

МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ  
УНИВЕРСИТЕТ ИМ Н.Э. БАУМАНА

Маркин Кирилл Вадимович

**Разработка метода тематического моделирования  
для новостей на русском языке**

Специальность 2301050065 —  
«Программное обеспечение вычислительной техники и  
автоматизированных систем»

Квалификационная работа бакалавра  
кандидата в бакалавры

Научный руководитель:  
доцент, кандидат технических наук  
Клышинский Эдуард Станиславович

Консультант:  
старший преподаватель  
Волкова Лилия Леонидовна

Москва — 2019

## 1 Техническое задание

Заменить эту страницу на подписанное ТЗ

## 2 Календарный план

Заменить эту страницу на подписанный календарный план

## Оглавление

<b>1</b>	<b>Техническое задание . . . . .</b>	<b>2</b>
<b>2</b>	<b>Календарный план . . . . .</b>	<b>3</b>
<b>3</b>	<b>Реферат . . . . .</b>	<b>6</b>
<b>4</b>	<b>Перечень условных обозначений . . . . .</b>	<b>7</b>
<b>5</b>	<b>Введение . . . . .</b>	<b>8</b>
5.1	// актуальность выбранной темы . . . . .	8
5.2	// подвести к предметной области и задаче . . . . .	8
<b>6</b>	<b>Аналитический раздел . . . . .</b>	<b>9</b>
6.1	// цель, перечисляются задачи, которые необходимо решить для достижения этой цели . . . . .	9
6.2	// проводится анализ предметной области и выделяется основной объект исследования . . . . .	9
6.3	// обзор существующих путей/методов/решений и алгоритмов решения . . . . .	9
6.4	// обосновывается необходимость разработки нового или адаптации существующего метода или алгоритма . . . . .	9
6.5	// формализованное описание проблемы предметной области . . . . .	9
6.5.1	описание входных и выходных данных . . . . .	9
6.5.2	описание критериев сравнения нескольких реализаций метода или алгоритма . . . . .	9
6.5.3	описание функциональных требований к разрабатываемому программному обеспечению . . . . .	9
<b>7</b>	<b>Конструкторский раздел . . . . .</b>	<b>10</b>
7.1	// обосновать последовательность этапов выполнения . . . . .	10
7.2	// Алгоритм сбора данных . . . . .	10
7.3	// Алгоритм анализа . . . . .	10

7.4	// ? Что делаем . . . . .	11
7.5	// Оценка . . . . .	11
7.6	// Требования к программе . . . . .	11
<b>8</b>	<b>Технологический раздел . . . . .</b>	<b>12</b>
8.1	// обоснованный выбор средств программной реализации .	12
8.2	// описание основных (нетривиальных) моментов разработки . . . . .	12
8.3	// методики тестирования созданного программного обеспечения . . . . .	12
8.4	// информация, необходимая для сборки и запуска разработанного программного обеспечения . . . . .	12
<b>9</b>	<b>Экспериментальный раздел . . . . .</b>	<b>13</b>
9.1	// эксперименты и их результаты . . . . .	13
9.1.1	// проводим апробацию . . . . .	13
9.1.2	// анализируем результаты . . . . .	13
9.2	// качественное и количественное сравнение с аналогами .	13
9.3	// даём рекомендации о применимости метода/софта . . .	13
<b>10</b>	<b>Заключение . . . . .</b>	<b>14</b>
10.1	// отчитаться по каждому пункту тз/по каждой задаче и цели . . . . .	14
10.2	// сказать про перспективы (мы все уже не умрём) . . . .	14
<b>11</b>	<b>Список источников . . . . .</b>	<b>15</b>
11.1	// Разобрать . . . . .	15
11.2	// Датасеты . . . . .	15
<b>12</b>	<b>Приложения . . . . .</b>	<b>16</b>
12.1	// . . . . .	16

### 3 Реферат

Объект исследования и разработки

Цель и задачи работы

Метод и методология проведения работы

Результаты работы

Основные конструктивные, технологические и технико-эксплуатационные характеристики объекта исследования

Степень внедрения

Рекомендации по внедрению

Область применения

Экономическая эффективность или значимость работы

Прогнозы и предположения о возможных направлениях развития объекта исследования

#### 4 Перечень условных обозначений

Добавить условные обозначения (только если встречается более 3 раз)

## 5 Введение

2 - 3 страницы

Костя пошарил свою работу - глянуть что тут должно быть

5.1 // актуальность выбранной темы

5.2 // подвести к предметной области и задаче



## 6 Аналитический раздел

25 – 30 страниц

6.1 // цель, перечисляются задачи, которые необходимо решить для достижения этой цели

6.2 // проводится анализ предметной области и выделяется основной объект исследования

