



Министерство науки и высшего образования Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Московский государственный технический университет
имени Н. Э. Баумана
(национальный исследовательский университет)»
(МГТУ им. Н. Э. Баумана)

ФАКУЛЬТЕТ «Информатика и системы управления»

КАФЕДРА «Программное обеспечение ЭВМ и информационные технологии»

РАСЧЕТНО-ПОЯСНИТЕЛЬНАЯ ЗАПИСКА

К КУРСОВОЙ РАБОТЕ

НА ТЕМУ:

*«Разработка базы данных для хранения и обработки
данных для приложения учета и аудита времени,
потраченного на рабочие и личные задачи»*

Студент ИУ7-66Б
(Группа)

(Подпись, дата)

М. Ю. Вольняга
(И. О. Фамилия)

Руководитель курсовой работы

(Подпись, дата)

К. Л. Тассов
(И. О. Фамилия)

2024 г.

СОДЕРЖАНИЕ

ВВЕДЕНИЕ	4
1 Аналитический раздел	5
1.1 Анализ предметной области	5
1.2 Формулировка требований к разрабатываемой базе данных и приложению	6
1.3 Формализация и описание информации	7
1.4 Формализация и описание пользователей проектируемого при- ложения к базе данных	9
1.5 Диаграмма вариантов использования	10
1.6 Анализ существующих баз данных на основе формализации данных	11
2 Конструкторский раздел	14
2.1 Проектирование базы данных	14
2.2 Требования к программному обеспечению	20
2.3 Описание используемых типов данных	21
2.4 Разработка алгоритмов	21
3 Технологический раздел	23
3.1 Средства реализации	23
3.2 Сведения о модулях программы	24
3.3 Реализация алгоритмов	24
4 Исследовательская часть	28
4.1 Демонстрация работы программы	28
4.2 Технические характеристики	28
4.3 Время выполнения реализаций алгоритмов	28
ЗАКЛЮЧЕНИЕ	31
ПРИЛОЖЕНИЕ А	33
ПРИЛОЖЕНИЕ Б	35

ВВЕДЕНИЕ

В современном мире, где время является ценным ресурсом, эффективное управление им является важной задачей как для отдельных людей, так и для организаций. Чтобы грамотно распоряжаться своим временем, необходимо знать, сколько времени реально требуется на выполнение каждой задачи. Именно поэтому разработка приложения для учета и аудита времени, потраченного на определенную задачу, является актуальной [1].

Приложение, которое объединит в себе удобство канбан-досок для визуального планирования задач с интегрированным таймером для учета времени и аналитики, станет уникальным решением и поможет эффективно планировать как личные, так и рабочие задачи. Оно позволит пользователям совместно вести учет времени, что сделает его незаменимым инструментом для широкого круга пользователей, от фрилансеров до компаний, стремящихся к оптимизации своих процессов управления временем и повышению общей продуктивности.

Цель курсовой работы — разработка базы данных для хранения и обработки данных приложения учета и аудита времени, потраченного на рабочие и личные задачи.

Для достижения поставленной цели необходимо выполнить следующие задачи:

- сформулировать описание пользователей проектируемого приложения по учету и аудиту времени, потраченного на рабочие и личные задачи для доступа к базе данных;
- спроектировать сущности базы данных и ограничения целостности учета и аудита времени, потраченного на рабочие и личные задачи;
- выбрать средства реализации и разработать базу данных и приложение;
- провести исследование зависимости времени выполнения запросов от количества записей в базе данных.

1 Аналитический раздел

В данном разделе будет рассмотрено следующее: анализ предметной области, формулировка требований к разрабатываемой базе данных и приложению, формализация и описание информации, анализ существующих баз данных на основе формализации данных, ER-диаграмма, формализация и описание пользователей проектируемого приложения к базе данных, диаграмма вариантов использования.

1.1 Анализ предметной области

Предметная область «эффективное управление временем и задачами» затрагивает процессы планирования, отслеживания и анализа времени, затрачиваемого на выполнение различных задач и проектов.

Существует множество теорий и методологий, способствующих эффективному управлению временем. Среди наиболее известных подходов выделяются Матрица Эйзенхауэра, которая помогает приоритизировать задачи по их срочности и важности, метод Помодоро, предлагающий разбивку работы на короткие интервалы с перерывами, и система GTD, нацеленная на снижение стресса за счет переноса задач из головы в внешнюю систему организации. Современные методологии планирования, такие как Agile, Scrum и Kanban, приносят гибкость и адаптивность в процесс управления проектами и задачами. Они позволяют командам реагировать на изменения и эффективно управлять рабочим процессом, поддерживая постоянную обратную связь и цикличность в деятельности [2; 3].

Анализ аналогичных решений

На рынке существует большое количество приложений, которые имеют функциональность замера времени для задач, наиболее известные: clockify, toggl, timesamp, также есть приложения по планированию задач в стиле канбан-досок, например как trello. Выделим следующие критерии для сравнения:

- наличие таймера для задачи;
- ограниченность в создание проектов;
- наличие аналитики и экспорт ее;
- ограниченность в количестве пользователей у проекта;

- наличие канбан-досок для удобной организации задач.

В таблице 1.1 представлено сравнение существующих решений для бесплатных версий приложений.

