



МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ
РОССИЙСКОЙ ФЕДЕРАЦИИ

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ БЮДЖЕТНОЕ
ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ
«МОСКОВСКИЙ АВИАЦИОННЫЙ ИНСТИТУТ
(национальный исследовательский университет)»

Институт (Филиал) № 8 «Компьютерные науки и прикладная математика» Кафедра 806

Группа М8О-407Б-19 Направление подготовки 01.02.03 «Прикладная математика и информатика»

Профиль Информатика

Квалификация: бакалавр

ВЫПУСКНАЯ КВАЛИФИКАЦИОННАЯ РАБОТА БАКАЛАВРА

на тему «Сервис поиска, классификации и обработки тематических изображений в интернете»

Автор ВКРБ: Волков Матвей Андреевич ()

Руководитель: Филимонов Николай Сергеевич ()

Консультант: ()

Консультант: ()

Рецензент: ()

К защите допустить

Заведующий кафедрой № 806 «Вычислительная математика
и программирование»

Крылов Сергей Сергеевич ()

____ мая 2023 года

Москва 2023

РЕФЕРАТ

Выпускная квалификационная работа бакалавра состоит из 30 страниц, 16 рисунков, 9 использованных источников, 2 приложений.

МАНГА, БАЗА ДАННЫХ, МИКРОСЕРВИС, ПАРСИНГ, САЙТ, API

Работа состоит из введения, двух глав и заключения. Темой работы является веб приложение, основанное на технологии парсинга данных для сбора тематических картинок (манги) и компилирование их в сборники в одном pdf файле.

Целью работы является написание API, которое возможно использовать в других клиентских сервисах, таких как мессенджер, сайт. API должен предоставлять интерфейс для отправки собранных изображений на электронную книгу через отправку тех по электронной почте.

Проведение работы состоит из исследования целевых сайтов, список которых задается заранее. Написание парсинг алгоритма. Конструирование микросервиса, применяя методiku чистой архитектуры для дальнейшего развития проекта.

Областью применения является любой клиентский сервис, способный общаться с сделанным API.

Итог внедрения считается успешно работающее API любым из возможных клиентских сервисов, таких как сайт, мессенджер.

Дальнейшее развитие проекта предполагает оформление подписки на сайты для автоматического информирования пользователя о новых главах комиксов или манги. Оптимизирование работы сервиса для успешного обслуживания больше десятка пользователей.

Появление подобного микросервиса должно принести в рынок Open source новое кроссплатформенное решение по оффлайн прочтению манги, в том числе, и на электронных устройствах с искусственными чернилами (электронные книги).

СОДЕРЖАНИЕ

ТЕРМИНЫ И ОПРЕДЕЛЕНИЯ	4
ПЕРЕЧЕНЬ СОКРАЩЕНИЙ И ОБОЗНАЧЕНИЙ	6
ВВЕДЕНИЕ	7
1 ТЕОРЕТИЧЕСКАЯ ЧАСТЬ	8
1.1 Обоснование библиотек	8
1.1.1 Pdfcpu	8
1.1.2 Goquery	9
1.1.3 Jet	10
1.2 Устройство сайтов	11
1.2.1 HTML	12
1.2.2 CSS	13
1.2.3 JavaScript	13
1.3 Методика парсинга	13
1.3.1 Взаимодействие с сервером	14
1.3.2 Изучение данных	15
1.3.3 Сайты с API	21
1.4 База данных в микросервисах	22
1.5 Чистая архитектура	23
1.5.1 Представление (Presentation Layer)	25
1.5.2 Использование (Use Case Layer)	25
1.5.3 Хранилище (Storage Layer)	25
1.5.4 Инфраструктура (Infrastructure Layer)	25
2 ПРАКТИЧЕСКАЯ ЧАСТЬ	26
2.1 Структура приложения	26
ЗАКЛЮЧЕНИЕ	27
СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ	28
ПРИЛОЖЕНИЕ А Приложение к изучению данных	29
ПРИЛОЖЕНИЕ Б Приложение к чистой архитектуре	30

ТЕРМИНЫ И ОПРЕДЕЛЕНИЯ

В настоящей выпускной квалификационной работе бакалавра применяют следующие термины с соответствующими определениями:

HTML — стандартизированный язык гипертекстовой разметки документов для просмотра веб-страниц в браузере.

HTTP — протокол передачи гипертекста, изначально в виде гипертекстовых документов в формате HTML, в настоящее время используется для передачи произвольных данных.

DOM-дерево — структурное представление HTML файла.

API — описание способов взаимодействия одной компьютерной программы с другими.

Метод HTTP — последовательность из любых символов, кроме управляющих и разделителей, указывающая на основную операцию над ресурсом. Обычно метод представляет собой короткое английское слово, записанное заглавными буквами.

