Часть правил будет основана на учете связи морфологических признаков с морфемной структурой слова. Так, например, для форм на *-ения* анализатор предлагает 5 гипотетических лемм:

```
напоения напоение?=N,n,inan=pl,nom/acc|sg,gen=N27 |  
напоения?=N,f,inan=sg,nom=N37 | напоения?=N,topn,f,inan=sg,nom=N37 | напоения?=N,persn,f,anim=sg,nom=N37 |  
напоений?=N,m,inan=sg,gen=N17
```

```
посмотрения посмотрение?=N,n,inan=pl,nom/acc|sg,gen=N27 | посмот  
рения?=N,f,inan=sg,nom=N37 | посмотрения?=N,topn,f,inan=sg,nom=N37 | посмотрения?=N,persn,f,anim=sg,nom=N37 | посмотрений?  
=N,m,inan=sg,gen=N17
```

При этом статистически формы распределяются следующим образом. Из 444 форм 439 (98%) являются формами род.п. существительных среднего рода, 2 формы — род.п. мужского рода и 3 — им.п. женского рода (имена собственные). Следовательно, формам на *-ения* целесообразно приписывать леммы существительных среднего рода как наиболее статистически значимые, а существительные женского и мужского рода задать списком и включить в словарь.

Аналогичная картина наблюдается у существительных с суффиксом *-ствиј* (*-ствие, -ствия, -ствиц, -ствию* и т.д.): 206 из 207 словоформ относятся к существительным среднего рода, 1 форма — предлог (*вследствии*). Следовательно, набор лемм, который предлагается при анализе словоформ этого достаточно продуктивного класса, можно сократить до одной.

**неблагодарствию** **неблагодарствие**?=N,n,inan=sg,loc=N27 |  
неблагодарствия?=N,topn,f,inan=pl,nom/  
acc|sg,dat|sg,gen|sg,loc=N37 |  
неблагодарствия?=N,persn,f,anim=pl,nom|sg,dat|sg,gen|sg,loc=N37 |  
неблагодарствий?=N,persn,m,anim=pl,nom|sg,loc=N17 |  
неблагодарствия?=N,f,inan=pl,nom/acc|sg,dat|sg,gen|sg,loc=N37

Среди небольших групп, в которых целесообразно уменьшить количество вероятных гипотез, можно назвать характерные для языка XVIII века германизмы со специфическими исходами: формы на *-берг* (топонимы и фамилии), *-бург* (топонимы), на *-ау* (топонимы славянского происхождения, например, *Бункау, Лаубау*); на *-мейстер, -мистр* (сущ. муж. рода, одушевленные *герольдмейстер, ширмейстер, вафмистр, виц-вахмистр, квартирмистр, секунд-ротмистр*) и др.

Другой источник сокращения количества гипотез, как уже говорилось в п. 4.3, видится в анализе орфографической вариативности, реально представленной в текстах, и применении орфографических преобразований. Приведем несколько примеров.

Формы с приставками на *-з*: формы с начальным *без-/в(о)з-/из-/низ-/раз-/роз-/ч(е)рез-/з-* перед глухими анализировать как формы с начальным *бес-/в(о)с-/ис-/нис-/рас-/рос-/ч(е)рес-/с-* (для приставки *без-* такое обобщение уже сделано).

Слитное написание *не* с глаголами: *небудет, неисповедует, ненайдет, нетребует*;

*сч => иц*: *поосчрение=поощрение, немосчный=немошный* и др., более 500 форм;

*жсы, шы => жи, ши*: *показавшые, ваши, стражсы, лжсы* (55);

конечное *-ья => ия, -ьи => ии*: *Францья, полицья, позижья* (10), *амуницьи, полицьи* (18), конечное *-иа => -ия*: *библиа, благополучиа* (145), *-иах => иях*: *баталшах, материах* (23);

*-лск(ий, ого, аго, ому и т.д.) => льск(ий, ого, аго, ому и т.д.)*: *посолский, тоболской, неприятелское, евангелский* (87), *-лст- => -лст-*: *лстивый, жителство, посолство, ловителствуя, началству, обстоятелство* (более 130). Список орфографических особенностей текстов в старой орфографии см. также в [Поляков 2012].