Таблица 1.1 – Сравнение существующих решений

Критерий	clockify	toggl	timecamp	trello
Таймер	Есть	Есть	Есть	Нет
Количество проектов	Неограниченно	Неограниченно	Ограничено	Ограничено
Экспорт аналитики	PDF, CSV, Excel	PDF, CSV	PDF	Нет
Кол-во пользователей	Неограниченно	Ограничено	Ограничено	Ограничено
Канбан-доски	Нет	Нет	Нет	Есть

Как видно из таблицы 1.1 ни одно из рассмотренных решений не удовлетворяет всем критериям сравнения, также данные приложения являются зарубежными. Таким образом мое решение будет актуальным и в отличие от других решений, будет иметь канбан-доски с интегрированным таймером для каждой задачи, облегчая тем самым учет времени и повышая эффективность планирования

1.2 Формулировка требований к разрабатываемой базе данных и приложению

В рамках курсовой работы необходимо разработать базу данных, в которой будет храниться информация о пользователях, их проектах, задачах, карточках, тегах, списках дел и группах. Также были выделены необходимые требования для приложения. Приложение должно включать следующие функциональности:

- регистрация и аутентификация пользователя;
- изменение списка дел, задач проектов и тегов;
- поиск задач, проектов и тегов;
- замер времени для конкретной задачи;
- просмотр и сохранение аналитики использования времени, замеренного для определенных задач;

- создание и изменение групп пользователей.

Необходимо предусмотреть работу с задачами в виде канбан-досок, так как данная функция является отличительной от других приложений.

1.3 Формализация и описание информации

Разрабатываемая база данных для приложения по учету и аудиту затраченного времени на определенные задачи, должна содержать информацию о пользователях, их проектах, задачах, карточках, тегах, списках дел и группах.

В таблице 1.2 представлены необходимые таблицы и информация, которая должна содержаться в этой таблице для разрабатываемой базы данных.

Таблица 1.2 – Таблиц и их информации

Таблица	Информация
Пользователь	id, имя, фамилия, роль, почта, дата последнего входа и дата регистрации
Проект	id, название, описание
Задача	id, название, статус, описание, дата начала и завершения задачи, сумма времени таймера,
Группа	id, название, описание
Тег	id, название, цвет
Список дел	id, статус, приоритет, дата завершения, контент

На рисунке 1.1 приведена ER диаграмма в нотации Чена.

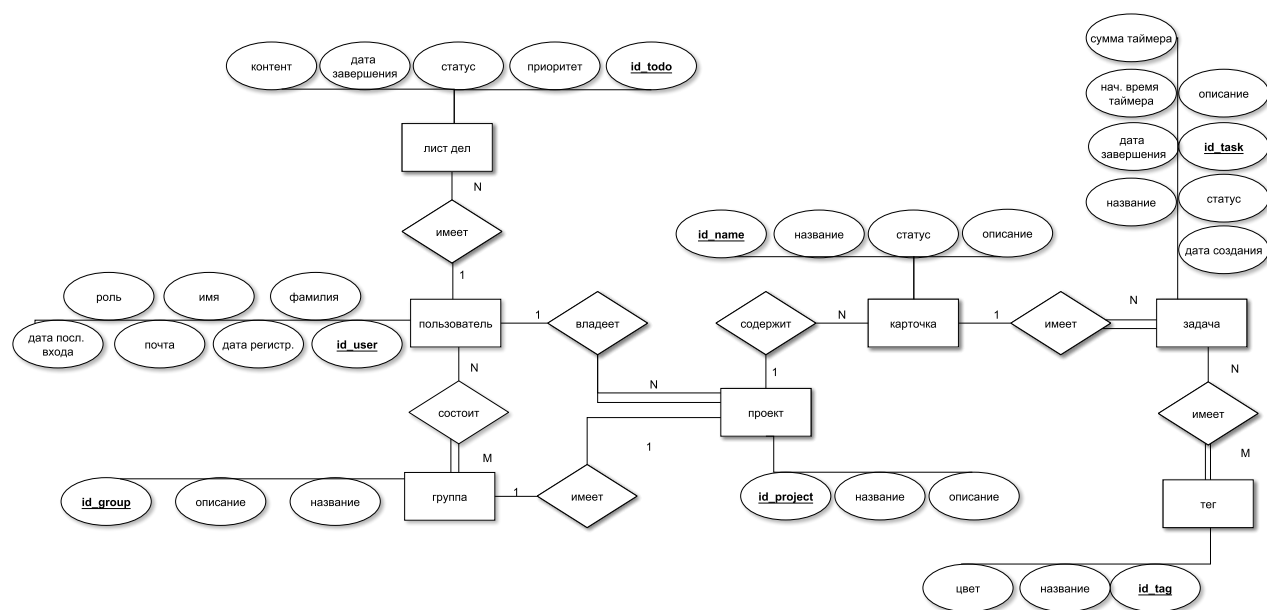


Рисунок 1.1 – ER диаграмма

1.4 Формализация и описание пользователей проектируемого приложения к базе данных

Для взаимодействия с приложением по учету и аудиту времени, было выделено три роли пользователей: неавторизированный, авторизованный и администратор.

В таблице 1.3 представлена функциональность для администратора, неавторизованного и авторизованного пользователя. Также в таблице 1.3 используются сокращения: 1 означает пользователь администратор, 2 — неавторизованный пользователь, 3 — авторизованный пользователь.

Таблица 1.3 – Функциональность администратора, неавторизованного и авторизованного пользователя

Функциональность	1	2	3
Зарегистрироваться и войти в систему	+	+	+
Изменить список дел		+	+
Создать и изменить личный проект, задачи и теги		+	+
Провести поиск задач, проектов и тегов		+	+
Замерить время для конкретной задачи		+	+
Просмотреть и сохранить аналитику		+	+
Создать группу			+
Изменить группу (если владелец)			+
Просмотреть и сохранить аналитику (группы)	+		+
Изменить группы пользователей	+		
Удалить пользователя из системы	+		

1.5 Диаграмма вариантов использования

На рисунке 1.2 приведена диаграмма вариантов использования.

