



НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ

Научно-учебная лаборатория учебных  
корпусов

# **РАЗРАБОТКА ПОДСИСТЕМЫ ЭКСПОРТА-ИМПОРТА КОРПУСОВ ТЕКСТОВ**

Подготовил: Зайцев Константин, стажёр-  
исследователь

Пермь, 2020

# ПОСТАНОВКА ЗАДАЧИ

Разработка хранилища корпусов текстов в гетерогенных форматах, поддерживающего возможности интеграции в сервисно-ориентированную архитектуру и импорта-экспорта корпусов в различных форматах.

# ПОНЯТИЕ КОРПУСНОЙ ЛИНГВИСТИКИ

## Определение

**Корпусная лингвистика** опирается на понятия корпуса.

**Корпус** – множество текстов, соответствующее определённым свойствам.

Свойства:

1. *Электронный*. Современный корпус должен храниться в электронном виде.
2. *Репрезентативный*. Корпус должен отражать изучаемый язык в полной мере.
3. *Размеченный*. Корпус должен содержать определённую разметку (синтаксическую, морфологическую, семантическую).
4. *Целеориентированный*. Каждый корпус создаётся для конкретных целей исследований.

# ПОНЯТИЕ КОРПУСНОЙ ЛИНГВИСТИКИ

## Корпусные менеджеры

Для работы с корпусами используются различные инструменты работы с естественным языком.

В число возможностей инструментов входят следующие операции:

1. Поиск словоформ по определённым параметрам.
2. Отображение информации о тексте.
3. Вывод результатов поиска.
4. Получение статистических данных.
5. Сохранение конкорданса (представление искомой словоформы в контексте и со статистическими данными).

# ПОНЯТИЕ КОРПУСНОЙ ЛИНГВИСТИКИ

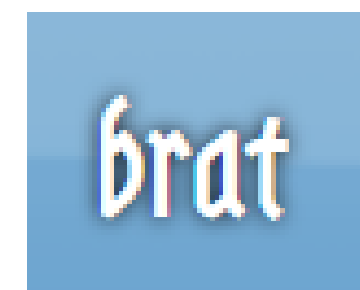
Примеры инструментов работы над корпусами

Для представления корпусов могут использоваться следующие инструменты:

1. Sketch Engine ([www.sketchengine.co.uk](http://www.sketchengine.co.uk)).
2. Gate Developer (<https://gate.ac.uk/family/developer.html>).
3. BRAT (<https://brat.nlplab.org/>).

Sketch  Engine

 general architecture  
for text engineering





# ПРИМЕРЫ КОРПУСОВ

Примеры русскоязычных корпусов

## Национальный корпус русского языка (НКРЯ)

Проект поддерживается лингвистами из разных городов (Москва, Воронеж, Санкт-Петербург, Казань и т.д.) и институтов.

В Национальный корпус включены литературные тексты (с начала 18 века), переводные сочинения, поэтические тексты, тексты, представляющие нелитературные формы языка – разговорную и диалектную.

Для каждого подкорпуса имеется возможность поиска.



# ПРИМЕРЫ КОРПУСОВ

## Примеры русскоязычных корпусов

Поиск точных форм ? А Б В

Слово или фраза

искать очистить

Лексико-грамматический поиск ?

|  |  |   |
|--|--|---|
| <p>Слово <span>?</span> <span>А Б В</span></p> <input type="text"/>              | <p>Грамм. признаки <span>?</span> <a href="#">выбрать</a></p> <input type="text"/> | <p>Семант. признаки <span>?</span> <a href="#">выбрать</a></p> <input type="text"/>   |
| <p>Доп. признаки <span>?</span> <a href="#">выбрать</a></p> <input type="text"/> | <input type="text"/>   | <p><input checked="" type="checkbox"/> 1-е знач. <input checked="" type="checkbox"/> др. знач. <input type="checkbox"/> фильтр 1 <input type="checkbox"/> фильтр 2 <span>?</span></p> |

Расстояние: от  до  ?

|  |  |   |
|--|--|---|
| <p>Слово <span>?</span> <span>А Б В</span></p> <input type="text"/>              | <p>Грамм. признаки <span>?</span> <a href="#">выбрать</a></p> <input type="text"/> | <p>Семант. признаки <span>?</span> <a href="#">выбрать</a></p> <input type="text"/>   |
| <p>Доп. признаки <span>?</span> <a href="#">выбрать</a></p> <input type="text"/> | <input type="text"/>   | <p><input checked="" type="checkbox"/> 1-е знач. <input checked="" type="checkbox"/> др. знач. <input type="checkbox"/> фильтр 1 <input type="checkbox"/> фильтр 2 <span>?</span></p> |

искать очистить





# ПРИМЕРЫ КОРПУСОВ

## Примеры русскоязычных корпусов

Объем всего корпуса: 124 565 документов, 321 712 061 слово.

"корпус"

Найдено 3 097 документов, 8 446 вхождений.

[Распределение по годам](#) [Статистика](#)

Поискать в других корпусах: [акцентологическом](#), [газетном](#), [диалектном](#), [мультимедийном](#), [обучающем](#), [параллельном](#), [поэтическом](#), [устном](#).