## 6. Заключение

Принципы формирования грамматического словаря для автоматического анализа текстов XVIII–XIX вв., описанные в настоящей статье, были опробованы на экспериментальном корпусе, составленном из текстов XVIII — 1-ой трети XIX в. Результаты автоматической морфологической разметки показали, что более 73% словоформ корпуса получили однозначные разборы, что свидетельствует о степени совпадения текстов данного периода и современных текстов по словарному составу. Однако следует учесть, что относительно высокий процент совпадения в значительной мере объясняется составом корпуса: в нем мало текстов начала XVIII в., кроме того, большая часть текстов, в соответствии

с ориентацией на авторитетные научные издания советского и постсоветского времени, представлена в современной орфографии или с некоторыми отклонениями от современной нормы. Очевидно, что увеличение доли оригинальных текстов XVII и первой трети XVIII в. изменило бы результат анализа в сторону уменьшения доли совпадающих словоформ и привело бы к росту количества ошибочных разборов.

Изучение и классификация ошибок анализатора позволило наметить пути дальнейшего расширения и совершенствования словаря. Это во-первых, ручной анализ неопознанных форм и включение соответствующих лемм в словарь; во-вторых, пополнение списка парадигм за счет церковнославянских и старорусских форм, а также включение вариативных грамматических форм в состав отдельных современных парадигм; в третьих, составление орфографических правил для алгоритмической нормализации. В ближайшие планы входит внедрение всех изменений в словарь, тестирование его новой версии на экспериментальном корпусе, откорректированном с учетом выявленных ошибок, а также на отдельных текстах более раннего периода, относящихся к разным жанрам.

## Литература

1. *Большаков И. А., Большакова Е. И.* Автоматический морфоклассификатор русских именных групп // Компьютерная лингвистика и интеллектуальные технологии. По материалам ежегодной Международной конференции «Диалог» (2012). Вып. 11. М., 2012. С. 81–92.
2. *Гришина Е. А., Савчук С. О.* Корпус устных текстов в НКРЯ: состав и структура // Национальный корпус русского языка: 2006–2008. Новые результаты и перспективы. СПб.: Нестор-История, 2009. С. 129–149.
3. *Грот Я. К.* Спорные вопросы русского правописания от Петра Великого до ныне. СПб, Типография императорской АН, 1873.
4. *Добровольский Д. О., Кретов А. А., Шаров С. А.* Корпус параллельных текстов: архитектура и возможности использования. // Национальный корпус русского языка: 2003–2005. — М.: Индрик, 2005. С. 263–296.
5. *Зализняк 1977/2003* — Зализняк А. А. Грамматический словарь русского языка. Изд. 1-е. М., 1977 (4-е изд., испр. и доп., М. 2003)
6. *Качинская И. Б.* Диалектный подкорпус НКРЯ. Новый стандарт подачи. Новое рабочее место // О. Ю. Крючкова и др. (ред.) Русская устная речь. Материалы международной научной конференции «Баранниковские чтения. Устная речь: русская диалектная и разговорно-просторечная культура общения» и межвузовского совещания «Проблемы создания и использования диалектных корпусов». Саратов, Издательский центр «Наука», 2011. С. 245–255
7. *Поляков А. Е.* Проблемы и методы анализа русских текстов в дореформенной орфографии // Компьютерная лингвистика и интеллектуальные технологии. По материалам ежегодной Международной конференции «Диалог» (2012). Вып. 11. М., 2012. С. 536–547.