HTTP запрос — сигнал серверу через http метод, всегда требующий какого-то ответа от сервера.

Хендлер — функция, которую выполняет веб-приложение, когда был сделан соответствующий http запрос.

База данных — набор информации, которая хранится упорядоченно в электронном виде.

Сервер — сервером в веб приложении называют программу, которая занимается обработкой запросов с клиентской части приложения. В основном на сервере происходит вся логика работы веб приложения. Обычно на серверной части запускается база данных.

Клиент — приложение, которое путем отправки запросов на сервер способно визуализировать полученную информацию. Обычно приложение клиент общается с пользователем. При этом не присутствует никакой логики, связанной с основным приложением. Например, браузер является клиентом.

Библиотека — сборник подпрограмм или объектов, используемых для разработки программного обеспечения.

Парсинг — это автоматизированный сбор и структурирование информации с сайтов при помощи программы или сервиса.

jQuery — набор функций JavaScript, фокусирующийся на взаимодействии

JavaScript и HTML. Библиотека jQuery помогает легко получать доступ к любому элементу DOM, обращаться к атрибутам и содержимому элементов DOM, манипулировать ими.

Консистентность — согласованность данных друг с другом, целостность данных, а также внутренняя непротиворечивость.

Миграция базы данных — переход от одной структуры базы данных к другой без потери консистентности.

Пайплайн — конвейер данных, поступающий из одной логической программы другой.

Чистая архитектура — понятие при конструировании микросервиса, созданное для разделения различных концепций, путем написания кода на нескольких уровнях с четким правилом, которое позволяет создать тестируемый и поддерживаемый проект.

GO — язык программирования, созданный компанией гугл преимущественно для создания бекэнд сервисов.

Модуль — называется пакет в проекте языка GO, который может быть переиспользован в других приложениях или проектах, как внешняя библиотека. Был добавлен в версии 1.11.

Пакет — понятие в языке GO, обозначающее коллекцию исходного кода в одной директории скомпилируемых вместе. Пакет также является подключаемым модулем.

Реверс-инженеринг — исследование некоторого готового устройства или программы, а также документации на него с целью понять принцип его работы. Также известно, как «обратная разработка».

SQL-инъекция — один из распространённых способов взлома сайтов и программ, работающих с базами данных, основанный на внедрении в запрос произвольного SQL-кода.

Dos-атака — хакерская атака на вычислительную систему с целью довести её до отказа. Если атака выполняется одновременно с большого числа компьютеров, говорят о DDoS-атаке

Манга — японские комиксы, иногда называемые комикку.

Транзакция — это набор операций по работе с базой данных, объединенных в одну атомарную пачку.

ПЕРЕЧЕНЬ СОКРАЩЕНИЙ И ОБОЗНАЧЕНИЙ

В настоящей выпускной квалификационной работе бакалавра применяют следующие сокращения и обозначения:

SQL — Structured Query Language

API — Application Programming Interface

БД — база данных

HTTP — HyperText Markup Language

DOM — Document Object Model

ЯП — язык программирования

PDF — Portable Document Format

ID — Identification

JSON — Javascript object notation

RegExp — Regular Expression

CSS — Cascading Style Sheets

JS — JavaScript

DDOS — Distributed Denial of Service

ВВЕДЕНИЕ

Сегодня манга является одним из самых популярных видов комиксов в мире, и ее популярность только растет. Однако, не всегда удобно читать мангу на экране компьютера или телефона. Именно поэтому возникает необходимость в сервисе, который позволяет скачивать мангу и отправлять ее на электронные книги. В данном случае мы говорим о микросервисе, который предоставляет API к скачиванию манги и отправке ее на устройства чтения. Задача данного сервиса заключается в том, чтобы обеспечить удобный и быстрый доступ к любимой манге, а также сделать ее чтение более комфортным для пользователей.

Плюсы API заключаются в том, что это, в первую очередь, кроссплатформенное решение. Исследовав рынок подобных сервисов, не было обнаружено чего-то удобного для операционных систем IOS. Из-за открытости системы android удалось найти похожее решение, однако удобство использования оставляло желать лучшего: есть привязанность к устройству, нет никакой синхронизации и прочих удобств современных приложений.

Чтение манги, используя электронную книгу — совершенно новое направление, которое открывает возможности для любителей манги. Так люди, которые обожают печатные издания, смогут насладиться технологией искусственных чернил, при этом не таская с собой в дорогу огромную библиотеку комиксов. Скачивать мангу на электронную книгу — очень трудная задача. С ней не справится пользователь, который не умеет писать даже базовые скрипты на компьютере. Сервис же будет предоставлять удобный API, который возможно подключить к любому клиентскому сервису.