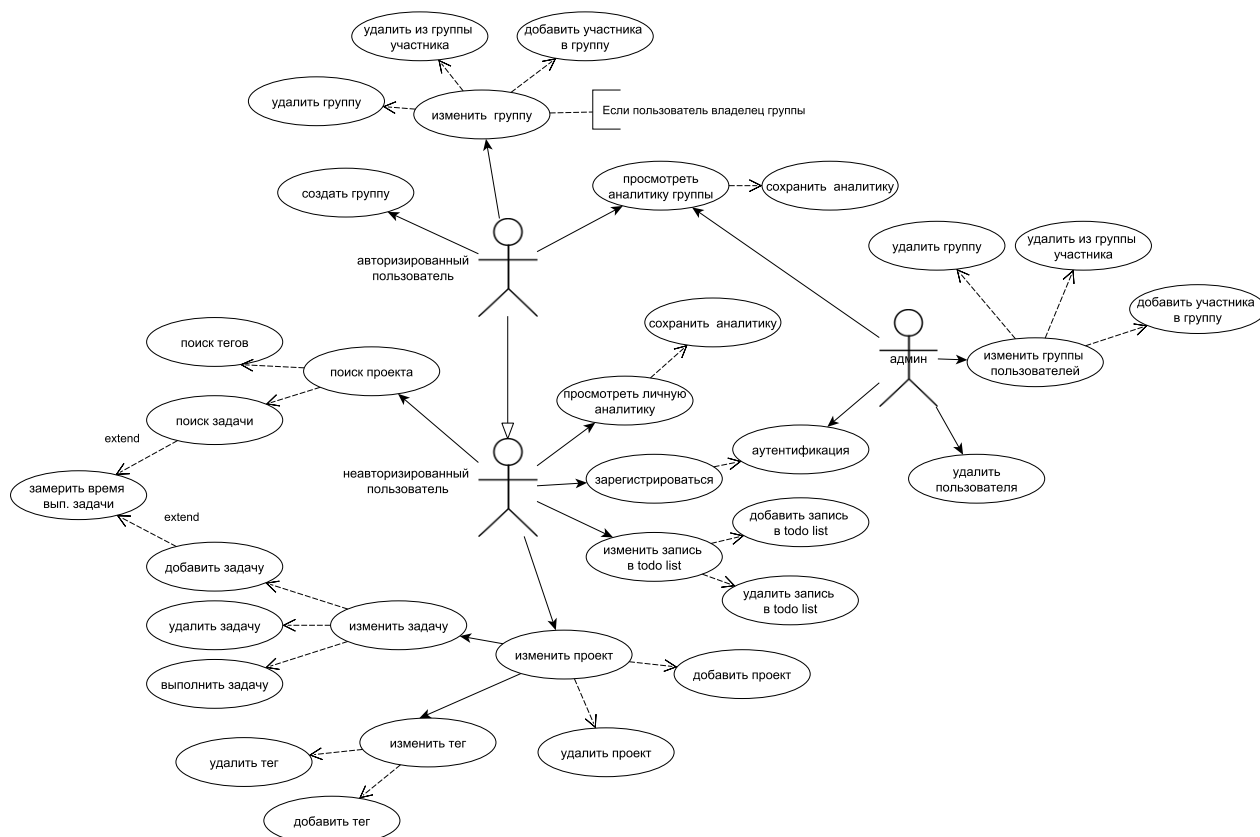


Рисунок 1.2 – Диаграмма использования

1.6 Анализ существующих баз данных на основе формализации данных

В контексте темы курсовой работы и проведенной формализации задачи, данных и пользователей, была выбрана реляционная модель данных.

Реляционная модель данных представляет собой подход к управлению данными, где все данные хранятся в таблицах (также называемых «отношениями»), состоящих из строк и столбцов. Каждая строка таблицы представляет собой запись с уникальным идентификатором (ключом), а каждый столбец — атрибутом данных. Реляционные базы данных обеспечивают гибкость в обработке данных, возможность выполнения сложных запросов и обеспечения целостности данных через систему отношений между таблицами.

SQLite

SQLite — это встраиваемая, высоконадежная библиотека СУБД, которая предоставляет реляционную базу данных в виде файла. Ее главное преимущество — не требует отдельного серверного процесса или системы управления, что делает SQLite идеальным выбором для мобильных приложений, приложений для настольных компьютеров и небольших веб-проектов, где требуется простота без жертвы функциональности. Она поддерживает основные операции SQL и отличается малым размером, быстродействием и надежностью.

PostgreSQL

PostgreSQL — это мощная, открытая и полностью бесплатная объектно-реляционная система управления базами данных (ОРСУБД). Она предоставляет расширенные возможности SQL, включая сложные запросы, внешние ключи, триггеры, представления, транзакции, индексирование полнотекстового поиска и сохраненные процедуры. PostgreSQL является хорошим выбором для крупных приложений и систем, требующих высокой надежности, масштабируемости и гибкости.

MySQL

MySQL — это популярная открытая реляционная система управления базами данных. Она широко используется в веб-разработке и является частью стека LAMP (Linux, Apache, MySQL, PHP/Python/Perl). MySQL предлагает высокую производительность, надежность и простоту в использовании, а также поддерживает широкий спектр операционных систем. Эта СУБД

идеально подходит для веб-сайтов и приложений среднего и большого размера, требующих эффективного управления данными.

Эти СУБД выбраны из-за их популярности, мощности, гибкости и разнообразия функций, которые они предлагают для разработчиков мобильных приложений, веб-приложений и крупных систем.

Таблица 1.4 – Сравнение реляционных СУБД

Критерий	SQLite	PostgreSQL	MySQL
Лицензия	Public Domain [sqlite]	PostgreSQL License, открытая [postgresql]	GPL для бесплатной версии или коммерческая лицензия [mysql]
Тип хранения	Встраиваемая	Клиент-сервер	Клиент-сервер
Поддержка JSON	Частичная	Полная	Полная
Репликация	Поддержка через сторонние инструменты	Нативная поддержка	Нативная поддержка
Полнотекстовый поиск	Да	Да	Да

Источники:

SQLite: Официальный сайт SQLite.

PostgreSQL: Официальный сайт PostgreSQL.

MySQL: Официальный сайт MySQL.

