

**Министерство образования и науки Российской Федерации**  
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ

**«САНКТ-ПЕТЕРБУРГСКИЙ НАЦИОНАЛЬНЫЙ  
ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ ИНФОРМАЦИОННЫХ  
ТЕХНОЛОГИЙ, МЕХАНИКИ И ОПТИКИ»**

**ПОЯСНИТЕЛЬНАЯ ЗАПИСКА  
ВЫПУСКНОЙ КВАЛИФИКАЦИОННОЙ РАБОТЫ**

**«Реализация эффективного взаимодействия между платформой  
для анализа экспрессии генов Morpheus и библиотекой  
вычислительных методов R/Bioconductor»**

Автор: Зенкова Дарья Михайловна \_\_\_\_\_

Направление подготовки (специальность): 01.03.02 Прикладная математика и  
информатика

Квалификация: Бакалавр

Руководитель: Сергушичев А.А., канд. техн. наук \_\_\_\_\_

**К защите допустить**

Зав. кафедрой Васильев В.Н., докт. техн. наук, проф. \_\_\_\_\_

« \_\_\_\_ » \_\_\_\_\_ 20 \_\_\_\_ г.

Санкт-Петербург, 2017 г.

**Студент** Зенкова Д.М. **Группа** М3436 **Кафедра** компьютерных технологий **Факультет** информационных технологий и программирования

**Направленность (профиль), специализация** Математические модели и алгоритмы разработки программного обеспечения

Квалификационная работа выполнена с оценкой \_\_\_\_\_

Дата защиты « \_\_\_\_ » \_\_\_\_\_ 20 \_\_\_\_ г.

Секретарь ГЭК \_\_\_\_\_

Листов хранения \_\_\_\_\_

Демонстрационных материалов/Чертежей хранения \_\_\_\_\_

**Министерство образования и науки Российской Федерации**  
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ  
**«САНКТ-ПЕТЕРБУРГСКИЙ НАЦИОНАЛЬНЫЙ  
ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ ИНФОРМАЦИОННЫХ  
ТЕХНОЛОГИЙ, МЕХАНИКИ И ОПТИКИ»**

**УТВЕРЖДАЮ**

Зав. каф. компьютерных технологий  
докт. техн. наук, проф.  
\_\_\_\_\_ Васильев В.Н.  
« \_\_\_\_ » \_\_\_\_\_ 20 \_\_\_\_ г.

**ЗАДАНИЕ  
НА ВЫПУСКНУЮ КВАЛИФИКАЦИОННУЮ РАБОТУ**

**Студент** Зенкова Д.М. **Группа** М3436 **Кафедра** компьютерных технологий **Факультет** информационных технологий и программирования  
**Руководитель** Сергушичев Алексей Александрович, канд. техн. наук, программист кафедры информационных систем

**1 Наименование темы:** Реализация эффективного взаимодействия между платформой для анализа экспрессии генов Morpheus и библиотекой вычислительных методов R/Bioconductor

**Направление подготовки (специальность):** 01.03.02 Прикладная математика и информатика

**Направленность (профиль):** Математические модели и алгоритмы разработки программного обеспечения

**Квалификация:** Бакалавр

**2 Срок сдачи студентом законченной работы:** «31» мая 2017 г.

**3 Техническое задание и исходные данные к работе.**

Разработать веб-приложение для анализа экспрессии генов, интегрирующее возможности визуального анализа morpheus.js и методы анализа библиотек R/Bioconductor. Веб-приложение должно быть легко дополняемо новыми методами для исследования и анализа экспрессии генов.

**4 Содержание выпускной квалификационной работы (перечень подлежащих разработке вопросов)**

- а) Обзор предметной области
- б) Архитектура проекта
- в) Практическая реализация и результаты

**5 Перечень графического материала (с указанием обязательного материала)**

Не предусмотрено

## 6 Исходные материалы и пособия

- а) Joshua Gould. Morpheus.js. JavaScript matrix visualization and analysis. [Электронный ресурс]. URL: <https://github.com/cmap/morpheus.js/>;
- б) Arora Sonali, Carlson Marc, Hayden Nate [и др.]. Bioconductor is an open source, open development software project to provide tools for the analysis and comprehension of high-throughput genomic data. [Электронный ресурс]. URL: <https://www.bioconductor.org/>;
- в) Ooms Jeroen. OpenCPU is a system for embedded scientific computing and reproducible research. [Электронный ресурс]. URL: <https://www.opencpu.org/>;
- г) Docker. Docker is the software container platform. [Электронный ресурс]. URL: <https://www.docker.com/>;

