



НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ

Лаборатория интернет-исследований

# ПРЕПРОЦЕССИНГ РУССКОЯЗЫЧНЫХ ТЕКСТОВ В

Кольцов С.Н

Санкт-Петербург, 2018



## ПРЕПРОЦЕССИНГ

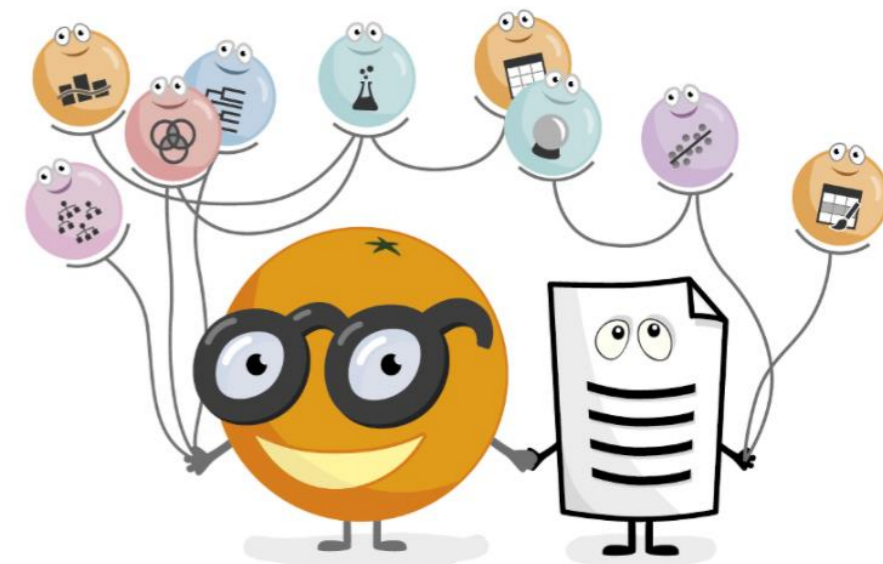
Препроцессинг состоит из следующих этапов

1. **Загрузка данных из источника данных**, например csv файл.
2. **Процедура лематизации**. Данная процедура заключается в процессе токинизации текстов.
3. **Удаление стоп слов**. Некоторые слова являются общими и при этом часто встречаются, поэтому их необходимо удалять.
4. **Визуализация текстовых данных в виде облака слов**. Наиболее частотные слова можно представить в виде облака слов, что упрощает некоторое понимание того какие темы внутри коллекции документов.