SQLite был выбран как целевая СУБД для разработки мобильного приложения на Kotlin по нескольким ключевым причинам, учитывая специфику задачи, данные и пользователей:

Встраиваемость и Легковесность: SQLite — это встраиваемая база данных, что означает, что она легко интегрируется непосредственно в мобильное приложение без необходимости во внешнем сервере баз данных. Это делает SQLite особенно подходящим для мобильных приложений, где необходимо минимизировать использование ресурсов и обеспечить высокую производительность.

Независимость и Портативность: SQLite хранит всю базу данных в

одном файле, что облегчает передачу, резервное копирование и синхронизацию данных между устройствами. Это важно для мобильных приложений, которые могут работать в автономном режиме или в условиях ненадежного сетевого соединения.

Простота Управления: SQLite не требует настройки, установки и управления сервером баз данных, что снижает затраты на поддержку и управление приложением. Это делает его идеальным выбором для малых и средних проектов, где может не хватать ресурсов для обслуживания сложной инфраструктуры баз данных.

Поддержка Транзакций и ACID: SQLite полностью поддерживает транзакции и соответствует принципам ACID (атомарность, согласованность, изоляция, долговечность), что гарантирует надежность и целостность данных в приложении.

Широкая Поддержка и Сообщество: SQLite имеет обширную документацию и поддерживается большим сообществом разработчиков. Это обеспечивает доступность ресурсов, библиотек и инструментов для разработки, тестирования и отладки приложений.

Гибкость: Несмотря на свою простоту, SQLite предлагает достаточный набор функций SQL, включая поддержку подзапросов, внешних ключей, триггеров и представлений, что позволяет эффективно управлять данными в мобильном приложении.

В совокупности, эти преимущества делают SQLite оптимальным выбором для мобильных приложений, где требуется надежное, эффективное и удобное в обслуживании решение для хранения данных.

Вывод

В данном разделе было рассмотрено следующее анализ предметной области, формулировка требований к разрабатываемой базе данных и приложению, формализация и описание информации, анализ существующих баз данных на основе формализации данных, ER-диаграмма, формализация и описание пользователей проектируемого приложения к базе данных, диаграмма вариантов использования.

2 Конструкторский раздел

В данном разделе будут представлены требования к программному обеспечению, описание используемых типов данных и схемы реализуемых алгоритмов.

2.1 Проектирование базы данных

Исходя из ER диаграммы, приведенной в предыдущем разделе (см. рисунок 1.1), были выделены следующие таблицы :

- таблица пользователей (user_app);
- таблица записей списка дел (todo_node);
- таблица проектов (project);
- таблица ролей (roles);
- таблица групп пользователей (group_user);
- таблица для связи группы с пользователем (group_member);
- таблица таблиц задач (table_app);
- таблица тегов (tag);
- таблица задач (task);
- таблица для связи задач с тегами (task_tag).

Описание сущностей и ограничений целостности в проектируемой базе данных

Информация о полях в таблице пользователей (user_app) представлена в таблице 2.1.

Таблица 2.1 – Описание таблицы пользователей

Поле	Тип данных	Ограничение	Описание
id_user	Целочисленный	Уникальное, обязательное	Первичный ключ таблицы
email	Символьный	Уникальное, обязательное, максимальная длина — 255 символов	Адрес электронной почты пользователя (логин пользователя)
password	Символьный	Обязательное, максимальная длина — 255 символов	Хэш пароля пользователя
user_name	Символьный	Обязательное, максимальная длина — 255	Имя пользователя
user_secondname	Символьный	Обязательное, максимальная длина — 255 символов	Фамилия пользователя
data_last_log_in	Дата	—	Дата последнего входа пользователя
data_sign_in	Дата	—	Дата регистрации пользователя
role_id	Целочисленный	Обязательное	Внешний ключ, для определения роли

Информация о полях в таблице ролей (roles) представлена в таблице 2.2.

Таблица 2.2 – Описание таблицы ролей

Поле	Тип данных	Ограничение	Описание
id_role	Целочисленный	Уникальное, обязательное	Первичный ключ таблицы
name	Символьный	Обязательное	Имя роли

Информация о полях в таблице записи списка дел (todo_node) представлена в таблице 2.3.

Таблица 2.3 – Описание таблицы записи списка дел

Поле	Тип данных	Ограничение	Описание
id_node	Целочисленный	Уникальное, обязательное	Первичный ключ таблицы
content	Символьный	Максимальная длина — 255 символов	Описание записи
priority	Символьный	Максимальная длина — 255 символов	Приоритет записи
status	Символьный	Максимальная длина — 255 символов, обязательное	Статус записи
due_data	Дата	—	Дата и время срока выполнения записи
user_id	Целочисленный	Обязательное	Идентификатор пользователя, создавшего запись, внешний ключ

Информация о полях в таблице проекта (project) представлена в таблице 2.4.

Таблица 2.4 – Описание таблицы проекта

Поле	Тип данных	Ограничение	Описание
id_project	Целочисленный	Уникальное, обязательное	Первичный ключ таблицы
name	Символьный	Обязательное, максимальная длина — 255 символов	Название проекта
description	Символьный	Обязательное, максимальная длина — 255 символов	Описание проекта
user_id	Целочисленный	Обязательное	Идентификатор пользователя, создавшего проект, внешний ключ

Информация о полях в таблице таблиц (table_app) представлена в таблице 2.5.

Таблица 2.5 – Описание таблицы таблиц

Поле	Тип данных	Ограничение	Описание
id_project	Целочисленный	Уникальное, обязательное	Первичный ключ таблицы
name	Символьный	Максимальная длина — 255 символов, обязательное	Название таблицы
description	Символьный	Максимальная длина — 255 символов	Описание таблицы
status	Символьный	Максимальная длина — 255 символов, обязательное	Статус таблицы
project_id	Целочисленный	Обязательное	Идентификатор проекта к которому относится таблица задач, внешний ключ

Информация о полях в таблице задачи (task) представлена в таблице 2.6.

