

Томи́та-парсер

Технологии Яндекса

<https://tech.yandex.ru/tomita/> Система, предназначенная для извлечения структурированных сущностей из текста: фактов, отношений, причинно-следственных связей и др.

- Использует механизм контекстно-свободных грамматик и словари ключевых слов.
- Позволяет добавлять свои расширения.

исходный код:

<https://github.com/yandex/tomita-parser/>

Основные понятия

- Газеттир — словарь ключевых слов
(пример статьи: «все города России»)
- Грамматика — множество правил (шаблонов) на языке КС-грамматик
- Факты — таблицы с колонками (полями)

Алгоритм работы парсера

- 1 Найти вхождения всех ключей из газеттира
- 2 Найти ключи, которые есть в грамматике (kwtype)
- 3 Покрыть предложение непересекающимися ключами
- 4 Отобразить терминалы грамматики на входные слова
- 5 Интерпретировать на построенном синтаксическом дереве

Пример:

$$S \rightarrow Noun$$

`t = 'механизм контекстно-свободных грамматик'`

`...`

`['механизм', 'грамматика']`

Простейшие правила:

`S -> Adj Noun`

`S -> Adj Word<h-reg1>`

`S -> Adj<gnc-agr[1]> Noun<gnc-agr[1]>`

`S -> A B C | A B* C+`

`S -> Adj "слово";`

Список всех помет:

<https://tech.yandex.ru/tomita/doc/dg/concept/all-labels-list-docpage/>

Начало работы

- Откуда скачать: <https://tech.yandex.ru/tomita/>
- Как запустить:
<https://tech.yandex.ru/tomita/doc/dg/concept/run-parser-docpage/>
- Как запустить на Mac iOS:

```
chmod a+x tomita-mac  
./tomita-mac config.proto
```

Файлы проекта

config.proto — конфигурационный файл парсера.
Сообщает парсеру, где искать все остальные файлы, как их интерпретировать и что делать.

dic.gz — корневой словарь. Содержит перечень всех используемых в проекте словарей и грамматик.

mygram.cxx — грамматика

facttypes.proto — описание типов фактов

kwtypes.proto — описания типов ключевых слов

Файл mydic.gzt – корневой словарь:

```
encoding "utf8";  
import "base.proto";          // описания protobuf-типов (TAuxDicA  
rticle и прочих)  
import "articles_base.proto"; // Файлы base.proto и articles_base.  
proto встроены в компилятор.  
                                // Их необходимо включать в начало л  
юбого gzt-словаря.  
// статья с нашей грамматикой:  
TAuxDicArticle "наша_первая_грамматика"  
{  
    key = { "tomita:first.cxx" type=CUSTOM }  
}
```

Описание простейшей грамматики:

```
#encoding "utf-8"
```

```
#GRAMMAR_ROOT S
```

```
S -> Noun;
```

Сохраняется в специальный файл (first.cxx)

Файл config.proto:

```
encoding "utf8";
TTextMinerConfig {
    Dictionary = "mydic.gzt"; // путь к корневому словарю
    PrettyOutput = "PrettyOutput.html"; // путь к файлу с отладочным
    выводом в удобном для чтения виде
    Input = {
        File = "test.txt"; // путь к входному файлу
    }
    Articles = [
        { Name = "наша_первая_грамматика" } // название статьи в корнев
    ом словаре,
                                                // которая содержит запус
    каемую грамматику
    ]
}
```

Осталось сделать файл test.txt и проверить:

запуск:

```
./tomitaparser config.proto
```

Вывод печатается в файл PrettyOutput.html.

Труд облагораживает человека . EOS

| Text | Type |
|-------------------------|---|
| труд | TAuxDicArticle [наша_первая_грамматика] |
| человек | TAuxDicArticle [наша_первая_грамматика] |

Задание

- 1) Выделить содержание витаминов в продуктах из текста:

<http://chem21.info/info/1069461/>

«Витамин А содержится в моркови»

«Жиры рыб богаты витамином О»

База – tutorial1

- 2) Переложить вывод в таблицу фактов.

База – tutorial4