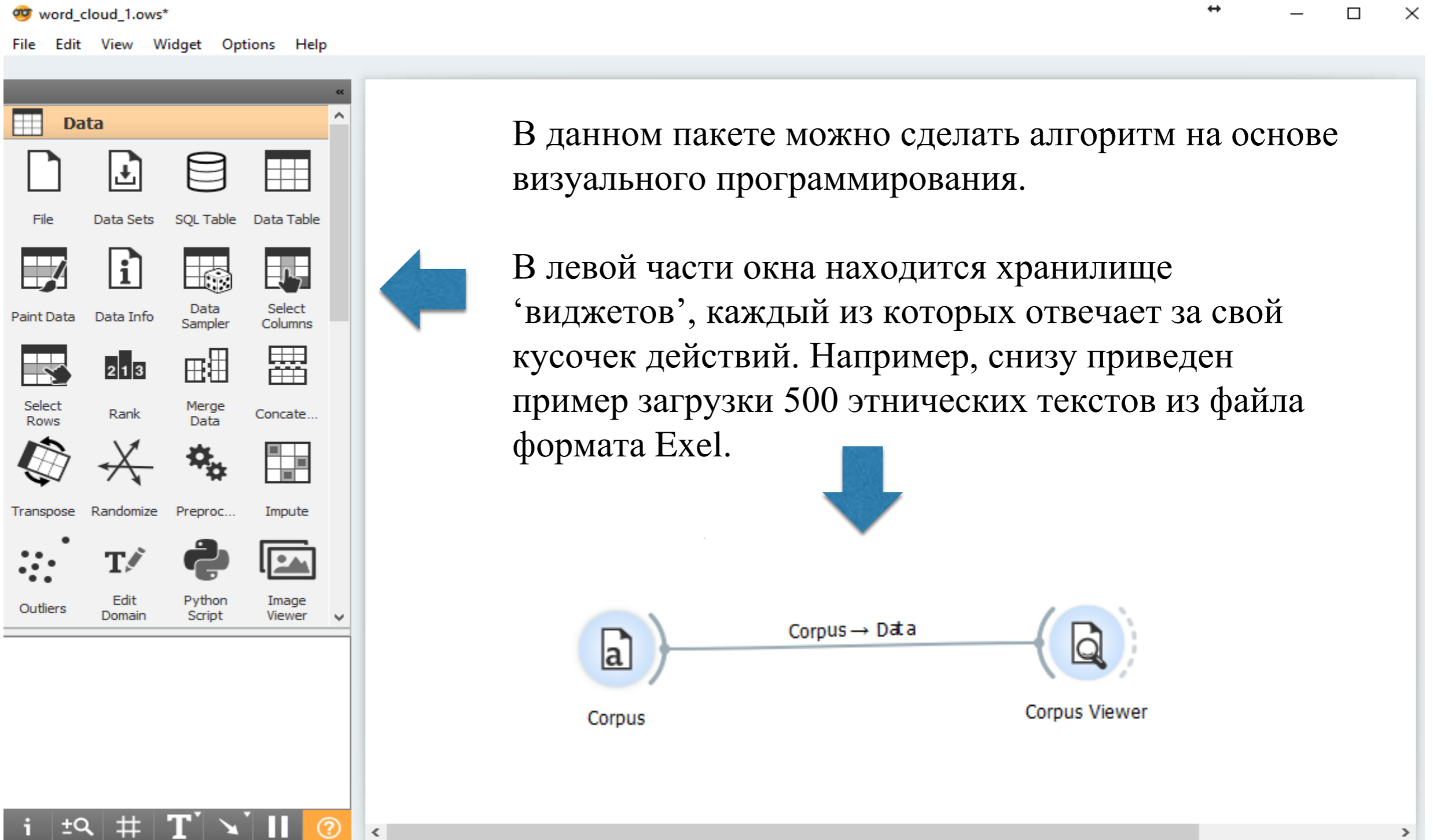
В дальнейшем использовать пакет Orange.

**Данный пакет можно скачать по адресу:**

<https://orange.biolab.si/>



## ORANGE



The screenshot shows the ORANGE software interface. The title bar reads "word\_cloud\_1.ows\*". The menu bar includes "File", "Edit", "View", "Widget", "Options", and "Help". On the left, a "Data" widget palette is visible, containing various widgets such as File, Data Sets, SQL Table, Data Table, Paint Data, Data Info, Data Sampler, Select Columns, Select Rows, Rank, Merge Data, Concatenate, Transpose, Randomize, Preproc..., Impute, Outliers, Edit Domain, Python Script, and Image Viewer. The main workspace on the right contains a workflow diagram with two widgets: "Corpus" (represented by a document icon with 'a') and "Corpus Viewer" (represented by a document icon with a magnifying glass). An arrow labeled "Corpus → Data" connects the two widgets. A blue arrow points from the text to the widget palette, and another blue arrow points from the text to the workflow diagram.