Страницы: [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [10](#) [11](#) [следующая страница](#)

1. [И. А. Барков. Автоматический синтез структурного описания конструкции // «Информационные технологии», 2004](#) [\[омонимия снята\]](#) [Все примеры](#)

В приведённом примере описания у конструкторского понятия "сборочная единица" раздел Состав присутствует, а у обозначающих детали понятий "**корпус**" "втулка" "вал" раздел Состав пуст. [И. А. Барков. Автоматический синтез структурного описания конструкции // «Информационные технологии», 2004] [\[омонимия снята\]](#) [←...→](#)  
В примере понятия "сборочная\_единица" "**Корпус**" "Втулка" "Вал" отражают конструкторские представления о свойствах изделия, а понятия "Размерная\_цепь" "Размер" "Оборудование" — представления о том же изделии инженера-технолога. [И. А. Барков. Автоматический синтез структурного описания конструкции // «Информационные технологии», 2004] [\[омонимия снята\]](#) [←...→](#)

2. [Михаил Кузьминский. Водопады Supermicro // «Computerworld», 2004](#) [\[омонимия снята\]](#) [Все примеры](#)

**Корпус** сервера-лезвия с формфактором 7 U и блоками питания мощностью 4, 5 либо 3 кВт может включать в себя и модули коммутаторов, в том числе InfiniBand. [Михаил Кузьминский. Водопады Supermicro // «Computerworld», 2004] [\[омонимия снята\]](#) [←...→](#)





# ПРИМЕРЫ КОРПУСОВ

Примеры русскоязычных корпусов

## Открытый корпус русского языка (OpenCorpora)

OpenCorpora – это проект по созданию размеченного корпуса силами сообщества.

Размеченные корпуса находятся в открытом доступе, каждый может скачать их и поработать с ними.

В проекте может поучаствовать любой заинтересованный человек, достаточно зарегистрироваться и выбрать задание.



# ПРИМЕРЫ КОРПУСОВ

Примеры иностранных корпусов

## British National Corpus (BNC)

Британский национальный корпус – собрание текстов на английском, включающий письменную и устную речь.

Корпус создан при помощи издательства Оксфордского университета.

Письменная часть корпуса представлена большим количеством жанров: новостными газетами, журналами, научными публикациями, эссе и многое другое.

Разговорный подкорпус содержит транскрибированные неформальные разговоры, формальные бизнес встречи или государственные мероприятия, передачи на радио и т.д.



# ПРИМЕРЫ КОРПУСОВ

Примеры иностранных корпусов

## Open American National Corpus (OANC)

Американский национальный корпус – корпус английского языка, в который входят письменные и устные источники – сообщения на почте, твиты, текст с веб-страниц.

Можно стать участником в формировании и аннотации корпуса.

# ПРИМЕРЫ КОРПУСОВ

Примеры иностранных корпусов

## **Corpus of American English (COCA)**

Корпус содержит огромное количество текстов (суммарно больше 25 миллионов слов с 1999 по 2019 гг.) на английском языке по разным жанрам: разговорный, художественный, популярные журналы, новостные газеты, научные тексты и, с недавнего времени, телевидение, субтитры, блоги и веб-страницы.

Имеется возможность проводить поиск по корпусу таким же образом, как в национальном корпусе русского языка.



# ПРИМЕРЫ КОРПУСОВ

## Примеры иностранных корпусов

[List](#) [Chart](#) [Word](#) [Browse](#) [+](#)

[POS] ?

[Find matching strings](#) [Reset](#)

☐ Sections [Texts/Virtual](#) [Sort/Limit](#) [Options](#)

(HIDE HELP) NOT LOGGED IN

[Take a Tour](#) [Take a guided tour of English-Corpora.org \(Nov 2020\)](#)

[Download the corpus \(and corpus-based frequency data\) for offline use](#)

The Corpus of Contemporary American English (COCA) is the only **large**, genre-balanced corpus of American English. COCA is probably the **most widely-used corpus of English**, and it is related to many other **corpora of English** that we have created, which offer unparalleled insight into **variation in English**.

The corpus contains more than **one billion words** of text (25+ million words each year 1990-2019) from eight genres: spoken, fiction, popular magazines, newspapers, academic texts, and (with the **update in March 2020**): TV and Movies subtitles, blogs, and other web pages.

Click on any of the links in the search form to the left for context-sensitive help, and to see the range of queries that the corpus offers.



