

Метод обработки естественного языка в условиях слабой структурированности исходного текста

Кобук М.Г.

Атаева О.М.

Тучкова Н.П.

Теймуразов К.Б.

03 декабря 2025

- Введение
- Цель и область исследования
- Решение::Транслятор
- Решение::Поиск
- Дальнейшие шаги

Базы знаний

- Δ Единое представление
- Δ Структурированный формат
- Δ Детерминированный подход
- Δ Автоматизация

Входные данные

- Δ Разное представление
- Δ Разные типы
- Δ Большие объёмы

Обработка

- Δ Онтология
- Δ Обучение моделей
- Δ Заполнение каталогов
- Δ Полнотекстовый поиск

Требования

- Δ Быстро
- Δ Качественно
- Δ Легковесно

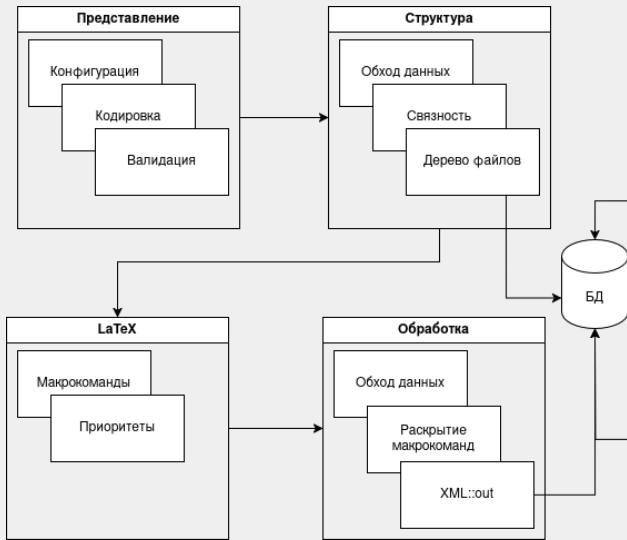
- В рамках проекта SciLibRu
- Выпуски периодических изданий
 - ▶ Многофайловый состав
 - ▶ Пользовательские макрокоманды
 - ▶ Нестабильный формат
 - ▶ Опора на редкие/устаревшие библиотеки
 - ▶ Нестандартные кодировки
- Отдельно стоящие файлы

! Нужно универсальное, быстро адаптируемое решение

Сравнение функций

	tex2xml	LaTeXML	TeX4ht	t2x3
Сложные файлы	—	+	+	+
Многофайловые док-ты	—	—	—	+
Простая конфигурация	*	—	+/-	+/-
Пользовательские макро	—	+	+/-	+
Стабильность	+	+	—	—
Строгий шаблон	*	+	+/-	-/+
Нестандартные файлы	—	—	—	+

Описание подхода [Трансляция]



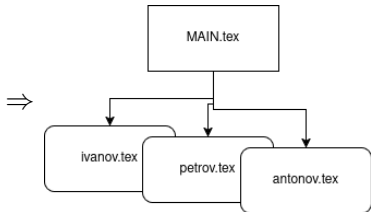
Представление и структура

Кодировка

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
```

Дерево

```
journal/
├── antonov.tex
├── ivanov.tex
├── MAIN.tex
└── petrov.tex
```



- Регистрация макрокоманд в два прохода
- Чёрный список и приоритет макрокоманд
- Итеративное раскрытие макрокоманд

Не всё надо раскрывать