Таблица 2.6 – Описание таблицы задачи

Поле	Тип данных	Ограничение	Описание
id_task	Целочисленный	Уникальное, обязательное	Первичный ключ таблицы
name	Символьный	Максимальная длина — 255 символов, обязательное	Название задачи
description	Символьный	Максимальная длина — 255 символов	Описание задачи
status	Символьный	Максимальная длина — 255 символов, обязательное	Статус задачи
sum_timer	Целочисленный	—	Общее время выполнения задачи
start_timer	Дата	—	Время начала выполнения задачи
timer_add_task	Дата	—	Дата начала выполнения задачи
time_end_task	Дата	—	Дата завершения выполнения задачи
table_id	Целочисленный	Обязательное	Идентификатор проекта к которому относится таблица задач, внешний ключ

Информация о полях в таблице тегов (tag) представлена в таблице 2.7.

Таблица 2.7 – Описание таблицы тегов

Поле	Тип данных	Ограничение	Описание
id_tag	Целочисленный	Уникальное, обязательное	Первичный ключ таблицы
name	Символьный	Максимальная длина — 255 символов, обязательное	Название тега
color	Символьный	Максимальная длина — 255 символов	цвет тега
project_id	Целочисленный	Обязательное	Идентификатор проекта к которой относится тег, внешний ключ

Информация о полях в таблице связи задач с тегами (task_tag) представлена в таблице 2.8.

Таблица 2.8 – Описание таблицы связи задачи и тега

Поле	Тип данных	Ограничение	Описание
id	Целочисленный	Уникальное, обязательное	Первичный ключ таблицы
tag_id	Целочисленный	Обязательное	Идентификатор тега, внешний ключ
task_id	Целочисленный	Обязательное	Идентификатор задачи, внешний ключ

Информация о полях в таблице групп (group_user) представлена в таблице 2.9.

Таблица 2.9 – Описание таблицы групп

Поле	Тип данных	Ограничение	Описание
id	Целочисленный	Уникальное, обязательное	Первичный ключ таблицы
name	Символьный	Максимальная длина — 255 символов, обязательное	Название группы
description	Символьный	Максимальная длина — 255 символов	Описание группы
project_id	Целочисленный	Уникальное, обязательное	Идентификатор проекта, внешний ключ

Информация о полях в таблице связи группы и пользователя (group_member) представлена в таблице 2.10.

Таблица 2.10 – Описание таблицы связи группы и пользователя

Поле	Тип данных	Ограничение	Описание
id	Целочисленный	Уникальное, обязательное	Первичный ключ таблицы
role	Символьный	Максимальная длина — 255 символов, обязательное	Роль пользователя в группе
user_id	Целочисленный	Обязательное	Идентификатор пользователя, внешний ключ
group_id	Целочисленный	Обязательное	Идентификатор группы, внешний ключ

На рисунке 2.1 приведена диаграмма разрабатываемой базы данных.

2.2 Требования к программному обеспечению

К программе предъявлен ряд требований:

- должен присутствовать интерфейс для выбора действий;

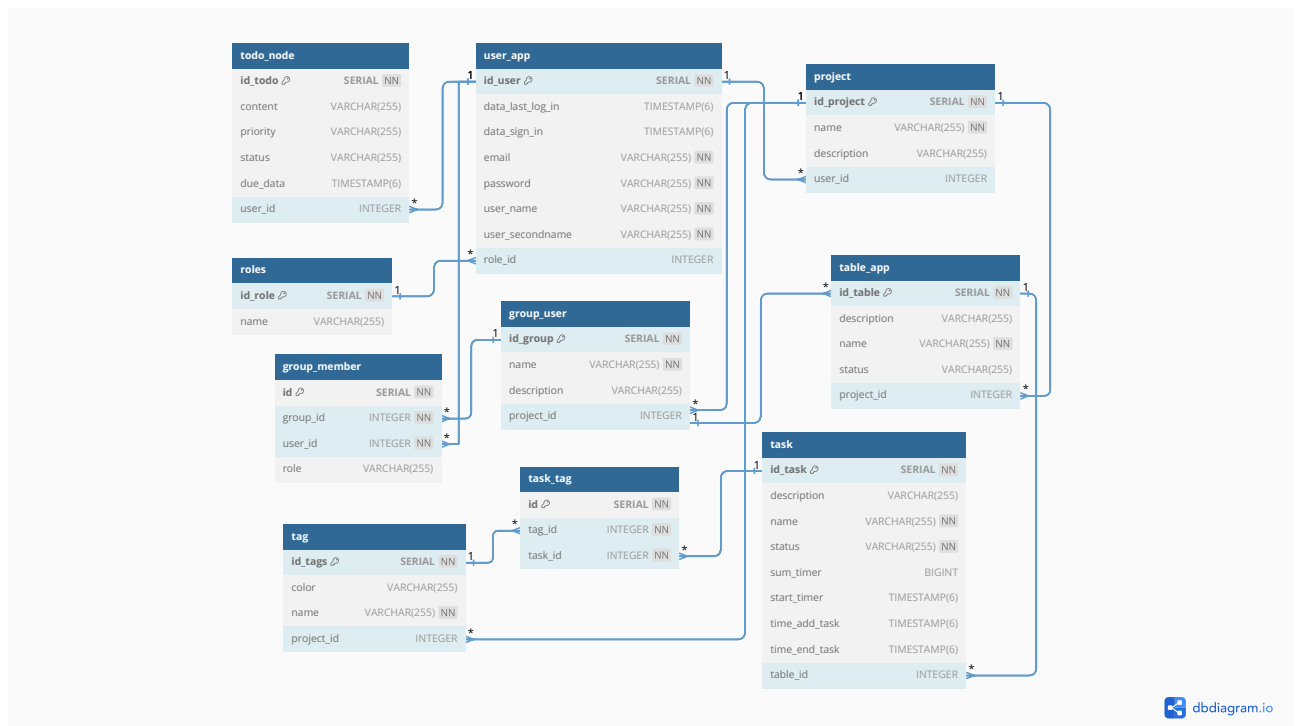


Рисунок 2.1 – Диаграмма разрабатываемой базы данных

- считывание данных должно производиться из файла;
- результат должен записываться в файл;
- должен присутствовать замер реального времени для реализаций алгоритмов;
- результат замера должен выводиться в виде таблицы.