# ПРИМЕРЫ КОРПУСОВ

## Примеры иностранных корпусов

FIND SAMPLE: [100](#) [200](#) [500](#) [1000](#)  
PAGE: << < 1 / 41 > >>

| CLICK FOR MORE CONTEXT |      |      |                         |   |   |   | [?]  |   | SHOW DUPLICATES  |
|------------------------|------|------|-------------------------|---|---|---|--|---|--|
| 1                      | 2012 | BLOG | crawfishboxes.com       | A | B | C | acquisition  | Oberholtzer   | strangely enough upped his strikeout rate from 6.6 to 9.2 after joining the Corpus Christi hooks, his ERA also jumped to 5.27 in 27. |
| 2                      | 2012 | BLOG | crawfishboxes.com       | A | B | C | rate (SO/9)  | and a 2.2 walks per nine rate (BB/9) between Corpus Christi and Oklahoma City. His ERA between the two levels was a less than |  |
| 3                      | 2012 | BLOG | crawfishboxes.com       | A | B | C | At Oklahoma he posted an average Game Score of 48.47 in 15 starts. At Corpus Christi he posted a Game Score of 50.62 in 13 starts. # The thing                 |   |  |
| 4                      | 2012 | BLOG | ...rspains.blogspot.com | A | B | C | headquarters, tried by a military court 6-7 May, denied a writ of habeas corpus, and sentenced to 2 years' confinement in a military prison. Following a       |   |  |
| 5                      | 2012 | BLOG | ...rspains.blogspot.com | A | B | C | Parte Vallandigham, which decreed that the Court could not issue a writ of habeas corpus in a military case, and for a Democratic campaign slogan he crea      |   |  |
| 6                      | 2012 | BLOG | ...review.wordpress.com | A | B | C | Illias Panagiotaros (an MP with the Golden Dawn who participated in the attacks outside Corpus Cristi last week) admitted to Newsnight's Paul Mason last       |   |  |
| 7                      | 2012 | BLOG | reason.com              | A | B | C | , I object more to the actions Lincoln took during the war (suspending habeas corpus, the draft, etc) than the war itself. And I'm in                          |   |  |
| 8                      | 2012 | BLOG | reason.com              | A | B | C | less casualties. # A hero who invaded a sovereign nation, illegally suspended habeas corpus and imprisoned thousands of political enemies, tacitly sanctio     |   |  |
| 9                      | 2012 | BLOG | chronicle.com           | A | B | C | BYU had previously created his own POS-tagged version of Google Ngrams as part of his corpus collection, he only had access to the publicly available data     |   |  |
| 10                     | 2012 | BLOG | chronicle.com           | A | B | C | threshold of 40 occurrences for inclusion) and thus wasn't able to parse the corpus in a systematic fashion. " # (The answer is putting only after he          |   |  |
| 11                     | 2012 | BLOG | laboratorium.net        | A | B | C | remind people that not everyone is pleased that Google did not come by their complete corpus by totally honest means, in other words the foundations of        |   |  |
| 12                     | 2012 | BLOG | laboratorium.net        | A | B | C | in-copyright works as currency to get a better deal for their subscription to the Google corpus. The University of Wisconsin is no better than a school yard l |   |  |
| 13                     | 2012 | BLOG | laboratorium.net        | A | B | C | # And then there is the question of who gets access to the complete Google corpus (which contains works from the world over) by subscriptions -- just Am       |   |  |
| 14                     | 2012 | BLOG | laboratorium.net        | A | B | C | At the D is for Digitize Symposium?) the point that Google owns the corpus and could if it wanted to sell it... " even to China ".                             |   |  |
| 16                     | 2012 | BLOG | laboratorium.net        | A | B | C | his promises. " The settlement " gives Google the right to sell the corpus, " she said. " Suppose they sell it to Rupert Murdoch... Wal-Mart                   |   |  |
| 17                     | 2012 | BLOG | laboratorium.net        | A | B | C | to engage in lawful large-scale book digitization projects, thereby denying the public the full corpus of 20th century books and, perhaps, unknown benefits    |   |  |





# ПРИМЕРЫ КОРПУСОВ

Примеры иностранных корпусов

## Czech National Corpus (CNC)

Корпус состоит из собраний текстов на чешском языке. В их число входят современный чешский, разговорный язык, корпус исторических текстов и параллельный корпус, содержащий перевод на более 30 других языков.

Корпус доступен для всех, кто заинтересован в его исследовании, в нём также доступна возможность поиска интересующих словоформ.

# ПРИМЕРЫ КОРПУСОВ

## Примеры иностранных корпусов

### Czech National Corpus



Search engine demonstration:

Search for Czech words:

**Country:**

Czech Republic

**CLARIN Centre:**

Charles University

**Description**

A Corpus is a collection of texts in electronic form used for linguistic research, using provided with digital tools to allow searching, analysis and research. Users can use these tools to find words and collocations in their original contexts, and determine their frequency in the corpus. **The Czech National Corpus** (CNC) is an academic project focusing on building a large electronic corpus of mainly written Czech. The Institute of the Czech National Corpus (ICNC), Faculty of Arts, Charles University in Prague oversees the development of the CNC, including its use in teaching, and advancing the field of the corpus linguistics.



# ПРИМЕРЫ КОРПУСОВ

## Примеры иностранных корпусов

|   |  |  |         |  |
|---|--|--|---------|--|
| Hits: 94   i.p.m.: 0.77 (related to the whole corpus)   ARF: 44.19   Result is shuffled |  |  |         | 1 / 3  |
| Line selection: simple  |  |  |         |  |
| <input type="checkbox"/>  | K znamená Kennedy                                | žádost odmítneme , musíme mít prezidentův podpis ke zrušení Habeas             | Corpus  | pro tento případ . Francisi , jsi ochoten podepsat mi                    |
| <input type="checkbox"/>  | Vynálezci z Guslaru                              | mu dáme peníze , chytíme ho in flagranti , s                                   | corpus  | delicti v ruce . " Ráno strávili velice pilnou prací                     |
| <input type="checkbox"/>  | Vzkříšení Olympie                                | . Vázových maleb se však zachovalo takové množství , že                        | Corpus  | vasorum antiquorum , který je zachycuje , má od roku                     |
| <input type="checkbox"/>  | Sport magazín, č. 12/2008                        | slovní spojení " Jagr and Granger " . Baseballová pálka                        | CORPUS  | DELICTI Historie : Kolem roku 1850 , v počátcích baseballu               |
| <input type="checkbox"/>  | Lidové noviny, 14. 2. 2006                       | stabilním stavu . Oznámila to včera mluvčí nemocnice v texaském                | Corpus  | Christi , kam byl osmasedmdesátiletý Harry Whittington po incidentu v    |
| <input type="checkbox"/>  | Chladnokrevně                                    | zjištěné skutečnosti opravňovaly , napadl platnost rozsudku podle zásady habes | corpus  | u kansaského Nejvyššího soudu , který ho nedávno potvrdil .              |
| <input type="checkbox"/>  | Svět středověkých gest                           | Zatím bylo v hlavním ohnisku pozornosti věřících Tělo Páně ,                   | Corpus  | Christi , jemuž patřila i zvláštní zbožnost . Vznikla v                  |
| <input type="checkbox"/>  | Moderní řízení, č. 10/2006                       | útvary , " říká Ing. Radek Bělina , konzultant společnosti                     | Corpus  | Solutions a. s. , a autor tohoto článku a pokračuje                      |
| <input type="checkbox"/>  | Živě.cz, 6. 1. 2005                              | zástupci lokálních partnerů ze společností Barracudanetworks , S & T a         | Corpus  | Solution , kteří ve svých příspěvcích využijí zkušeností se specifickými |
| <input type="checkbox"/>  | Hold Katalánsku a Ohlédnutí za Španělskou válkou | Ve Španělsku neexistuje - alespoň se nepraktikuje - žádný habeas               | corpus  | a ve vězení vás mohou držet nepřetržitě několik měsíců ,                 |
| <input type="checkbox"/>  | Jedinec: Gen - prostředí - vývoj                 | se v důsledku prudkého snížení krevní hladiny lh aktivita buněk                | corpus  | luteum a kolem 12 . dne po ovulaci ustane .                              |
| <input type="checkbox"/>  | Encyklopedie mystiky II.                         | , v tomto textu popsáno není ; jiný text z                                     | Corpus  | Hermeticum je označí za " božské mlčení " . Připomeňme                   |
| <input type="checkbox"/>  | Dům z úsvitu                                     | šlo o hada . Že prej zabil pitomýho hada !                                     | Corpus  | delicti , drž se , vyhrožoval , že se promění                            |
| <input type="checkbox"/>  | Silnější než slova                               | je to tak . Jak potom vysvětlíte tu práci v                                    | Corpusu | ? " " Dělal jste tam totéž . EastPark je                                 |

# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

## Стандарты разметки

В ходе исследования разметок были обнаружены свои стандарты разметок и форматы представления для различных корпусов.

1. Стандарт BNC. Присутствие метаинформации, текст разбивается на предложения, предложения – на слова с частеречной разметкой.
2. Формат базы данных. Корпус разбивается на текстовые файлы. Для текста, словаря и метаинформации созданы отдельные таблицы.
3. Стандарт GATE Developer. Текст разбивается на три секции. Первая - метаинформация, текст с узлами, третья – аннотация узла (слова).



# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

## Стандарты разметки

4. Стандарт НКРЯ. Используется HTML-разметка, текст разбивается на абзацы, которые имеют определяющий их атрибут. У каждой лексической единицы есть свой тег, который их обозначает.
5. Стандарт OpenCorpora. Присутствует метаинформация, текст разбивается на параграфы, параграфы – на предложения, предложения – на токены. Токены описываются специальными атрибутами.

# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

## Форматы хранения корпусов

Исходя из просмотренных стандартов и форматов, можно выделить следующие форматы хранения корпусов:

1. Формат базы данных.
2. Формат плоских файлов.
3. Размеченный формат (XML, HTML, SGML).



# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

Форматы хранения корпусов

## Формат базы данных и плоских файлов

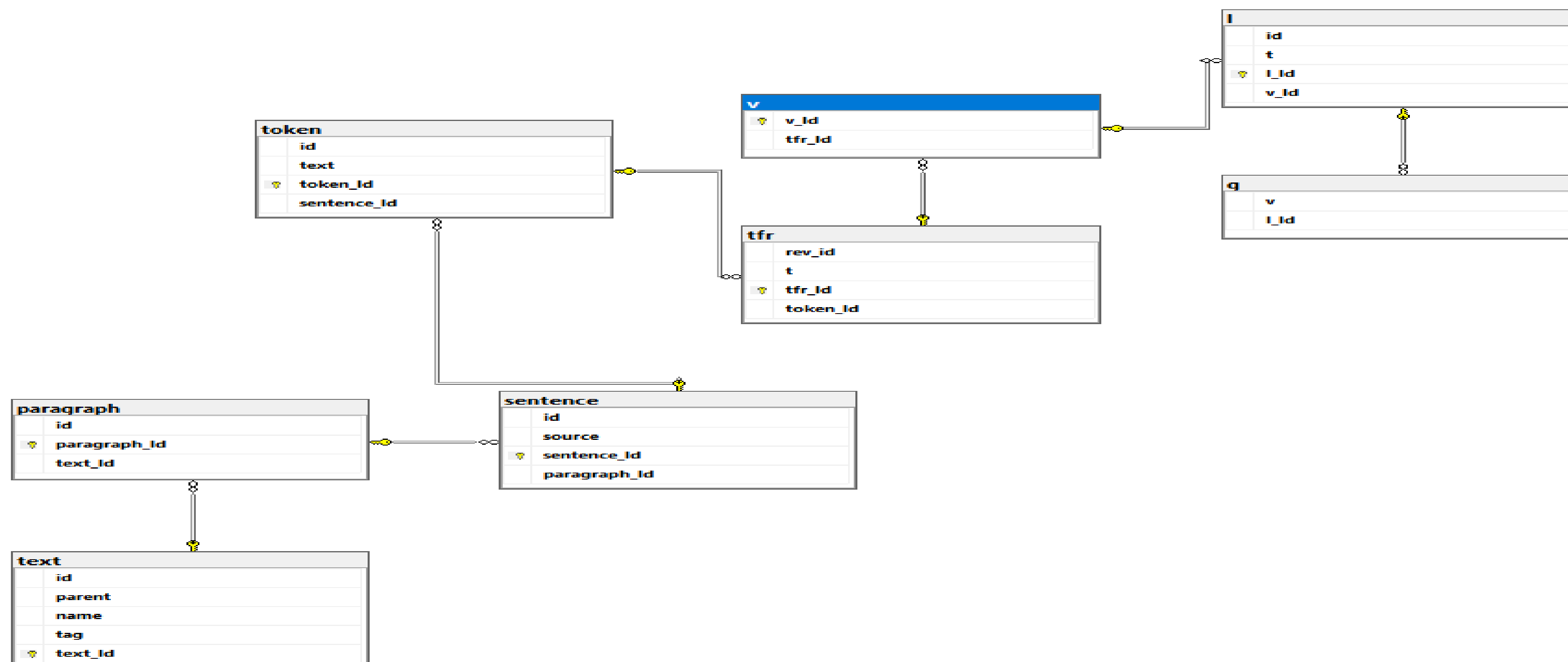
Оба формата довольно похожи, поскольку они обычно представляются в виде таблиц или файлов.