word\_cloud\_1.ows\*

File Edit View Widget Options Help

Data

File Data Sets SQL Table Data Table

Paint Data Data Info Data Sampler Select Columns

Select Rows Rank Merge Data Concatenate

Transpose Randomize Preproc... Impute

Outliers Edit Domain Python Script Image Viewer

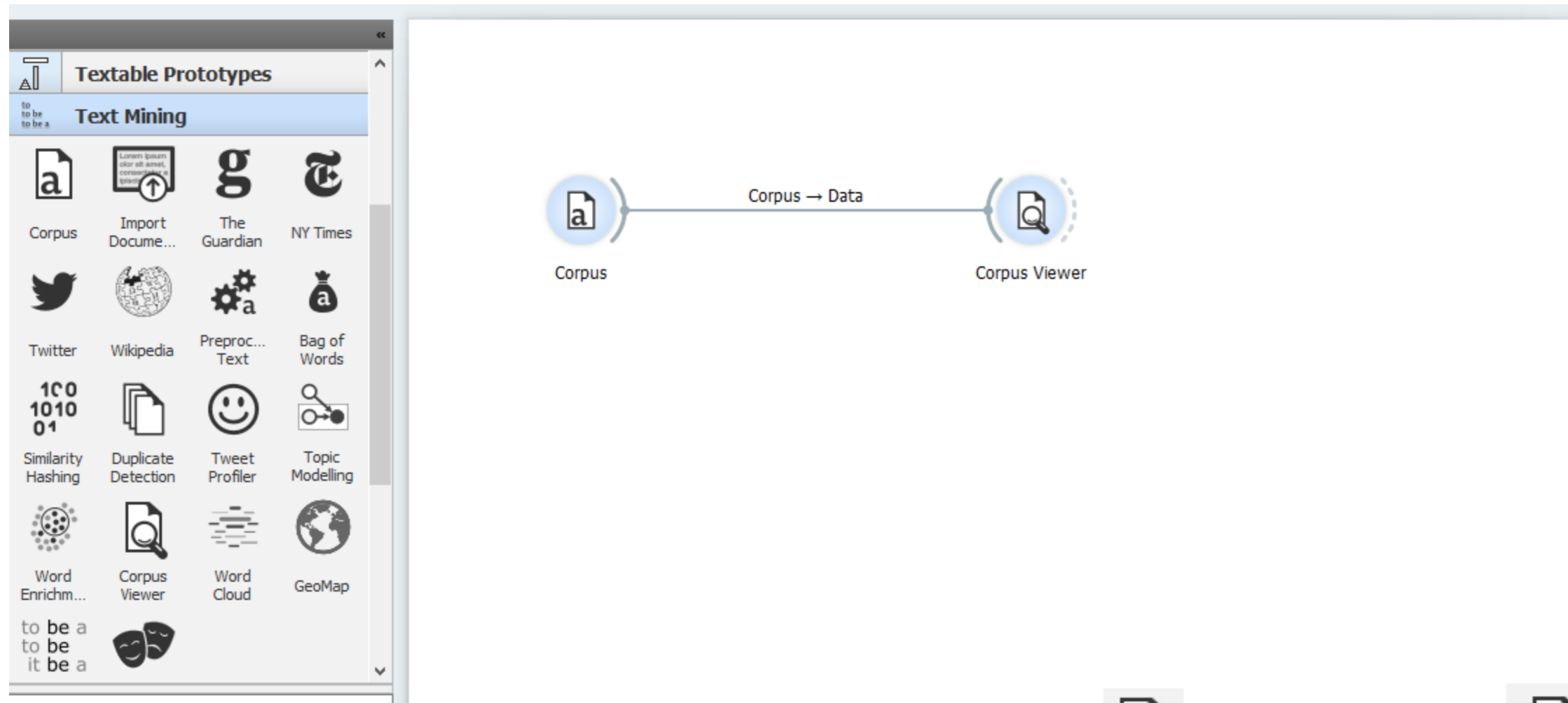
В данном пакете можно сделать алгоритм на основе визуального программирования.

В левой части окна находится хранилище ‘виджетов’, каждый из которых отвечает за свой кусочек действий. Например, снизу приведен пример загрузки 500 этнических текстов из файла формата Excel.