2.3 Описание используемых типов данных

При реализации алгоритмов будут использованы следующие структуры и типы данных:

- массив символов для хранения терма;
- вещественное число для хранения $TF-IDF$ терма;
- мьютекс — примитив синхронизации.

2.4 Разработка алгоритмов

На рисунке ?? приведена схема дивизимной иерархической кластеризации. На рисунке ?? приведена схема алгоритма k-средних. На рисунке ??

представлена схема алгоритма создания потоков для алгоритма дивизимной иерархической кластеризации. На рисунке ?? представлена схема многопоточного алгоритма кластеризации документов.

Вывод

В данном разделе были представлены требования к программному обеспечению, описание используемых типов данных и схемы реализуемых алгоритмов.

3 Технологический раздел

В данном разделе будут представлены средства реализации, сведения о модулях программы и листинги кода реализации алгоритмов.

3.1 Средства реализации

В качестве языка программирования для разработки программного обеспечения был выбран язык *C++* [**cpp**].

Данный выбор обусловлен следующим:

- язык позволяет реализовать все алгоритмы, выбранные в результате проектирования;
- язык поддерживает все типы и структуры данных, которые были выбраны в результате проектирования;
- язык позволяет работать с нативными потоками [**thread**].

Время выполнения реализаций было замерено с помощью функции *clock* [**clock**]. Для хранения термов использовалась структура данных *string* [**wstring**], в качестве массивов использовалась структура данных *vector* [**vector**]. В качестве примитива синхронизации использовался *mutex* [**mutex**].

Для создания потоков и работы с ними был использован класс *thread* из стандартной библиотеки выбранного языка [**thread**]. В листинге 3.1, приведен пример работы с описанным классом, каждый объект класса представляет собой поток операционной системы, что позволяет нескольким функциям выполняться параллельно [**thread**].

Листинг 3.1 – Пример работы с классом thread

```
1 #include <iostream>
2 #include <thread>
3
4 void foo(int a){
5     std::cout << a << '\n';
6 }
7
8 int main(){
9     std::thread thread(foo, 10);
10    thread.join();
11
12    return 0;
13 }
```

3.2 Сведения о модулях программы

Данная программа разбита на следующие модули:

- *main.cpp* — файл, который содержит точку входа в программу;
- *Cluster.cpp* и *Cluster.h* — модуль, который реализует класс *Cluster*;
- *kMeans.cpp* и *kMeans.h* — модуль, содержащий реализацию функции *kMeans*;
- *Klustering.cpp* и *Klustering.h* — модуль, содержащий реализации функций дивизимной иерархической кластеризации: *hierarchicalClustering* и многопоточная версия кластеризации *parallelDocClustering*;
- *io.cpp* и *io.h* — модуль, содержащий реализации функций для работы входными и выходными файлами;
- *Document.cpp* и *Document.h* — модуль, который реализует класс *Document*;
- *Term.cpp* и *Term.h* — модуль, который реализует класс *Term*;

3.3 Реализация алгоритмов

В листинге 3.2 приведена реализация дивизимной иерархической кластеризации. В листинге A.1 приведена реализация алгоритма k-средних. В

листинге 3.3 приведена реализация алгоритма создания потоков для алгоритма дивизимной иерархической кластеризации. В листинге 3.4 приведена реализация многопоточного алгоритма кластеризации документов.

Листинг 3.2 – Реализация дивизимной иерархической кластеризации

```
1 void hierarchicalClustering(const
    std::vector<std::vector<double>> &docs,
2         int depth,
3         int maxDepth,
4         std::vector<Cluster> &clusters) {
5     if (depth == maxDepth) {
6         Cluster cluster;
7         cluster.docs = docs;
8         cluster.calculateCentroid();
9         clusters.push_back(cluster);
10        return;
11    }
12
13    // Применяем k-средних для разделения на 2 кластера
14    std::vector<std::vector<double>> centroids = kMeans(docs, 2);
15    std::vector<std::vector<double>> cluster1Docs, cluster2Docs;
16
17    for (const auto &doc: docs) {
18        if (cosineDistance(doc, centroids[0]) <
19            cosineDistance(doc, centroids[1])) {
20            cluster1Docs.push_back(doc);
21        }
22        else {
23            cluster2Docs.push_back(doc);
24        }
25    }
26
27    // Рекурсивно разделяем каждый подкластер
28    hierarchicalClustering(cluster1Docs, depth + 1, maxDepth,
29        clusters);
30    hierarchicalClustering(cluster2Docs, depth + 1, maxDepth,
31        clusters);
32 }
```

Листинг 3.3 – Реализация алгоритма создания потоков

```
1 void parallelDocClustering(const std::vector<std::string>
    &filenames,
2
    const std::vector<std::string>
    &filenames_out,
3
    int num_threads) {
4     int count_f = filenames.size();
5     std::mutex io_mutex;
6
7     // Расчет размера подгруппы для каждого потока
8     int chunk_size = count_f / num_threads;
9     std::thread threads[num_threads]; // Массив потоков
10
11     int start_index = 0;
12     for (int i = 0; i < num_threads; ++i) {
13         int end_index = start_index + chunk_size + (i < count_f
            % num_threads ? 1 : 0);
14
15         // Создание потока для обработки подгруппы документов
16         threads[i] = std::thread(processDocuments,
            std::ref(filenames), std::ref(filenames_out),
17
            start_index, end_index,
            std::ref(io_mutex));
18
19         start_index = end_index;
20     }
21
22     // Дождаться завершения всех потоков
23     for (int i = 0; i < num_threads; ++i) {
24         threads[i].join();
25     }
```