6.3 // обзор существующих путей/методов/решений и алгоритмов решения

Классификация

Кластеризация

PLSA

ARTM

dwl.kiev.ua - Дмитрия Владимировича Ландэ

6.4 // обосновывается необходимость разработки нового или адаптации существующего метода или алгоритма

выводы из обзора (лучше сравнительную таблицу) отсюда актуальность (никто не делал так/улучшаем то-то и то-то)

6.5 // формализованное описание проблемы предметной области

Необходимая существующая математика

6.5.1 описание входных и выходных данных

Откуда брать данные и какие они бывают

6.5.2 описание критериев сравнения нескольких реализаций метода или алгоритма

6.5.3 описание функциональных требований к разрабатываемому программному обеспечению

Что мы хотим получить (это и будет "мостиком" к конструкторской)

## 7 Конструкторский раздел

25 – 30 страниц

**7.1 // обосновать последовательность этапов выполнения**

**7.2 // Алгоритм сбора данных**

как будем извлекать данные (без кода пока)

Мой написанный код для парсинга

Уже предварительно собранные открытые данные

<https://newspaper.readthedocs.io/en/latest/> - возможный инструмент для парсинга

25 500 новостей (там суммарно 9 000 000 слов - я посчитал) за все время существования media.zone (я сам написал парсер, могу его же натравить на любой другой новостной ресурс) - уже скачены и лежат на моем компьютере

statmt.org - это не совсем подходит нам, тут новости короткие совсем. Но тоже скачал на всякий случай поиграться - тут суммарно 8,4 гигабайта чистого текста - уже скачены и лежат на моем компьютере

webhose.io - 290 000 новостей - уже скачены и лежат на моем компьютере

Можно сделать сервис на РИА новости

Можно сделать сервис на агрегаторы новостей

**7.3 // Алгоритм анализа**

разработка метода

Базовый алгоритм: ARTM (bigartm.readthedocs.io)

Предобработка текста: лемматизация, удаление стоп-слов, ngrams

Используем модальности (дата публикации, ссылки на другие документы, авторы)

Используем производные от статьи данные по различным алгоритмам (записываем в модальности) - алгоритмы еще не выбраны

IDEF0 метода

#### 7.4 // ? Что делаем

Можно попробовать обучаться на месяце/неделе/дне (и это в теории можно вынести в эксперимент) и выдавать как меняются темы

решить иерархически ли хотим строить темы или многое ко многим

#### 7.5 // Оценка

как будем оценивать (без кода)

Разбиение на 2 части и замеры разницы оценки - устойчивость - Через предложение разбивать статью можно попробовать

Толока - описание теста - выбрать лишнее слово, подумать что еще можно

#### 7.6 // Требования к программе

## 8 Технологический раздел

20 - 25 страниц

- 8.1 // обоснованный выбор средств программной реализации
- 8.2 // описание основных (нетривиальных) моментов  
разработки
- 8.3 // методики тестирования созданного программного  
обеспечения
- 8.4 // информация, необходимая для сборки и запуска  
разработанного программного обеспечения

## 9 Экспериментальный раздел

10 - 15 страниц

### 9.1 // эксперименты и их результаты

Можно поиграть с периодом обучения и сравнения данных (месяц/неделя/день) и посмотреть где лучше (?что лучше)

Можно поиграть с размером новости и посмотреть как от этого зависят результаты

#### 9.1.1 // проводим апробацию

#### 9.1.2 // анализируем результаты

9.2 // качественное и количественное сравнение с аналогами  
оцениваем адекватность и качество

9.3 // даём рекомендации о применимости метода/софта

## 10 Заключение

10.1 // отчитаться по каждому пункту тз/по каждой задаче  
и цели

10.2 // сказать про перспективы (мы все уже не умрём)

## 11 Список источников

### 11.1 // Разобрать

Ссылка на записи с datafest

Воронцов - книги и лекции

Ученики Воронцова - доклады и статьи

Анастасия Янина - работала с Воронцовым - посмотреть ее доклады и статьи

Потапенко Анна - работала с Воронцовым - посмотреть ее доклады и статьи

"Диалог NLP Конференция

курсы на курсере

dwl.kiev.ua - Дмитрия Владимировича Ландэ

### 11.2 // Датасеты

25 500 новостей (там суммарно 9 000 000 слов - я посчитал) за все время существования media.zone (я сам написал парсер, могу его же натравить на любой другой новостной ресурс) - уже скачены и лежат на моем компьютере

statmt.org - это не совсем подходит нам, тут новости короткие совсем. Но тоже скачал на всякий случай поиграться - тут суммарно 8,4 гигабайта чистого текста - уже скачены и лежат на моем компьютере

webhose.io - 290 000 новостей - уже скачены и лежат на моем компьютере

Можно сделать сервис на РИА новости

Можно сделать сервис на агрегаторы новостей

## 12 Приложения

добавить схемы, листинги программного кода, наборы тестов и др

12.1 //