Corpus → Data

Corpus Corpus Viewer

## ЗАГРУЗКА РУССКОЯЗЫЧНЫХ ДАННЫХ В ORANGE





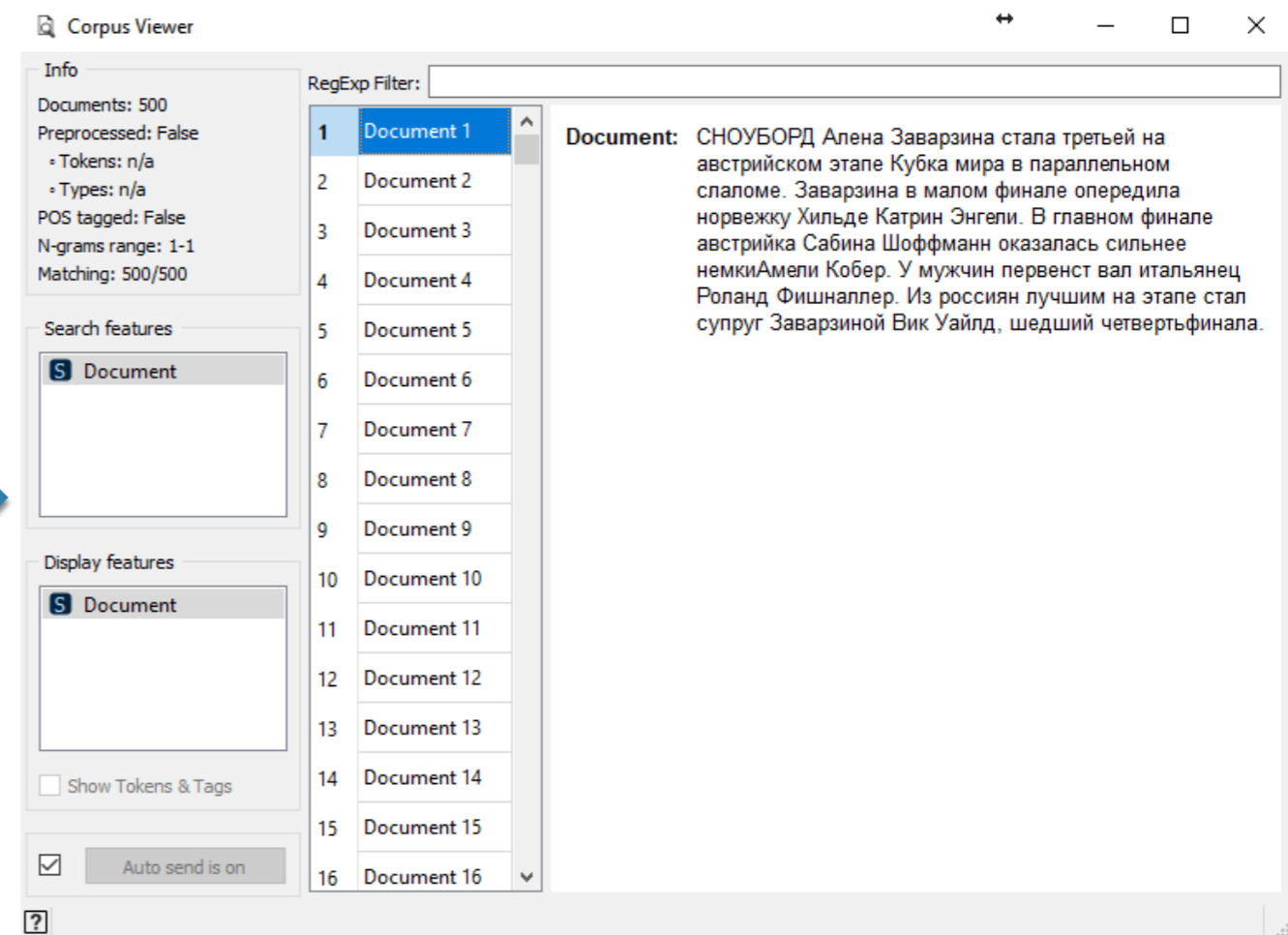