8. Поляков А. Е., Добрушина Е. Р., Иванова-Алленова Т. Ю. Корпус церковнославянских текстов в составе НКРЯ, первая версия: проблемы и решения. Информационные технологии и письменное наследие. Материалы международной научной конференции. — Петрозаводск, 2012. С. 211–215.
9. Савчук С. О., Сичинава Д. В., Гарипов И. Подкорпус текстов XVIII века в составе Национального корпуса русского языка: из опыта работы // Web Journal of Formal, Computational & Cognitive Linguistics. Специальный выпуск (Труды Российского научно-образовательного центра по лингвистике им И. А. Бодуэна де Куртенэ), 2006.
10. Савчук С. О. Корпус русских текстов XVIII века в составе Национального корпуса русского языка: проблемы и перспективы // Информационные технологии и письменное наследие. Материалы международной научной конференции. Казань, 2008. С. 241–244.
11. Савчук С. О., Сичинава Д. В. Корпус русских текстов XVIII века в составе НКРЯ: проблемы и перспективы // Национальный корпус русского языка: 2006–2008. Новые результаты и перспективы. СПб.: Нестор-История, 2009. С. 52–70.
12. Савчук С. О. Электронный словарь вариантов на основе текстов XVIII в. Информационные технологии и письменное наследие. Материалы международной научной конференции. — Петрозаводск, 2012. С. 241–244.
13. Соловьев В. Д., Ахтямов Р. Б. Корпус русского языка XVIII века: текущее состояние/ Материалы международной научной конференции Ижевск, 13–17 июля 2006 г. Ижевск, 2006. С. 156–160.
14. Baron, A., Raison, P. (2009) Automatic standardization of texts containing spelling variations. How much training data do you need? // In M. Mahlberg, V. González-Díaz, and C. Smith (eds.), Proceedings of the Corpus Linguistics Conference, CL2009. University of Liverpool, UK, pp. 20–23, [http://ucrel.lancs.ac.uk/publications/CL2009/314\\_FullPaper.pdf](http://ucrel.lancs.ac.uk/publications/CL2009/314_FullPaper.pdf)
15. Lay, M. H. (2012) VariaLog: how to locate words in a French Renaissance Virtual Library // Digital Humanities Conference, University of Hamburg, Germany, 2012, available at: <http://www.dh2012.uni-hamburg.de/conference/programme/abstracts/varialog-how-to-locate-words-in-a-french-renaissance-virtual-library/>
16. Meyer, R. (2009) Semi-automatic morphosyntactic tagging of a diachronic corpus of Russian // In M. Mahlberg, V. González-Díaz, and C. Smith (eds.), Proceedings of the Corpus Linguistics Conference, CL2009. University of Liverpool, UK, pp. 20–23. <http://ucrel.lancs.ac.uk/publications/cl2009/abstracts.htm#347>
17. Sitchinava D. (2012) Parallel corpora within the Russian National Copus // *Prace Filologiczne*, LXIII, 2012. С. 271–278.

## References

1. *Baron, A., Raison, P.* (2009), Automatic standardization of texts containing spelling variations. How much training data do you need? M. Mahlberg, V. González-Díaz, and C. Smith (eds.), Proceedings of the Corpus Linguistics Conference, CL2009. University of Liverpool, UK, pp. 20–23, available at: [http://ucrel.lancs.ac.uk/publications/CL2009/314\\_FullPaper.pdf](http://ucrel.lancs.ac.uk/publications/CL2009/314_FullPaper.pdf)
2. *Bolshakov, I. A., Bolshakova, E. I.* (2012), An Automatic morphological classifier of noun phrases in Russian [Avtomaticheskij morfoklassifikator russkih imennyh grupp]. Computational Linguistics and Intellectual Technologies: Papers from the Annual International Conference “Dialog” 2012 [Komp’iuternaia Lingvistika i Intellektual’nye Tehnologii: Po materialam ezhegodnoj mezhdunarodnoj konferencii “Dialog” 2012]. Bekasovo, pp. 81–92.
3. *Dobrovol’skij D. O., Kretov A. A., Sharov S. A.* (2005), Parallel Corpus: architecture and usability [Korpus parallel’nyh tekstov: arhitektura i vozmozhnosti ispol’zovanija], in Russian National Corpus: 2003–2005 [Nacional’nyj korpus russkogo jazyka: 2003–2005], Indrik, Moscow, pp. 263–296.
4. *Grishina E. A., Savchuk S. O.* (2009), Spoken texts in the RNC: composition and structure [Korpus ustnyh tekstov v NKRJa: sostav i struktura], in Russian National Corpus: 2006–2008. New Results and Perspectives [Nacional’nyj korpus russkogo jazyka: 2006–2008. Novye rezul’taty i perspektivy]. Nestor-Istorija, SPb, pp. 129–149.
5. *Grot Ja. K.* (1873) Controversial issues of Russian spelling since Peter the Great until now [Spornye voprosy russkogo pravopisanija ot Petra Velikogo donyne]. Tipografija imperatorskoj AN, SPb.
6. *Kachinskaja I. B.* (2011), Dialectal subcorpus of the RNC. The new standards. New workplace [Dialektnyj podkorpus NKRJa. Novyj standart podachi. Novoe rabochee mesto], in O. Ju. Krjuchkova i dr. (red.) Russian speech. Proceedings of the International Conference “Barannikovskie reading. Spoken speech: Russian dialect and colloquial vernacular culture of communication” and intercollegiate conference “Development and Use of dialect corpora” [Russkaja ustnaja rech’. Materialy mezhdunarodnoj nauchnoj konferencii “Barannikovskie chtenija. Ustnaja rech’: russkaja dialektnaja i razgovorno-prostorechnaja kul’tura obshhenija” i mezhvuzovskogo soveshhanija “Problemy sozdanija i ispol’zovanija dialektnyh korpusov”], Izdatel’skij centr “Nauka”, Saratov, pp. 245–255.
7. *Lay, M. H.* (2012), VariaLog: how to locate words in a French Renaissance Virtual Library, Digital Humanities Conference, University of Hamburg, Germany, 2012, available at: <http://www.dh2012.uni-hamburg.de/conference/programme/abstracts/varialog-how-to-locate-words-in-a-french-renaissance-virtual-library/>
8. *Meyer, R.* (2009), Semi-automatic morphosyntactic tagging of a diachronic corpus of Russian, M. Mahlberg, V. González-Díaz, and C. Smith (eds.), Proceedings of the Corpus Linguistics Conference, CL2009. University of Liverpool, UK, pp. 20–23, available at: <http://ucrel.lancs.ac.uk/publications/cl2009/abstracts.htm#347>