## 7 Календарный план

№№ пп.	Наименование этапов выпускной квалификационной работы	Срок выполнения этапов работы	Отметка о выполнении, подпись руков.
1	Ознакомление с предметной областью	30.09.2016	
2	Изучение исходного кода morpheus.js	31.10.2016	
3	Проектирование метода взаимодействия	30.11.2016	
4	Внедрение и тестирование нового функционала	31.03.2017	
5	Запуск веб-приложения в публичное пользование	28.04.2017	
6	Обработка результатов, написание пояснительной записки	31.05.2017	

**8 Дата выдачи задания:** «01» сентября 2016 г.

Руководитель \_\_\_\_\_

Задание принял к исполнению \_\_\_\_\_ «01» сентября 2016 г.

**«САНКТ-ПЕТЕРБУРГСКИЙ НАЦИОНАЛЬНЫЙ  
ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ ИНФОРМАЦИОННЫХ  
ТЕХНОЛОГИЙ, МЕХАНИКИ И ОПТИКИ»**

**АННОТАЦИЯ  
ВЫПУСКНОЙ КВАЛИФИКАЦИОННОЙ РАБОТЫ**

**Студент:** Зенкова Дарья Михайловна

**Наименование темы работы:** Реализация эффективного взаимодействия между платформой для анализа экспрессии генов Morpheus и библиотекой вычислительных методов R/Bioconductor

**Наименование организации, где выполнена работа:** Университет ИТМО

**ХАРАКТЕРИСТИКА ВЫПУСКНОЙ КВАЛИФИКАЦИОННОЙ РАБОТЫ**

**1 Цель исследования:** Создать веб-приложение, интегрирующее существующие возможности веб-приложения morpheus.js и методы анализа, реализованные в Bioconductor.

**2 Задачи, решаемые в работе:**

- а) разработка способа взаимодействия между js-клиентом и R и встраивание его в morpheus.js;
- б) создание графического интерфейса в js-клиенте и серверной реализации в R-пакете;
- в) объединение всех составляющих в единое веб-приложение phantassus;
- г) запуск веб-приложения в открытый доступ для исследователей.

**3 Число источников, использованных при составлении обзора:** \_\_\_\_\_

**4 Полное число источников, использованных в работе:** 0

**5 В том числе источников по годам**

Отечественных			Иностранных		
Последние 5 лет	От 5 до 10 лет	Более 10 лет	Последние 5 лет	От 5 до 10 лет	Более 10 лет

**6 Использование информационных ресурсов Internet:** \_\_\_\_\_

**7 Использование современных пакетов компьютерных программ и технологий:**

Были использованы следующие программы и технологии: язык программирования JavaScript, фреймворк Node.js, веб-приложение morpheus.js, язык программирования R, библиотека биоинформатических алгоритмов Bioconductor, система интеграции R OpenCPU, механизм для сериализации данных Protocol Buffers, репозитория геномных данных Gene Expression Omnibus, программное обеспечение для запуска приложений в контейнерах Docker, веб-сервер Apache, среда разработки WebStorm, среда разработки RStudio, система контроля версий git, система компьютерной верстки L<sup>A</sup>T<sub>E</sub>X.

**8 Краткая характеристика полученных результатов:** Реализовано веб-приложение phantusus, отвечающее всем поставленным требованиям. Веб-приложение было запущено в публичный доступ, используется в лаборатории Максима Артемова в Washington University in St. Louis. Демонстрация приложения входит в программу семинара по системной биологии в Сиднее (10-13 апреля 2017) и в Санкт-Петербурге (14-19 мая 2017).

**9 Гранты, полученные при выполнении работы:** Работа над данной инженерной разработкой велась без поддержки грантами.

**10 Наличие публикаций и выступлений на конференциях по теме работы:** По данной инженерной разработке не имеется публикаций и она не была представлена на конференциях.

Выпускник: Зенкова Д.М. \_\_\_\_\_

Руководитель: Сергушичев А.А. \_\_\_\_\_

« \_\_\_\_ » \_\_\_\_\_ 20 \_\_\_\_ г.