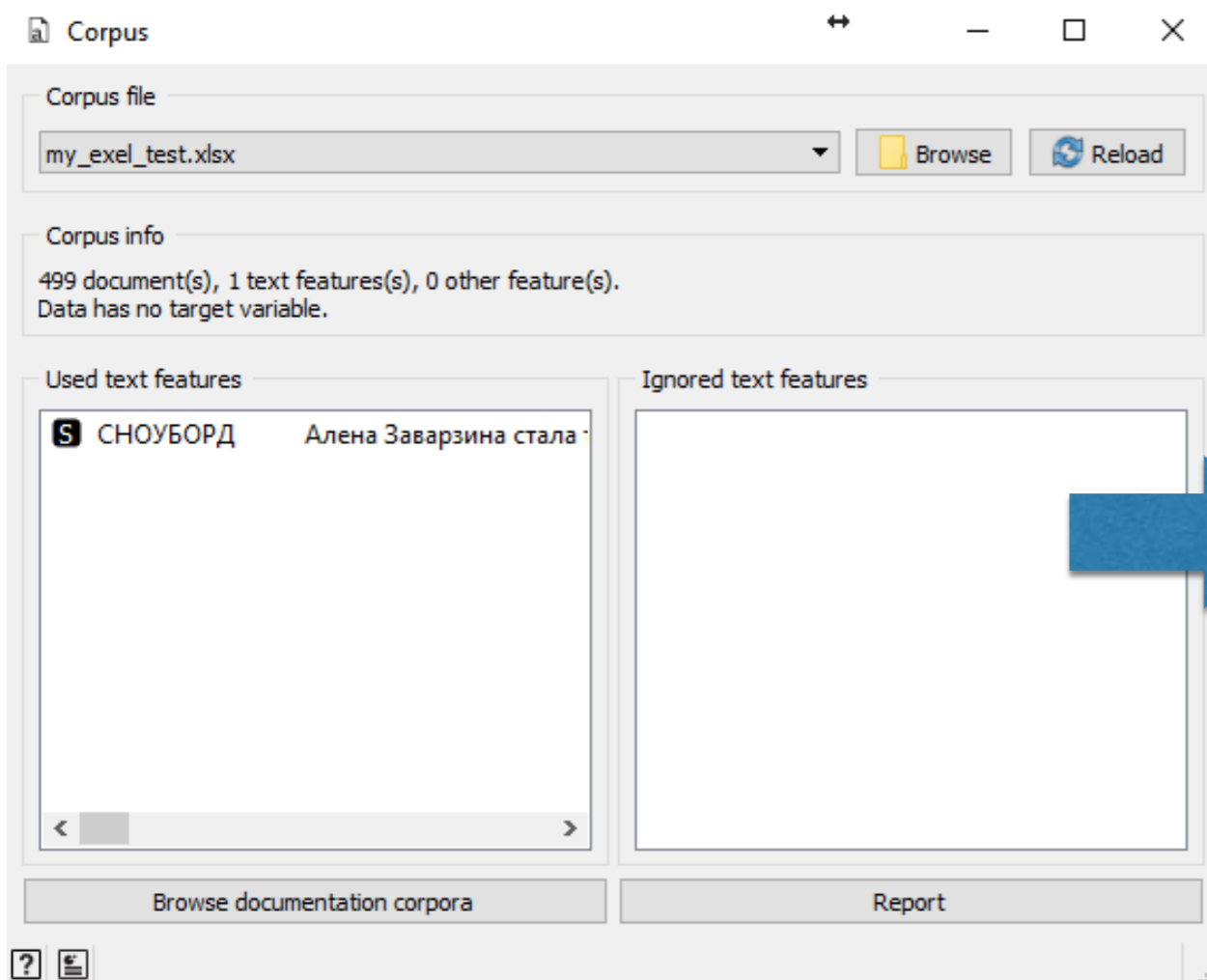
На вкладке ‘Text Mining’ нужно взять виджеты: ‘Corpus’  ‘Corpus Viewer’   
Далее нужно соединить два виджета линией.