Листинг 3.4 – Реализация многопоточного алгоритма кластеризации документов

```
1 void processDocuments(const std::vector<std::string> &filenames,
2                       const std::vector<std::string>
3                         &filenames_out,
4                       int start_index, int end_index,
5                       std::mutex &io_mutex) {
6     for (int i = start_index; i < end_index; ++i) {
7         // Загрузка и обработка документа
8         std::vector<Document> documents =
9             readDocumentsFromCSV(filenames[i]);
10        standardizeDocumentLengths(documents);
11
12        // Выполнение иерархической кластеризации
13        std::vector<Cluster> clusters =
14            hierarchicalClustering(documents, MAX_DEPTH);
15
16        // Синхронизированное сохранение результатов
17        std::lock_guard<std::mutex> lock(io_mutex);
18        saveClustersToJson(clusters, documents,
19                            filenames_out[i]);
20    }
21 }
```

Вывод

В данном разделе были представлены средства реализации, сведения о модулях программы и листинги кода реализации алгоритмов.

4 Исследовательская часть

В данном разделе будут приведены демонстрация работы программы, технические характеристики устройства, сравнительный анализ времени выполнения реализуемых алгоритмов.

4.1 Демонстрация работы программы

На рисунке ?? представлена демонстрация работы разработанного программного обеспечения.

4.2 Технические характеристики

Технические характеристики устройства, на котором выполнялись замеры по времени, представлены далее.

- Процессор: Ryzen 5 4600H, 6 процессорных ядер архитектуры Zen 2 и 12 потоков, работающих на базовой частоте в 3.0 ГГц (до 4.0 ГГц в Turbo режиме), 12 логических ядер [ryzen]
- Оперативная память: 16 ГБайт.
- Операционная система: Windows 10 Pro 64-разрядная система [windows].

При замерах времени ноутбук был включен в сеть электропитания и был нагружен только системными приложениями.

4.3 Время выполнения реализаций алгоритмов

Результаты замеров времени выполнения реализации алгоритма иерархической кластеризации документов в зависимости от числа потоков приведены в таблице Б.1. Каждый замер проводился 100 раз, после чего рассчитывалось их среднее арифметическое значение.

На рисунке ?? изображен график зависимостей времени выполнения реализаций от числа потоков.

Результаты замеров времени выполнения реализации алгоритма иерархической кластеризации документов в зависимости от числа входных файлов приведены в таблице Б.2 В таблице Б.2 в многопоточной реализации был

использован 1 вспомогательный поток, данное число было выбрано для демонстрации уменьшения времени получения результата при использовании многопоточности с минимальным числом вспомогательных потоков. Каждый замер проводился 100 раз, после чего рассчитывалось их среднее арифметическое значение.

В таблице Б.2 используются следующие обозначения:

- «ПР» — время выполнения (в мс) последовательной реализации алгоритма;
- «ПРОДП» — время выполнения (в мс) параллельной реализации алгоритма при одном дополнительном потоке.

На рисунке ?? изображен график зависимости времени выполнения реализации от числа файлов.

Вывод

В результате исследования реализуемых алгоритмов по времени выполнения можно сделать следующие выводы.

- Время выполнения последовательной реализации алгоритма остается постоянным и равным 4200 мс (см. рисунок ??). На графике видно, что для небольшого числа потоков (до 8), параллельная реализация работает быстрее последовательной. Однако, с увеличением количества потоков, параллельная реализация алгоритма становится менее эффективной по времени выполнения, чем последовательная реализация алгоритма. Это можно объяснить тем, что на машине, на которой проводились исследования 12 логических ядер. При увеличении количества потоков сверх количества логических ядер, происходит конкуренция за процессорные ресурсы, что приводит к переключению контекста и, как следствие, к увеличению времени выполнения. Кроме того, на создание и управление потоками тратится дополнительное время, что также влияет на общую производительность.
- Из таблицы Б.2 видно, что при малом количестве входных файлов разница во времени выполнения между последовательной реализацией и

параллельной реализацией с одним дополнительным потоком, не очень заметна, но с увеличением количества файлов параллельная версия с одним дополнительным потоком становится намного эффективней по времени выполнения. Параллельная реализация алгоритма с одним дополнительным потоком на 38% быстрее по сравнению с последовательной реализацией алгоритма, при обработке 50 файлов (см. таблицу Б.2).

- Исходя из данных (см. таблицу Б.2), оптимальное количество потоков для данной задачи находится в диапазоне от 2 до 8 (см. таблицу 4.1). В этом диапазоне время выполнения минимально и значительно меньше, чем у последовательной реализации алгоритма.

Таблица 4.1 – Относительное улучшение во времени выполнения параллельной реализации алгоритма по сравнению с последовательной реализацией алгоритма

Количество потоков	Относительное улучшение
2	2.6%
3	2.9%
4	3.65%
5	3.21%
6	3.18%
7	2.97%
8	1.8%

ЗАКЛЮЧЕНИЕ

Цель лабораторной работы достигнута, исследованы параллельные вычисления на основе нативных потоков.

Для достижения поставленной цели были выполнены следующие задачи:

- описан алгоритм иерархической кластеризации;
- спроектировано программное обеспечение, реализующее алгоритм и его параллельную версию;
- разработано программное обеспечение, реализующее алгоритм и его параллельную версию;
- проанализированы затраты реализаций алгоритмов по времени.

В результате исследования реализуемых алгоритмов по времени выполнения можно сделать следующие выводы.