Таким образом, скаченную мангу возможно хранить в мессенджере, электронной книге, наслаждаясь всеми прелестями технологии электронных чернил, а также, возможно сделать свое приложение со своим удобным оформлением.

1 ТЕОРЕТИЧЕСКАЯ ЧАСТЬ

1.1 Обоснование библиотек

В микросервисе были использованы внешние модули для упрощения частоиспользуемых действий в коде. Каждая из них представляет собой очень полезную функциональность, которая будет перечислена ниже с обоснованием использования.

1.1.1 Pdfcpu

Pdfcpu [1] — пакет, написанный на ЯП GO, являющийся мощным сборником pdf файлов. Представляет интерфейс для создания сложноструктурированных pdf файлов.

В данном пакете возможно создавать абсолютно любые pdf файлы с различными метаданными документа, начиная с содержания, заканчивая указанием автора написания документа. Библиотека способна также использовать различные картинки различных разрешений, позволяя менять их размер, цвет, прозрачность и другие параметры.

Пакет был взят в использование потому, что предоставляется разработчиками удобный и понятный интерфейс для собирания большого количество картинок в pdf файл. Также модуль отличается своей быстротой исполнения поставленных задач. Рисунок 1 может продемонстрировать простоту использования пакета.


```

1 func CreatePDFFromImagesDir(imagesDirPath string, outputPath string
  ) error {
2   imagesPath, err := GetImagesPathStr(imagesDirPath)
3   if err != nil {
4     return errors.Wrap(err, "something_wrong_with_creating_images_
      path")
5   }
6   err = api.ImportImagesFile(imagesPath, outputPath, nil, nil)
7   if err != nil {
8     return errors.Wrap(err, "pdf_creation_filure")
9   }
10  return nil
11 }

```

Рисунок 1 – Демонстрация работы интерфейса pdfcpu

1.1.2 Goquery

Goquery [2] — пакет, написанный на ЯП GO, который способен выполнять задачи библиотеки jQuery [3].

Данный пакет используется для того, чтобы парсить DOM дерево полученного HTML файла. С ним возможно быстро и удобно получить нужную информацию с любого сайта. После реверс-инженеринга, благодаря этому пакету, быстро и просто достается информация о картинках, классифицируется на разделы и подразделы (например, группировка на главы).

```

1 func (mlb mangaLibController) getChapterID(doc *goquery.Document)
  (string, error) {
2   result, ok := doc.Find("#comments").Attr("data-post-id")
3   if !ok {
4     return "", errors.Wrap(customerrors.ErrEmptyAttr, "mangalib:_
      data-post-id")
5   }
6   return result, nil
7 }
8

```

Рисунок 2 – Демонстрация работы пакета goquery

Как видно из рисунка 2, используются команды jQuery для получения информации о главе. Впоследствии эта глава используется для того, чтобы достать из неё картинки.

1.1.3 Jet

Go-jet [4] — пакет, написанный на ЯП GO, являющийся конструктором SQL запросов в БД.

С помощью этой библиотеки возможно собирать любые SQL запросы, избегая уязвимость — SQL-инъекцию. Этот пакет подключается к БД и генерирует код на языке GO, при помощи которого можно собирать различные SQL запросы. От инъекции помогает избавиться тот факт, что абсолютно весь запрос делается при помощи конструктора, в процессе сбора которого есть постоянные проверки на соответствие данных, поступающих в процесс создания запроса. Также стоит заметить, что сама возможность собирать запрос из различных кусочков, может очень сильно помочь при составлении динамического запроса, который так или иначе зависит от поступивших на вход параметров.

Именно по этой причине был выбран данный пакет. Как можно видеть из рисунка 3, собирается SQL запрос в зависимости от пришедших на вход параметров. Таким образом, условная секция запроса может варьироваться.

```

1      func (perC personController) GetEmailByID(ctx context.Context,
person domain.PersonInfo) (string, error) {
2          selectStmt := table.Persons.SELECT(table.Persons.Email)
3          var whereStmt postgres.BoolExpression
4
5          personID, _ := uuid.Parse(person.PersonID)
6          switch {
7              case personID != uuid.UUID{} && personID != uuid.Nil:
8                  whereStmt = table.Persons.ID.EQ(postgres.UUID(personID))
9              case person.TelegramID > 0:
10                 whereStmt = table.Persons.TelegramID.EQ(postgres.Int64(
person.TelegramID))
11             default:
12                 return "", customerrors.ErrEmailsNotFound
13             }
14             stmt, args := selectStmt.WHERE(whereStmt).Sql()
15
16             ...
17
18             return email, nil
19         }
20

```

Рисунок 3 – Демонстрация работы пакета jet

1.2 Устройство сайтов

Предметом исследования диплома был ограниченный список сайтов, участвующих в сборе картинок с последующей их сортировкой и группировкой в pdf файл для последующего их удобного просмотра или прочтения. Суть изучения состояло в том, чтобы произвести реверс-инженеринг сайта, понять как он работает. Затем использовать полученные знания в поиске нужной мне информации, путем парсинга страницы сайта, разбиения полученной структуры сайта на составляющие, с последующей структуризацией в нужный мне формат для правильной отдачи итогового результата.