Каждый атрибут, токен, единица разметки корпуса могут быть представлены в виде отдельной таблицы. Между таблицами/файлами существует определённая связь, определяющая структуру корпуса.

В формате баз данных удобно находить слова при помощи SQL-запросов.

# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

Форматы хранения корпусов



# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

Форматы хранения корпусов

| 1 | id | parent | name                    | tag  | text_Id   | corpus_Id |   |  |
|---|----|--------|-------------------------|------|---|-----------|---|--|
| 2 | 1  | 0      | "Частный корреспондент" | url: | <a href="http://www.chaskor.ru">http://www.chaskor.ru</a> | 0         | 1 |  |
| 3 | 2  | 1      | 00021 Школа злословия   | Год: | 2008  | 1         | 1 |  |

# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

Форматы хранения корпусов

## Размеченный формат

Обычно формат представлен в виде XML-файлов, но есть варианты HTML и SGML-разметки.

Слова, предложения, метаинформация заключается в тэги. У тэгов могут быть атрибуты, которые более подробно описывают заключенную в них информацию. Например, для слов через атрибут может передаваться часть речи, падеж и прочие категориальные признаки.

# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

## Форматы хранения корпусов

```
1 <text id="int" parent="int" name="string">
2   <tags>
3     <tag>string</tag>
4   </tags>
5   <paragraphs>
6     <paragraph id="int">
7       <sentence id="int">
8         <source>string</source>
9         <tokens>
10          <token id="int" text="string">
11            <tfr rev_id="id" t="string">
12              <v>
13                <l id="int" t="string">
14                  <g v="string"/>
15                </l>
16              </v>
17            </tfr>
18          </token>
19        </tokens>
20      </sentence>
21    </paragraph>
22  </paragraphs>
23 </text>
```

# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

Трансформация размеченных и плоских файлов, экспорт и импорт

Суть экспорта или импорта – преобразование разметки в текст, базу данных или другой стандарт.

В целом трансформацию можно разделить на два класса, которые соответствуют определённым форматам:

1. XSL-преобразование.
2. XSD-схема.



# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

Трансформация размеченных и плоских файлов, экспорт и импорт

## XSL-преобразование

XSLT – преобразование файла одного формата или стандарта в другой при помощи xsl-файла, в котором прописываются правила трансформации.

Для xslt-файлов определён собственный синтаксис, в котором прописываются различные шаблоны входного файла.

На выходе получается файл или несколько файлов с той разметкой, которая определена правилом.

# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

Трансформация размеченных и плоских файлов, экспорт и импорт

```
1 <xsl:stylesheet version="1.0"
2 xmlns:xsl="http://www.w3.org/1999/XSL/Transform">
3
4     <xsl:output method="text" indent="yes"/>
5     <xsl:strip-space elements="*" />
6     <xsl:template match="/">
7         <xsl:text>id&#009;text</xsl:text>
8         <xsl:value-of select="'&#13;'" />
9         <xsl:apply-templates select="annotation/text/paragraphs/paragraph/sentence/tokens" />
10    </xsl:template>
11    <xsl:template match="token">
12        <xsl:value-of select="concat(@id, '&#009;', @text, '&#13;')" />
13    </xsl:template>
14
15 </xsl:stylesheet>
```

# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

Трансформация размеченных и плоских файлов, экспорт и импорт

## XSD-схема

Для трансформации файла с XML-разметкой можно воспользоваться XSD-схемой.

Она описывает структуру XML-файла, его тэги, атрибуты, типы данных и т.д.

Соответственно эта структура должна иметь такую же структуру, какая присутствует у базы данных.

Через комплексные типы образуются связи между таблицами.

# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

Трансформация размеченных и плоских файлов, экспорт и импорт

```
1 <xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
2   <xs:element name="text">
3     <xs:complexType>
4       <xs:sequence>
5         <xs:element name="tag" type="xs:string"/>
6         <xs:element name="paragraph">
7           <xs:complexType>
8             <xs:sequence>
9               <xs:element name="sentence">
10                <xs:complexType>
11                  <xs:sequence>
12                    <xs:element name="source" type="xs:string"/>
13                    <xs:element name="token">
```

# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

## ETL-процесс

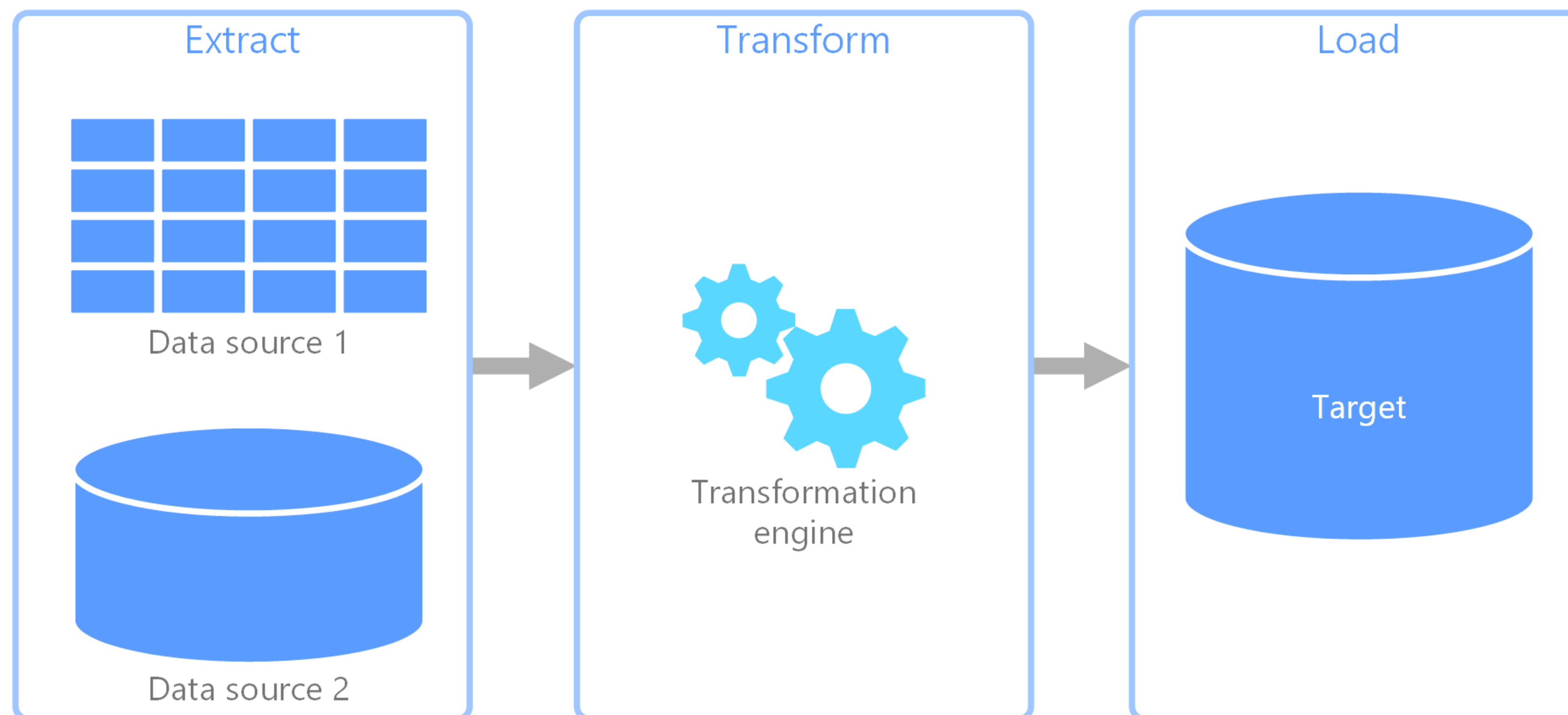
ETL - инструменты позволяют экспортировать или импортировать источники данных в назначения данных. Например, можно импортировать данные в excel-файле в базу данных.

ETL расшифровывается как Extract, Transform, Load. Подход можно описать следующим процессом:

1. Извлечение данных из источников.
2. Обработка, очистка данных, преобразование, интеграция.
3. Загрузка данных в назначение данных.

# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

## ETL-процесс







# СУЩЕСТВУЮЩИЕ ПОДХОДЫ

Примеры ETL-инструментов

ETL-процессы поддерживаются следующими инструментами:

1. SQL Server Integration Services (SSIS).
2. Oracle Data Integration.
3. GATE Developer.

# ПРЕДЛАГАЕМЫЙ ПОДХОД

## Основные моменты

Суть – объединение описанных подходов в решаемой проблеме.

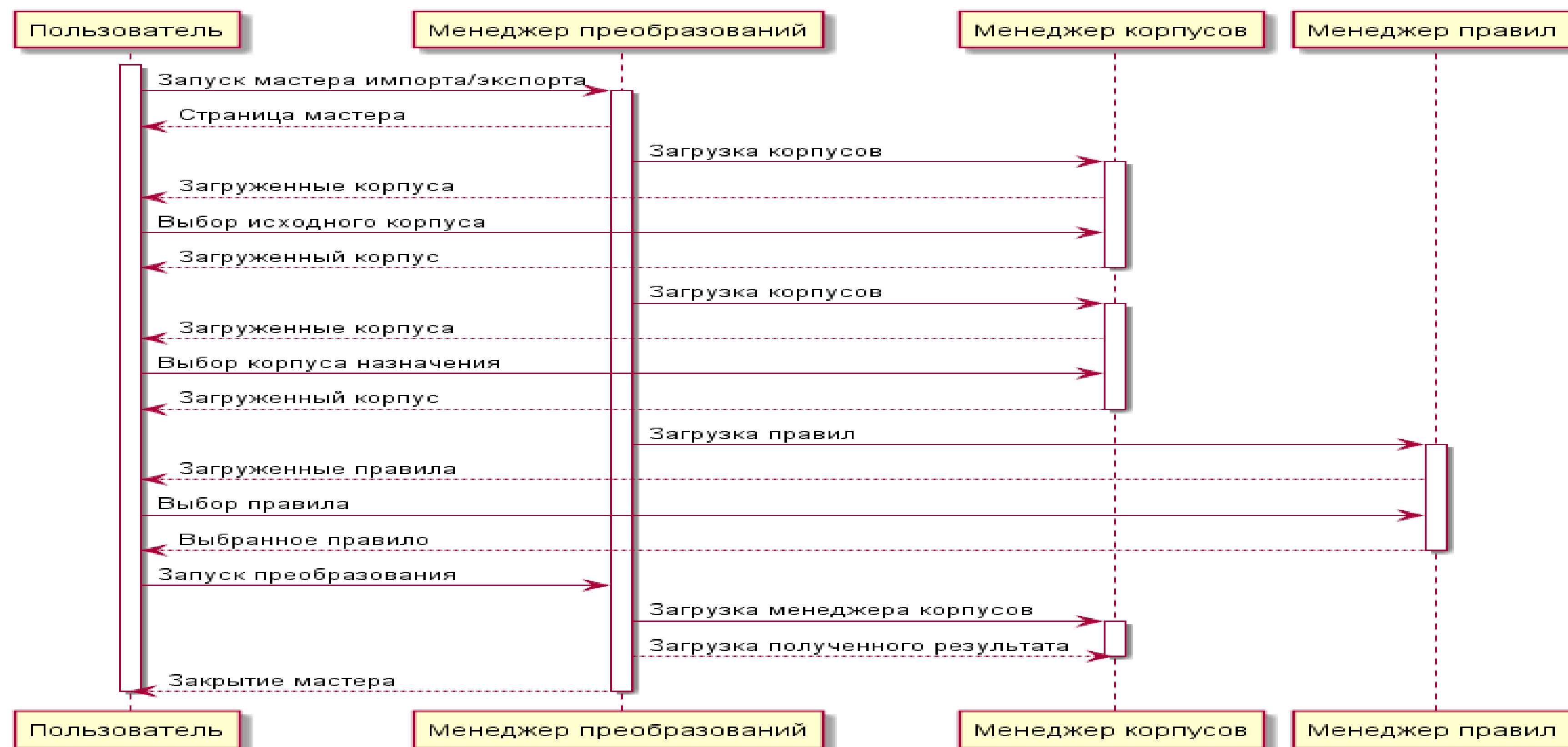
Для базы данных была разработана собственная схема, базирующаяся на стандарте OpenCorora. На нём же основывается XML-формат.