## ОГЛАВЛЕНИЕ

ВВЕДЕНИЕ.....	5
1. Обзор предметной области .....	6
1.1. Биоинформатика .....	6
1.1.1. Анализ экспрессии генов .....	6
1.1.2. Используемые методы .....	6
1.2. Существующие решения для анализа экспрессии генов .....	7
1.2.1. R/Bioconductor .....	7
1.2.2. GENE-E .....	7
1.2.3. morpheus.js .....	8
1.3. Инструменты, которые могут быть применены .....	8
1.3.1. Язык R и библиотека Bioconductor .....	8
1.3.2. JavaScript и Node.js .....	8
1.3.3. R shiny.....	8
1.3.4. OpenCPU .....	8
1.3.5. Gene Expression Omnibus .....	9
1.3.6. Docker .....	9
1.3.7. JSON .....	10
1.3.8. Protocol Buffers .....	10
1.3.9. Apache.....	10
1.3.10.HTML.....	10
1.4. Постановка задачи.....	10
1.4.1. Цель работы .....	10
1.4.2. Основные задачи.....	10
1.4.3. Требования к веб-приложению phantasus.....	10
Выводы по главе 1.....	11
2. Архитектура проекта phantasus.....	12
2.1. morpheus.js .....	12
2.1.1. Чтение данных .....	12
2.1.2. Класс Dataset.....	12
2.1.3. Класс SlicedDatasetView .....	13
2.1.4. Класс HeatMap .....	13
2.1.5. Реализованные методы .....	13

2.2. R-пакет phantusus.....	14
2.3. Связь через OpenCPU API .....	14
2.4. Неотсортировано .....	14
Выводы по главе 2.....	14
3. Реализация .....	15
Выводы по главе 3.....	15
ЗАКЛЮЧЕНИЕ.....	16



**ВВЕДЕНИЕ**

В данном разделе размещается введение.

## ГЛАВА 1. ОБЗОР ПРЕДМЕТНОЙ ОБЛАСТИ

### 1.1. Биоинформатика

**Биоинформатика** — наука, объединяющая в себе методы прикладной математики, статистики, информатики для создания новых методов и алгоритмов для анализа разного рода биологических данных.

Биоинформатика занимается биохимией, биофизикой, экологией и многими другими областями биологии. Однако в данной работе фокус направлен на геномную биоинформатику и на конкретную ее задачу — анализ экспрессии генов.

#### 1.1.1. Анализ экспрессии генов

**Экспрессия генов** — процесс преобразования наследственной информации от гена (в виде последовательности нуклеотидов ДНК) в функциональный продукт (РНК или белок).

Анализ экспрессии генов позволяет выяснить как ведет себя отдельный ген в разных условиях, тканях или организмах. Так, например, можно исследовать экспрессию вирусных белков или экспрессию онкогенов.

#### 1.1.2. Используемые методы

Как было сказано ранее, биоинформатика использует в себе математику, информатику и статистику. Соответственно, задача анализа экспрессии генов сводится к исследованию путем статистических методов и алгоритмов числовой двумерной матрицы, как, например, небольшой срез матрицы в таблице 1.

Таблица 1 – Срез матрицы GSE14308. Строки матрицы соответствуют генам, столбцы — образцам.

	GSM357839	GSM357841	GSM357842	GSM357843	GSM357844
Rps29	16.32	16.30	16.25	16.32	16.30
Rpl13a	16.27	16.23	16.32	16.30	16.27
Rps3a1	16.23	16.19	16.30	16.25	16.25
Rpl38	16.21	16.25	16.27	16.27	16.21
Tmsb4x	16.30	16.32	16.23	16.21	16.32

На рисунке можно увидеть визуализацию матрицы экспрессии в виде тепловой карты.

Также к основным методам анализа относятся:

- Иерархическая и вероятностная кластеризации;
- Дифференциальная экспрессия;
- Метод главных компонент и визуализация его результатов.

## 1.2. Существующие решения для анализа экспрессии генов

### 1.2.1. R/Bioconductor

**R** - язык программирования для статистического анализа данных и работы с графикой [**rproject**].

**Bioconductor** - библиотека, содержащая в себе множество реализаций биоинформатических алгоритмов и методов обработки биологических данных на R. Она постоянно обновляется, пополняется новыми библиотеками, модерируется сообществом [**bioconductor**]. R и Bioconductor очень популярны в биоинформатической среде ввиду предоставляемых возможностей.

Однако для качественного и полноценного анализа с помощью этих инструментов, нужно иметь навыки программирования на R, что весьма неудобно для исследователей биологических специальностей.

### 1.2.2. GENE-E

**GENE-E** - Платформа для анализа данных и визуального исследования данных, созданная на Java и R [**genee**]. Содержит в себе множество полезных для исследования инструментов: тепловые карты, кластеризацию, фильтрацию, построение графиков и т.д. Позволяет исследовать любые данные в виде матрицы. К тому же, содержит дополнительные инструменты для геномных данных.

Недостатки:

- Чтобы использовать, необходимо устанавливать на свой компьютер;
- Поддержка данного приложения прекратилась в связи с созданием morpheus.js [**morpheus**];
- Не имеет открытого исходного кода, а только API для взаимодействия и создания новых приложений на его основе.