```
\Abst{АННОТАЦИЯ}

\begin{center}\small\nwt
\parbox{150mm}{%\baselineskip=2.5ex
\textbf{Аннотация:}\ \
АННОТАЦИЯ}\end{center}}
```


1. Лексер

- ▶ Текст
- ▶ LaTeX-команда
- ▶ Управляющий символ
- ▶ [Формула]

2. Парсер

- ▶ Итеративное раскрытие макрокоманд согласно конфигурации
- ▶ Разные режимы строгости обработки ввода

3. Построение и обход AST

4. Оценка качества ($\text{integrity} < 0.1$; $\text{quality} > 2.75$)

$$\text{integrity} = \frac{\text{TeX}_{\text{xml}}}{\text{TeX}_{\text{tex}}}; \text{quality} = \frac{\text{Heads}_{\text{xml}}}{\text{Heads}_{\text{tex}}} + \frac{\text{Tails}_{\text{xml}}}{\text{Tails}_{\text{tex}}} + \frac{\text{Formula}_{\text{xml}}}{\text{Formula}_{\text{tex}}}$$

Результаты обработки тестовой выборки (1010 файлов)



- 89 технических/нецелевых файлов исключены из общего числа (1099)
- В $\approx 12\%$ файлов присутствуют ошибки структуры выходных данных
- В $\approx 33\%$ файлов наблюдаются остаточные конструкции LaTeX

1. Лексер/Парсер

- ▶ Расширение покрытия LaTeX-примитивов и встроенных команд
- ▶ Рекурсивный спуск \rightarrow LALR
- ▶ Автомат состояний

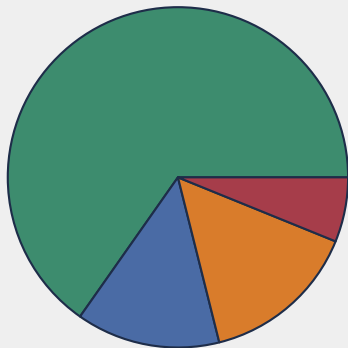
2. Анализ структуры документа

- ▶ Упрощение конфигурации

3. Индикация ошибок

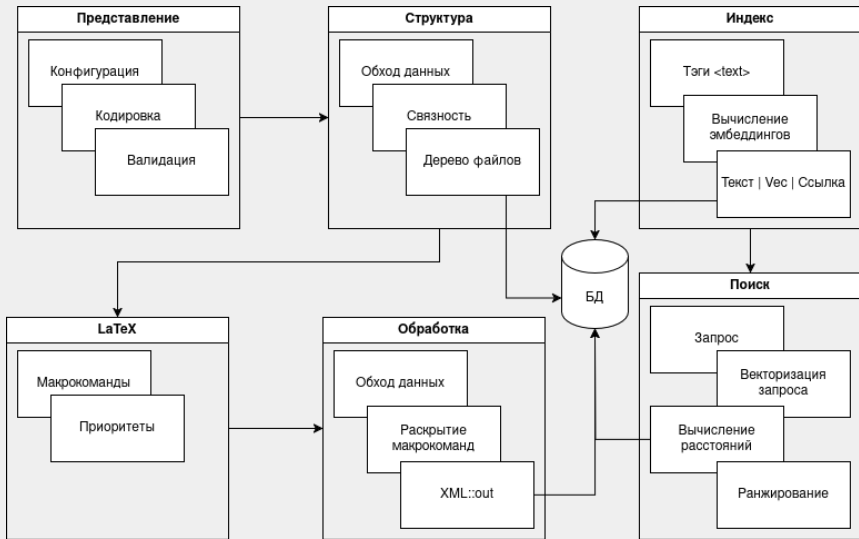
- ▶ Расширение логирования

Результаты доработки



- Без ошибок (65.2%)
- Доработка конфигурации (13.4%)
- Фатальные ошибки (15%)
- Расширение покрытия TeX (6.1%)

Описание подхода



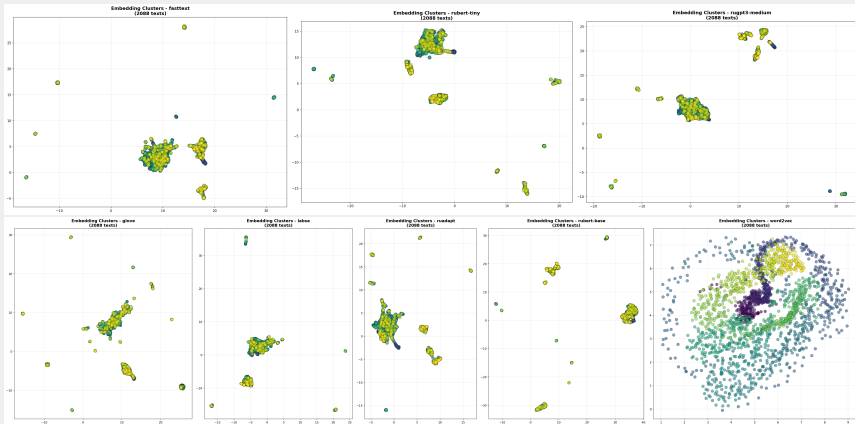
■ Модели

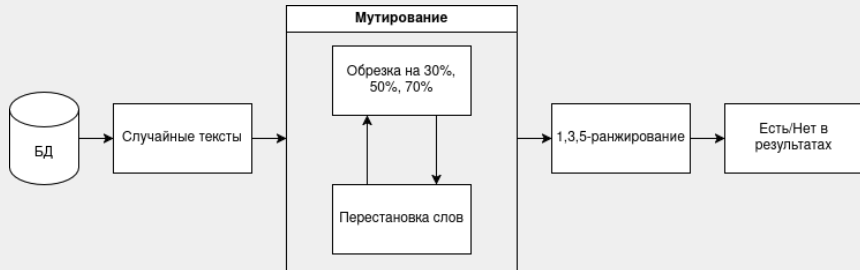
- ▶ fasttext
- ▶ rubert-tiny
- ▶ rugpt3-medium
- ▶ glove
- ▶ labse
- ▶ ruadapt
- ▶ rubert-base
- ▶ word2vec
- ▶ (tf-idf)
- ▶ (elmo)

■ Метрики

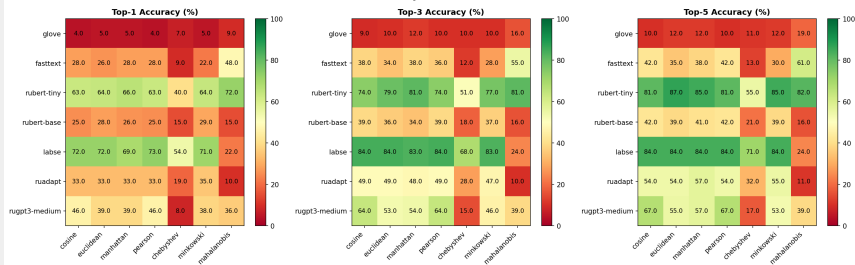
- ▶ Косинусное сходство
- ▶ Евклидово расстояние
- ▶ Манхэттенское расстояние
- ▶ Расстояние Пирсона
- ▶ Мера Чебышёва
- ▶ Мера Минковского
- ▶ Мера Махаланобиса

Распределение векторов





Search Quality Benchmark Results



1. Расширение покрытия TeX, доработка конфигурации
2. Интеграция с моделями для семантической верификации
3. Дообучение моделей на корпусе для повышения точности поиска
4. Двухуровневая модель поиска
5. Тестирование ELMo

Спасибо за внимание



Лирическое отступление



Лирическое отступление

Вот таблица со списком памятников и их наименованиями. Проанализируй таблицу и вставь слово ран/РАН там, где требуется. Будь внимателен к регистру.

Название памятника Местоположение

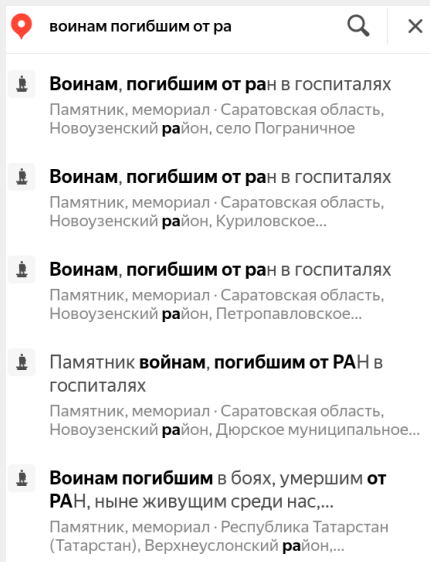
Кладбище советских воинов, погибших от ? в госпиталях Новосибирска (1941-1945) г.
Новосибирск

Мемориал солдатам, погибшим от ? в госпиталях г. Казань



Кладбище советских воинов, погибших от ран в госпиталях Новосибирска (1941-1945) г. Новосибирск
Мемориал солдатам, погибшим от РАН в госпиталях г. Казань

Лирическое отступление



Объективное отступление

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-z^2/2} dz, \quad x \in \mathbb{R}.$$