На основе ETL подхода разработан основной процесс импорта и экспорта корпусов.

Для каждой части процесса написан свой модуль.

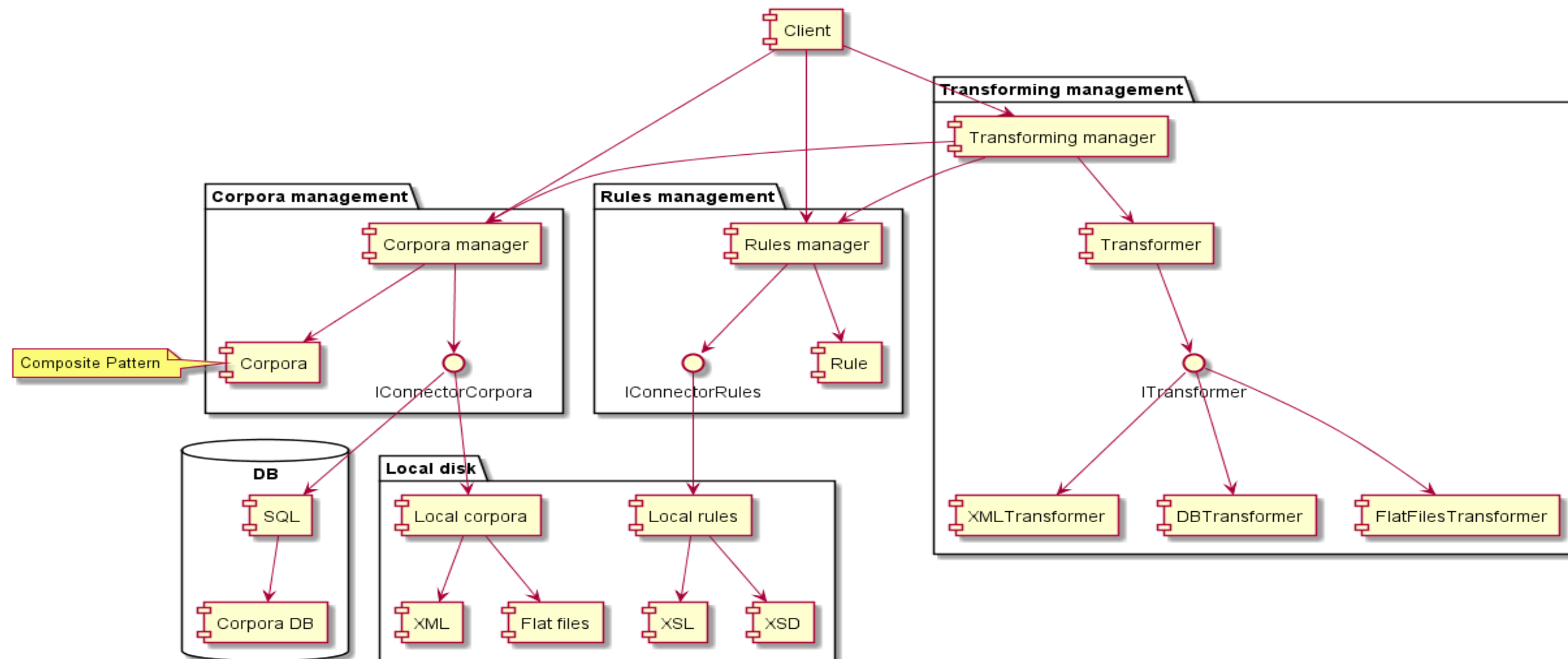
# ПРЕДЛАГАЕМЫЙ ПОДХОД

## Процесс экспорта и импорта



# ПРЕДЛАГАЕМЫЙ ПОДХОД

## Архитектура системы





# ПРЕДЛАГАЕМЫЙ ПОДХОД

Интерфейс системы. Менеджер корпусов

Менеджер корпусов

Менеджер правил

Загрузить корпус

Добавить корпус

Удалить корпус

Редактировать корпус

Добавить файл

Удалить файл

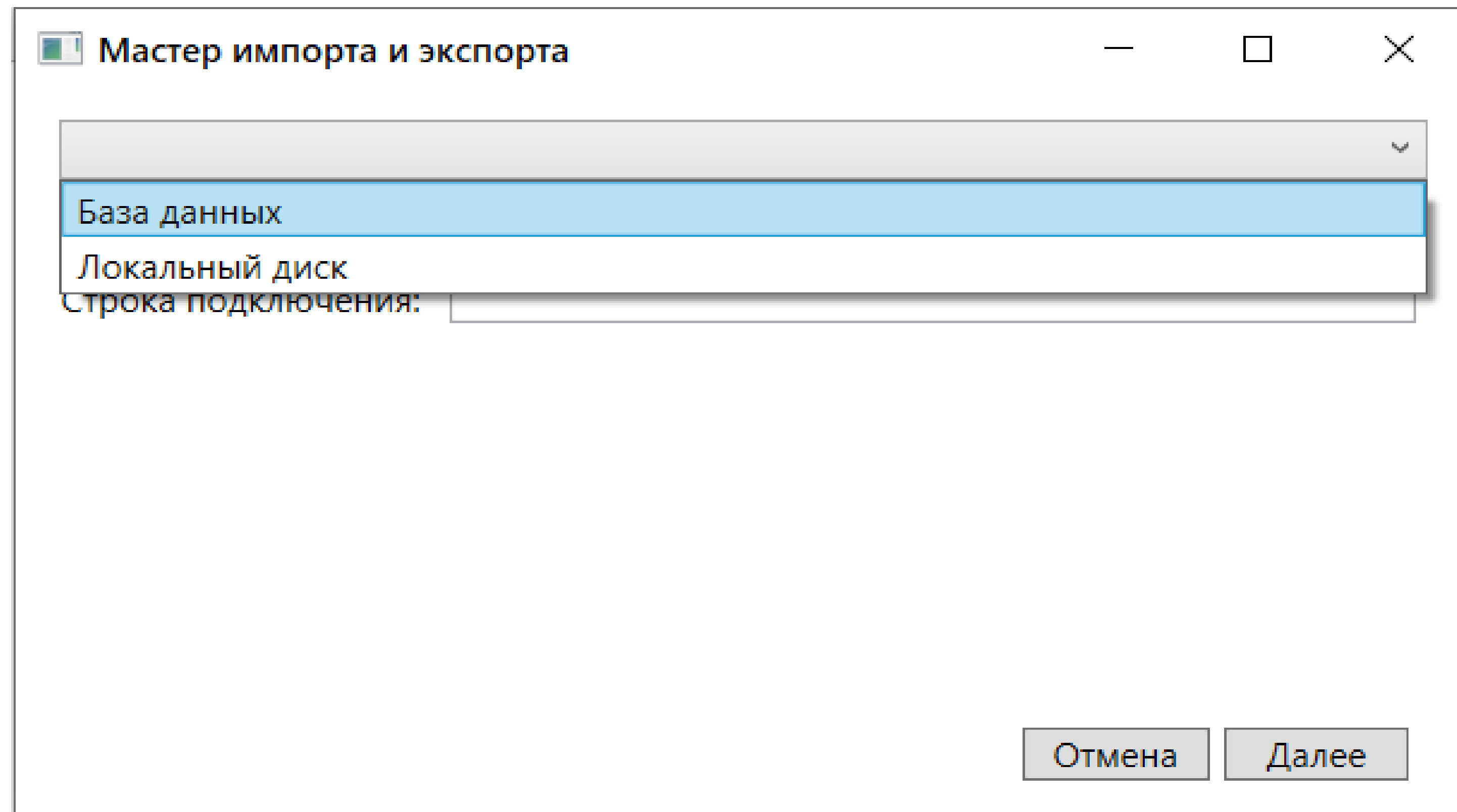
Редактировать файл

Преобразовать



# ПРЕДЛАГАЕМЫЙ ПОДХОД

Интерфейс системы. Мастер импорта и экспорта. Выбор типа хранилища





# ПРЕДЛАГАЕМЫЙ ПОДХОД

Интерфейс системы. Мастер импорта и экспорта. Выбор исходного корпуса

Мастер импорта и экспорта

Выберите исходный корпус:

DBCorpora

gdfg

InititalCorpus

OtherStore

Отмена

Назад

Далее





# ПРЕДЛАГАЕМЫЙ ПОДХОД

Интерфейс системы. Мастер импорта и экспорта. Выбор правила

Мастер импорта и экспорта

Выберите правило:

- Rule.xsl
- rule1.xsd

Отмена Назад Далее



# ПРЕДЛАГАЕМЫЙ ПОДХОД

Интерфейс системы. Мастер импорта и экспорта. Последняя стадия процесса преобразования

Мастер импорта и экспорта

Корпус:  
DBCorpora

Затронутые файлы:  
opcTest.xml  
opcTest1.xml

Правило преобразования:  
Rule.xsl

ОтменаНазадПреобразовать

# ЧТО ДАЛЬШЕ

Проблемы и будущие наработки

Проблемы:

1. Импорт в базу данных работает довольно долго.
2. Каждый файл означает текст корпуса. Некоторые системы аннотирования хранят файлы с метаданными (например, BRAT).

Что планируется:

Описанная система станет частью системы управления хранилищем корпусов на основе использования NoSQL базы данных.



НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ

Спасибо за внимание!