### 1.2.3. morpheus.js

**Morpheus.js** - веб-приложение для визуализации и анализа матриц от создателя GENE-E [**morpheus**]. Создано уже на JavaScript и с открытым исходным кодом. Удобно для использования исследователями без навыков программирования и так же, как и GENE-E, применимо к любым матрицам.

Недостатки:

- Ограниченный набор функций, которых недостаточно для полноценного анализа;
- Для расширения биоинформатическими алгоритмами, не прибегая к дополнительным инструментам, требуется реализовывать их заново на JavaScript.

## 1.3. Инструменты, которые могут быть применены

### 1.3.1. Язык R и библиотека Bioconductor

Как было сказано ранее, Bioconductor полон актуальными и широко используемыми биоинформатическими алгоритмами, в том числе и для анализа экспрессии генов. Соответственно, реализовывать их заново обычно нет необходимости и можно использовать их для достижения целей этой работы.

### 1.3.2. JavaScript и Node.js

**JavaScript** — мультипарадигменный скриптовый язык программирования, широко используемый для создания веб-приложений.

**Node.js** — [**nodejs**]

### 1.3.3. R shiny

### 1.3.4. OpenCPU

#### 1.3.4.1. Необходимые определения

HTTP HTTP API RPC Веб-сервер

**OpenCPU** — система для встроенных научных вычислений и воспроизводимых исследований, предоставляющая HTTP API для взаимодействия с R-серверами [**opencpu**].

Имеется также библиотека *opencpu.js* для интеграции JavaScript и R.

### 1.3.5. Gene Expression Omnibus

**Gene Expression Omnibus (GEO)** — международный публичный репозиторий, агрегирующий и распространяющий различные формы геномных данных от исследовательского сообщества [geo].

В библиотеке Bioconductor есть R-пакет *GEOquery* для удобной загрузки данных из GEO [geoquery].

### 1.3.6. Docker

**Docker** — программное обеспечение для автоматизации запуска и внедрения приложений внутри контейнеров [docker].

Для дальнейшего описания данного инструмента введем несколько определений.

*Образ* — отдельный исполняемый пакет, включающий себя все необходимое для запуска единицы программного обеспечения, в том числе исходный код, библиотеки, переменные окружения, конфигурационные файлы. Зачастую образ построен на основе другого образа с дополнительной конфигурацией. Образ компилируется по *Dockerfile*, каждая команда в котором соответствует новому слою. При перекомпиляции обновляются только те слои, которые изменились.

*Контейнер* — запущенный экземпляр образа. Контейнер обычно выполняется изолированно от окружения, имея доступ к файлам или портам хост-системы только при наличии соответствующей конфигурации.

В отличие от виртуальных машин, которые запускают гостевую операционную систему в каждом экземпляре, контейнеры могут разделять общее ядро, и вся информация, которая должна быть в контейнере, это исполняемый процесс и его зависимости. Исполняемые процессы из контейнеров работают как нативные процессы, и могут управляться по отдельности.

Для контроля версий и хранения образов в открытом доступе используется Docker Hub [dhub]. В этом хранилище можно как добавлять репозитории, управляемые вручную, так и поддерживать автоматические сборки (*Automated Build*), которые привязаны к репозиториям на в популярных системах контроля версий: GitHub [github] и Bitbucket [bitbucket].

### 1.3.7. JSON

### 1.3.8. Protocol Buffers

**Protocol Buffers (Protobuf)** — гибкий, универсальный и автоматизированный механизм для сериализации структурированных данных [protobuf].

Структура информации задается с помощью \*.proto файлов в форме сообщений (Message).

### 1.3.9. Apache

### 1.3.10. HTML

## 1.4. Постановка задачи

Рассмотрев существующие решения для анализа экспрессии генов и инструментов, которые могли бы пригодиться для будущих решений, можно сформулировать цель и основные задачи данной работы

### 1.4.1. Цель работы

Создать веб-приложение, интегрирующее существующие возможности веб-приложения morpheus.js и методы анализа, реализованные в Bioconductor.

### 1.4.2. Основные задачи

- а) Разработать способ взаимодействия между js-клиентом и R и встроить его в morpheus.js, чтобы избежать реализации с нуля уже существующих алгоритмов;
- б) Реализовать графический интерфейс в js-клиенте и серверную реализацию в R-пакете;
- в) Соединить все составляющие в одном веб-приложении phantasia;
- г) Запустить веб-приложение в открытый доступ для исследователей.