```
<formula id="6">\Phi(x)=\displaystyle\frac{\displaystyle 1\mathstrut}{\displaystyle \sqrt{2\pi}\mathstrut}\int\limits_{-\infty}^xe^{-z^2/2}\,dz\,,\,\\ \enskip x\in\mathbb{R}\,.\</formula>
```

Объективное отступление

```
<section id="1" name="Introduction. The Weibull distribution">
```

```
<text>
```

```
<value>In probability theory and mathematical statistics, the Weibull  
distribution concentrated on the nonnegative half-axis, whose tail decreases in an
```

```
<value>It is named after the Swedish scientist Waloddi Weibull (1887  
in the analysis of the strength of materials and described and studied it in de  
his distribution in describing many statistical patterns.</value>
```

```
<value>Let  $\gamma > 0$ . The Weibull distribution with shape parameter  $\gamma$  is t
```

```
<formula id="1">\sf P\left(W_{\gamma}<x\right)= \left[1-e^{-x^{\gamma}}\right]  
R\,,.</formula>
```


Here and galee,

■ поверхность шар

1. Таким образом, в однодипольной модели обратная задача решается следующим способом. Нужно найти **на поверхности сферы две точки с** максимальным значением радиальной компоненты магнитного поля...
2. т. е. функция ширины экрана непрерывна и в окрестности точки...
3. меньшим значениям модуля градиента на рис. 5 соответствуют более светлые серые пиксели. Так, при наименьшем значении модуля градиента (...)
4. во всей плоскости...
5. Расчет количества рабочих мест, которые могут быть размещены в области комфортного наблюдения, должен опираться на способ их размещения. Например, для построения ситуационного зала, где в центре будет **находиться овальный стол**, расчет будет заключаться в определении максимального размера такого стола...

■ Распределение в пределе

1. используем предел...
2. распределение Вейбулла называется распределением Рэля
3. На основе полученных результатов была написана программа, позволяющая вычислять распределения времен выхода из множества состояний, совместное стационарное распределение числа заявок в системе и состояния системы и связанные с ним характеристики, а также исследовать поведение рассматриваемой СМО в з...
4. (т. е. «изначально» **предельным распределением минимальной порядковой статистики должно быть экспоненциальное**), но объем выборки имеет вид...

■ Распределение в пределе

1. распределение Вейбулла называется распределением Рэлея...
2. (т. е. «изначально» предельным распределением минимальной порядковой статистики должно быть экспоненциальное), но объем выборки имеет вид...
3. При этом аналогичная нижняя доверительная граница, вычисляемая указанным выше упрощенным методом (Ллойда и Липова), основанным на непосредственном использовании частных доверительных границ для параметров отдельных элементов, равна