Схема загрузки готова, теперь нужно лишь указать имя файла для загрузки и просмотра.



## ЗАГРУЗКА И ПРОСМОТР РУССКОЯЗЫЧНЫХ ТЕКСТОВ

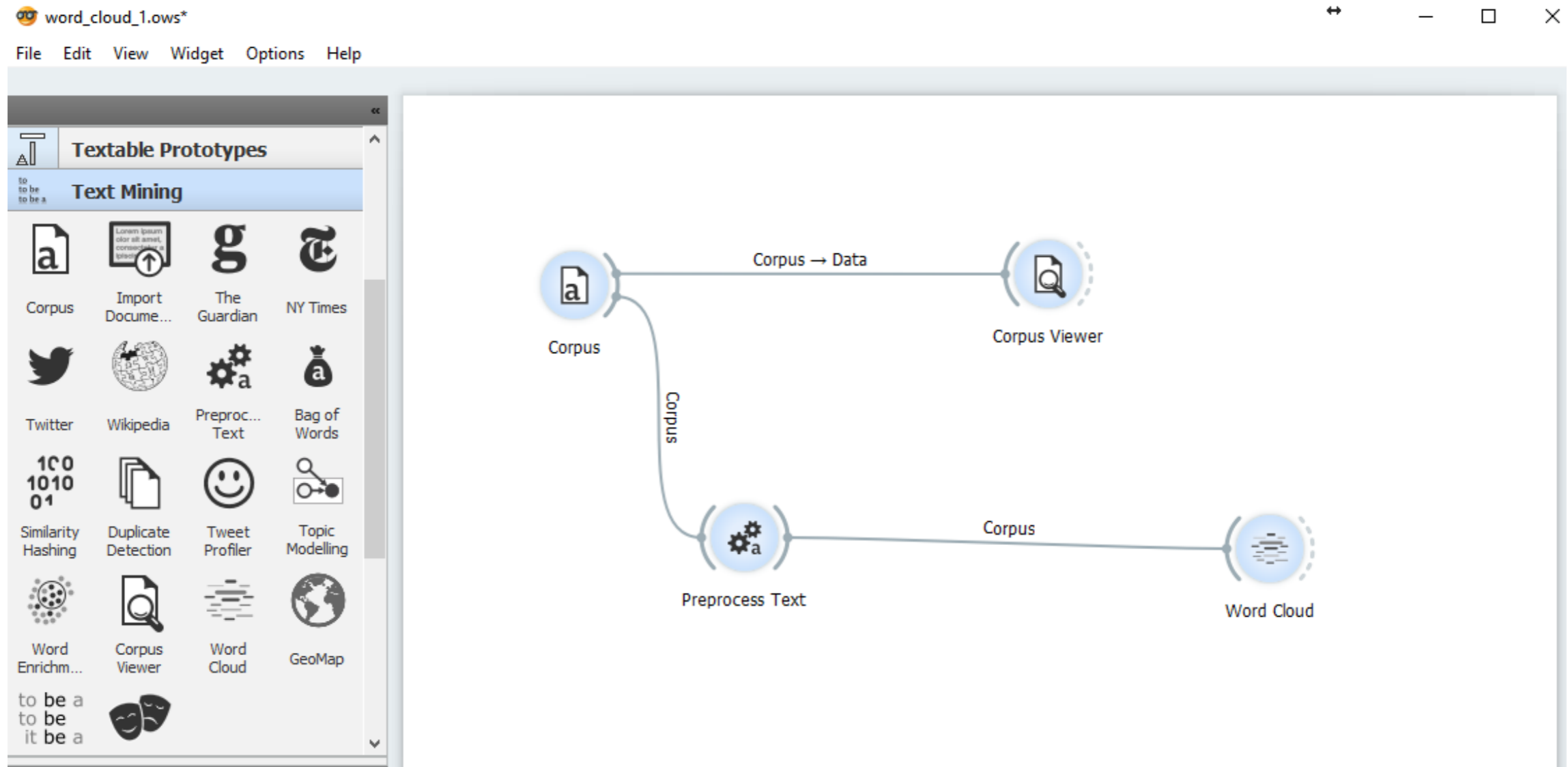
Для того что бы загрузить данные нужно кликнуть на виджете ‘Corpus’ и указать имя файла в формате книги Excel.



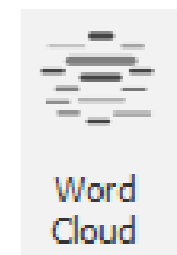
Для того что бы просмотреть результаты загрузки нужно кликнуть на виджете ‘Corpus Viewer’.



# ПРЕПРОЦЕССИНГ И ОБЛАКО ЛЕМАТИЗИРОВАННЫХ СЛОВ



Для удаления стоп слов, пунктуации (и так далее), а также для построения облака слов, достаточно добавить еще два виджета:





# ПРЕПРОЦЕССИНГ В ORANGE

Preprocess Text

Info

Document count: 500  
Total tokens: 21303  
Total types: 10017

Transformation

Lowercase     Remove accents     Parse html     Remove urls

Tokenization

Word & Punctuation  
 Whitespace  
 Sentence  
 Regexp    Pattern:   
 Tweet

Normalization [disabled]

Filtering

Stopwords    Russian    stop\_words.txt  
 Lexicon    (none)  
 Regexp    \.,!;:;!@?|X|V|N|+|\_|'|\"/>  
 Document frequency    0.10    0.90  
 Most frequent tokens    100