- Время выполнения последовательной реализации алгоритма остается постоянным и равным 4200 мс (см. рисунок ??). На графике видно, что для небольшого числа потоков (до 8), параллельная реализация работает быстрее последовательной. Однако, с увеличением количества потоков, параллельная реализация алгоритма становится менее эффективной по времени выполнения, чем последовательная реализация алгоритма. Это можно объяснить тем, что на машине, на которой проводились исследования 12 логических ядер. При увеличении количества потоков сверх количества логических ядер, происходит конкуренция за процессорные ресурсы, что приводит к переключению контекста и, как следствие, к увеличению времени выполнения. Кроме того, на создание и управление потоками тратится дополнительное время, что также влияет на общую производительность.
- Из таблицы Б.2 видно, что при малом количестве входных файлов разница во времени выполнения между последовательной реализацией и параллельной реализацией с одним дополнительным потоком, не очень заметна, но с увеличением количества файлов параллельная версия с одним дополнительным потоком становится намного эффективней по

времени выполнения. Параллельная реализация алгоритма с одним дополнительным потоком на 38% быстрее по сравнению с последовательной реализацией алгоритма, при обработке 50 файлов (см. таблицу Б.2).

- Исходя из данных (см. таблицу Б.2), оптимальное количество потоков для данной задачи находится в диапазоне от 2 до 8 (см. таблицу 4.1). В этом диапазоне время выполнения минимально и значительно меньше, чем у последовательной реализации алгоритма.

ПРИЛОЖЕНИЕ А

Листинг А.1 – Реализация алгоритма k-средних

```
1 std::vector<std::vector<double>> kMeans(const
    std::vector<std::vector<double>> &docs, int k) {
2     int n = docs.size(); //1
3     std::vector<std::vector<double>> centroids(k,
        std::vector<double>(docs[0].size())); //2
4     std::vector<int> assignments(n, 0); //3
5     // Инициализация центроидов
6     for (int i = 0; i < k; ++i) { // 4
7         centroids[i] = docs[rand() % n]; //5
8     }
9     bool changed; // 6
10    do { // 7
11        changed = false; // 8
12        // Назначение точек кластерам
13        for (int i = 0; i < n; ++i) { //9
14            double bestDist = -1.0; //10
15            int bestCluster = 0; //11
16            for (int j = 0; j < k; ++j) { //12
17                double dist = cosineDistance(docs[i],
                    centroids[j]); // 13
18                if (dist > bestDist) { //14
19                    bestDist = dist; //15
20                    bestCluster = j; //16
21                }
22            }
23            if (assignments[i] != bestCluster) { //17
24                assignments[i] = bestCluster; //18
25                changed = true; //19
26            }
27        }
28        // Обновление центроидов
29        std::vector<int> counts(k, 0); //20
30        std::vector<std::vector<double>> newCentroids(k,
            std::vector<double>(docs[0].size(), 0.0)); //21
31        for (int i = 0; i < n; ++i) { //22
32            for (size_t j = 0; j < docs[i].size(); ++j) { //23
33                newCentroids[assignments[i]][j] +=
```

```

34         docs[i][j]; //24
35     }
36     counts[assignments[i]]++; //25
37 }
38 for (int j = 0; j < k; ++j) { //26
39     if (counts[j] != 0) { //27
40         for (size_t m = 0; m < newCentroids[j].size();
41             ++m) { //28
42             newCentroids[j][m] /= counts[j]; //29
43         }
44     }
45     else { //30
46         newCentroids[j] = centroids[j]; //31
47     }
48     centroids = newCentroids; //32
49 } while (changed); //33
50 return centroids; //34
51 }

```

ПРИЛОЖЕНИЕ Б

Таблица Б.1 – Зависимость времени выполнения (в мс) от количества потоков

Количество потоков	Время выполнения (мс)
1	4271
2	4160
3	4146
4	4115
5	4134
6	4135
7	4144
8	4194
9	4200
10	4178
11	4207
12	4209
13	4242
14	4205
15	4240
16	4269
17	4240
18	4283
19	4287
20	4314
21	4317
22	4304
23	4454
24	4371
25	4426
26	4355
27	4376

28	4368
29	4426
30	4447
31	4503
32	4439
33	4468
34	4453
35	4476
36	4498
37	4521
38	4610
39	4615
40	4598
41	4577
42	4537
43	4564
44	4614
45	4606
46	4579
47	4893
48	4799
49	4819
50	4823
51	4825
52	4812
53	4818
54	4836
55	4849
56	4838
57	4834
58	4840
59	4815
60	4820

61	4772
62	4802
63	4806
64	4810
65	4826
66	4836
67	4823
68	4834
69	4850
70	4842
71	4857
72	4870
73	4874
74	4900
75	5056
76	5012
77	5004
78	5002
79	4947
80	4934
81	4915
82	4940
83	4944
84	4948
85	4959
86	4966
87	4978
88	4992
89	5149
90	5138
91	5126
92	5154
93	5158

94	5159
95	5161
96	5174

Таблица Б.2 – Зависимость времени выполнения (в мс) от количества входных файлов

Количество файлов	ПР(мс)	ПРОДП(мс)
1	44	47
2	88	70
3	132	108
4	178	118
5	225	153
6	271	178
7	327	220
8	359	220
9	403	249
10	445	273
11	489	310
12	542	327
13	582	363
14	623	375
15	666	424
16	729	426
17	765	480
18	811	470
19	850	516
20	909	529
21	940	564
22	987	598
23	1028	603
24	1076	631

25	1129	664
26	1159	682
27	1207	718
28	1247	729
29	1289	759
30	1342	766
31	1375	823
32	1433	803
33	1483	982
34	1531	924
35	1561	906
36	1600	928
37	1646	993
38	1693	994
39	1735	1042
40	1785	1090
41	1829	1105
42	1871	1149
43	1942	1155
44	1977	1151
45	2023	1206
46	2061	1218
47	2109	1297
48	2181	1305
49	2490	1441
50	2256	1398

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. Time Management - A Case Study [Электронный ресурс]. — Режим доступа: https://www.researchgate.net/publication/323825275_Time_Management_-_A_Case_Study (дата обращения: 17.03.2024).
2. Introducing the Eisenhower Matrix [Электронный ресурс]. — Режим доступа: <https://www.eisenhower.me/eisenhower-matrix/> (дата обращения: 17.03.2024).
3. *Брайан Трейси*. Тайм-менеджмент. — 2015. — С. 1—144.