### 1.4.3. Требования к веб-приложению phantasia

#### 1.4.3.1. Доступность

Необходимо, чтобы веб-приложение phantasia было доступно для исследователей независимо от их местоположения и времени суток. Варианты действий:

- а) Сделать его доступным по определенному веб-адресу, и тогда пользователь сможет продолжать исследования из любой точки, где есть подключение к интернету;

- б) Предоставить возможность запускать приложение локально, например, с помощью Docker или внутри R.

#### **1.4.3.2. Возможность дальнейшего расширения функционала**

Как уже было сказано выше, библиотека Bioconductor постоянно обновляется и пополняется новыми алгоритмами, а исследователи находят новые методы для анализа экспрессии генов, так что необходимо не только реализовать дополнительные методы, но и отладить и описать алгоритм действий для добавления новых.

### **Выводы по главе 1**

## ГЛАВА 2. АРХИТЕКТУРА ПРОЕКТА PHANTASUS

### 2.1. morpheus.js

Как уже было рассказано в обзоре, morpheus.js — веб-приложение, полностью созданное на JavaScript, для визуализации и анализа матриц.

В этом разделе будут описаны основные классы и функции, реализованные в исходном коде morpheus.js, которые будут в дальнейшем необходимы для расширения функционала.

#### 2.1.1. Чтение данных

В morpheus.js данные могут быть загружены из файла, полученного одним из следующих путей:

- Из компьютера;
- По URL-ссылке;
- Из Dropbox.

Допустимые форматы загружаемых файлов:

- txt-файл с tab-разделителями;
- Excel-таблица;
- MAF [**maf**];
- GCT [**gct**];
- GMT [**gmt**].

Для каждого формата файла в исходном коде morpheus.js присутствует соответствующий обработчик данных.

Также, morpheus.js предлагает набор предзагруженных данных из базы TCGA [**tcga**].

#### 2.1.2. Класс Dataset

Одним из ключевых классов всего веб-приложения является класс Dataset. В каждом экземпляре этого класса хранится вся необходимая информация о данных, в которую входят:

- Числовая матрица, характеризующая уровень экспрессии всех генов во всех образцах;
- Количество строк и столбцов в матрице;
- Аннотация к образцам, например:
  - пол, возраст, контактную информацию испытуемых, если образцы были взяты с людей;



- есть или нет инфекция в данном образце;
- способ лечения;
- контакты ответственного за взятие данного образца и пр.;
- Аннотация к генам, например:
  - Идентификатор гена в том или ином стандарте;
  - Числовые характеристики гена (средний уровень экспрессии по образцам, номер кластера) и пр.

Аннотация реализована в классе `MetadataModel`, который представляет собой не что иное, как набор именованных векторов с характеристиками. В каждом векторе хранятся:

- Название;
- Формат (строка, число);
- Массив значений.

Для векторов так же предусмотрены утилиты для визуализации. Так, например, есть возможность показать аннотацию в виде текста и/или цветом, что удобно для категориальных характеристик.

### 2.1.3. Класс `SlicedDatasetView`

Чаще всего во время работы программы экземпляры класса `Dataset` становятся обернуты в оболочку из `SlicedDatasetView`. Этот дополнительный класс дает возможность не пересоздавать каждый раз `Dataset`, а просто добавляет к данным информацию о том, какие индексы строк и столбцов выбраны и используются в данный момент.

### 2.1.4. Класс `HeatMap`

Данный класс предназначен для визуализации данных, обернутых в класс `Dataset` или `SlicedDatasetView`. Он дает возможность выбирать, какая аннотация будет представлена на экране, цветовой код, выбирать строки и столбцы, с которыми будут работать те или иные инструменты.

### 2.1.5. Реализованные методы

В `morpheus.js` имеются реализации следующих методов:

- `Adjust` — инструмент для корректировки данных:
  - $\log_2$ ;
  - $\log_2^{-1}$ ;

- Квантиль-нормализация;
- Z-тест;
- Устойчивый Z-тест;
- Collapse — инструмент, позволяющий агрегировать строки или столбцы с одинаковыми значениями с помощью функции: min, max, mean, median, sum, maximum of 25 and 75 percentiles;
- Создать вычисленную аннотацию для строк или столбцов;
- Similarity Matrix;
- Transpose;
- t-SNE;
- Построение графиков.

Также присутствуют фильтрация и сортировка.

## **2.2. R-пакет phantasus**

## **2.3. Связь через OpenCPU API**

## **2.4. Неотсортировано**

## **Выводы по главе 2**

## **ГЛАВА 3. РЕАЛИЗАЦИЯ**

### **Выводы по главе 3**

**ЗАКЛЮЧЕНИЕ**

В данном разделе размещается заключение.