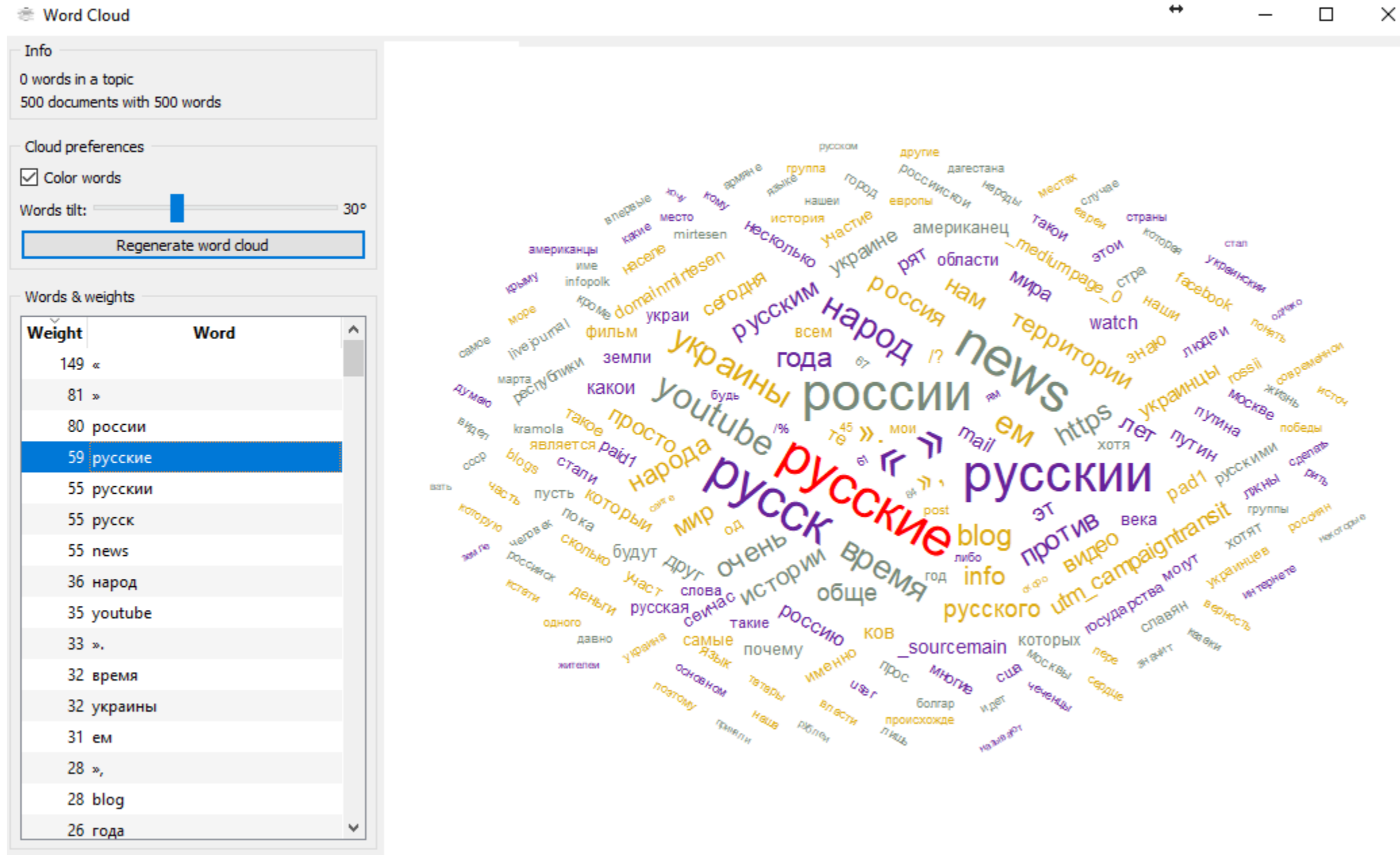
N-grams Range [disabled]

POS Tagger [disabled]

Commit Automatically

1. Можно указать параметры препроцессинга
2. Указать файл со списком стоп слов.
3. Указать язык.