9. *Poljakov A. E.* (2012), Problems and methods in analysis of Russian texts in the pre-reform spelling [Problemy i metody analiza russkikh tekstov v doreformennoj orfografii], *Computational Linguistics and Intellectual Technologies: Papers from the Annual International Conference "Dialog" 2012* [Komp'uternaia Lingvistika i Intellektual'nye Tehnologii: Po materialam ezhegodnoj mezhdunarodnoj konferencii "Dialog" 2012]. Bekasovo, pp. 536–547.
10. *Poljakov A. E., Dobrushina E. R., Ivanova-Allenova T. Ju.* (2012), Corpus of Church Slavonic texts in the RNC, the first version: problems and solutions. [Korpus cerkovnoslavjanskih tekstov v sostave nkrja, pervaja versija: problemy i reshenija], *Information technology and the written heritage. Proceedings of the International Conference* [Informacionnye tehnologii i pis'mennoe nasledie. Materialy mezhdunarodnoj nauchnoj konferencii], Petrozavodsk, pp. 211–215.
11. *Savchuk S. O., Sichinava D. V., Garipov I.* (2006), Subcorpus of the XVIIIth century texts in the Russian National Corpus [Podkorpus tekstov XVIII veka v sostave Nacional'nogo korpusa russkogo jazyka: iz opyta raboty], *Web Journal of Formal, Computational & Cognitive Linguistics. Special Issue (Proceedings of the Baudouin de Courtenay Russian Research and Educational Center in linguistics)* [Special'nyj vypusk (Trudy Rossijskogo nauchno-obrazovatel'nogo centra po lingvistike im. I. A. Boduena de Kurtene)], available at: <http://fcl.ksu.ru/fclpap.htm>.
12. *Savchuk S. O.* (2008), Corpus of the Russian XVIIIth century texts in the Russian National Corpus: problems and prospects [Korpus russkikh tekstov XVIII veka v sostave Nacional'nogo korpusa russkogo jazyka: problemy i perspektivy]. *Information technologies and written heritage. Proceedings of the International Conference* [Informacionnye tehnologii i pis'mennoe nasledie. Materialy mezhdunarodnoj nauchnoj konferencii], Kazan, p. 241–244.
13. *Savchuk S. O., Sichinava D. V.* (2009), Corpus of the Russian XVIIIth century texts in the RNC: Problems and Perspectives [Korpus russkikh tekstov XVIII veka v sostave NKRJa: problemy i perspektivy], in *Russian National Corpus: 2006–2008. New Results and Perspectives* [Nacional'nyj korpus russkogo jazyka: 2006–2008. Novye rezul'taty i perspektivy], Nestor-Istorija, SPb, pp. 52–70.
14. *Savchuk S. O.* (2012), Electronic dictionary of variants based on the 18th century texts [Elektronnyj slovar' variantov na osnove tekstov XVIII v.], *Information technology and the written heritage. Proceedings of the International Conference* [Informacionnye tehnologii i pis'mennoe nasledie. Materialy mezhdunarodnoj nauchnoj konferencii], Petrozavodsk, pp. 241–244.
15. *Sitchinava D.* (2012), Parallel corpora within the Russian National Corpus. *Prace Filologiczne*, LXIII, 2012, pp. 271–278
16. *Solovyev V. D., Akhtyamov R. B.* (2006), Corpus of the XVIIIth century Russian: the present state of affairs. [Korpus russkogo jazyka XVIII veka: tekushhee sostojanie]. *Proceedings of the International Conference* [Materialy mezhdunarodnoj nauchnoj konferencii]. Izhevsk, pp. 156–160.
17. *Zalznjak, A. A.* (1977/2003), *Grammatical dictionary of the Russian language* [Grammaticheskij slovar' russkogo jazyka], Moscow (4 ed. 2003).