# ОБЛАКО НЕЛЕМАТИЗИРОВАННЫХ СЛОВ



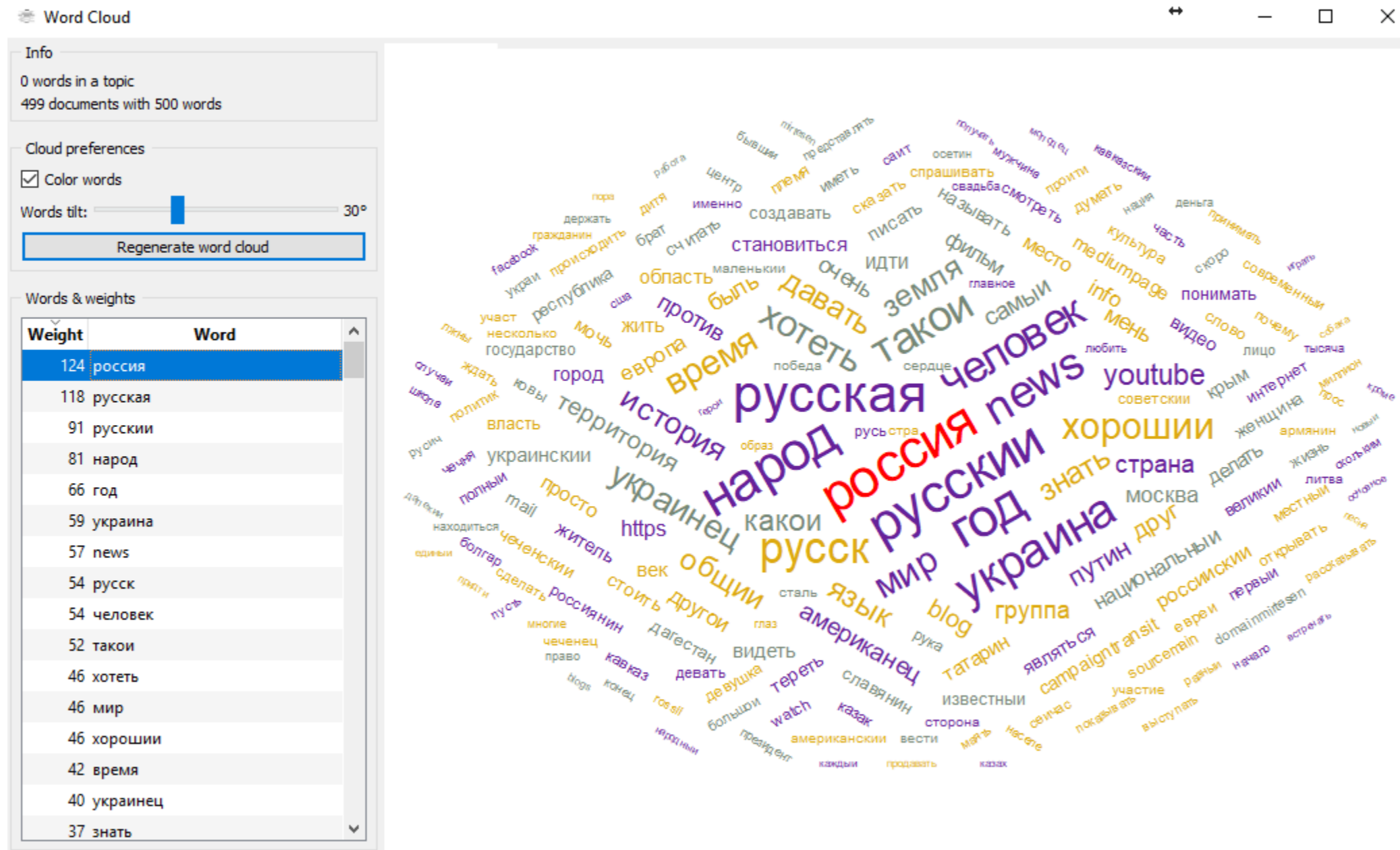
Таким образом, за счет списка стоп слов можно построит нормальное облако слов.





## TOPICMINER - MYSTEM

Если сделать препроцессинг в TopicMiner, и выгрузить лематизированные данные в csv формате, и загрузить в Orange, то получится наиболее качественный вариант облака.





НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ

<https://linis.hse.ru/>

Phone: +7 (911) 981 9165

Email: [skoltsov@hse.ru](mailto:skoltsov@hse.ru)