Абсолютно любой сайт состоит из 3-х компонент:

- a) HTML
- б) CSS

в) JavaScript

1.2.1 HTML

В данной компоненте описывается различная информация о странице. Будь то текст статьи, картинки и прочего. Затем, браузер, применяя свои движки визуализации, преобразует информацию, содержащуюся в компоненте, в понятное для человека визуализированное представление. HTML имеет свою структуру, называемой DOM деревом. Эта структура по сути своей является набором различных тэгов. Для идеального примера можно сказать, что HTML похож на LATEX. Таким образом, HTML — это то, откуда будет получаться нужная информация.

Как было сказано ранее, HTML состоит из тэгов. Есть следующие виды тэгов:

- а) Двойные
- б) Ординарные

Двойной тэг отличается от ординарного тем, что у двойного тэга есть закрывающий тэг. Между открывающим тэгом и закрывающим пишется текстовая информация. Например, на рисунке 4 продемонстрирован двойной тэг, обозначающий параграф. Обычно такой тэг используется в различных статьях. На рисунке 5 продемонстрирован пример ординарного тэга, обозначающий картинку. Когда браузер видит этот тэг, он понимает, что перед ним картинка. За кадром происходит запрос на сервер по ссылке, расположенной в мета информации, с последующим отображением изображения.

```
1 <p class="paragraf">Текст параграфа</p>
2
```

Рисунок 4 – Пример двойного тэга

```
1 
2
```

Рисунок 5 – Пример ординарного тэга

На картинках 4 и 5 были продемонстрированы тэги, которые требуют от

браузера некоторых действий. Но, как было выше отмечено, HTML очень похож на LATEX. И похож он тем, что в основном страницы сайта нужно верстать. Поэтому люди придумали тэг, который не требует никаких действий со стороны браузера, но при этом будет как-то объединять информацию.

На рисунке 6 изображен тэг `div`. Этот тэг нужен исключительно для того, чтобы как-то объединять информацию, присваивая ей какую-то общую метainформацию, например, класс.

```
1 <div class="pritty"></div>  
2
```

Рисунок 6 – Тэг `div`

Из совокупности тэгов состоит страница сайта. А структура, которая в конечном итоге образуется, называется DOM деревом.

1.2.2 CSS

Эта компонента, в свою очередь, отвечает за то как будет выглядеть информация, содержащаяся в HTML. По сути своей, CSS — оформление страницы. После появления данной компоненты, сайты в интернете стали выглядеть так, как мы их сейчас видим.

Для формирования различных стилей, используется метainформация тэгов. Удобнее всего для таких нужд оказалось использовать тэг, изображенный на рисунке 6, так как этот тэг как никто лучше служит для объединения информации, которую можно как-то по-особенному расположить и украсить.

Именно поэтому на большинстве сайтов можно увидеть именно этот тэг.

1.2.3 JavaScript

Эта компонента отвечает за логику на страницах. Логика может использоваться при анимировании сайтов и общении с сервером. Общение с сервером нужно для того, чтобы возможно было у него получить какую-то информацию.

1.3 Методика парсинга

Понимание сайта — важная часть абсолютно любого парсинга данных.

Как понятно из определения, необходимо собирать информацию. Для того, чтобы собирать информацию, необходимо знать как она расположена. Для того, чтобы понимать как расположена информация, необходимо исследовать пути появления этой информации перед пользователем и выяснить откуда она берется.

Плавню мы подбираемся к тому, что знание устройства сайта, клиентской части и сервера — необходимы для ранее сказанного парсинга данных. Другими словами, для парсинга данных с сайта необходимо полностью понимать как тот устроен. Если же не получается полностью изучить, то необходимо понять, какие данные должны использоваться, чтобы в конечном итоге достичь нужного результата. Нужно цепляться буквально за любую крупину данных, которая может так или иначе помочь.

1.3.1 Взаимодействие с сервером

Для поиска подобной крупинки информации нужно понимать, что ничего не происходит из ниоткуда. Везде есть какие-то следы. Они могут быть зашифрованы, или спрятаны за тонной других запросов. Но информация откуда-то получается. Нужно сделать оговорку, потому что есть сайты, которые не используют сервер для получения какой-либо информации, а она уже сразу закодирована в верстке сайта. Такими сайтами называются «визиткой». Как понятно из названия, такие сайты нужны для того, чтобы красиво продемонстрировать род деятельности, прорекламировать компанию, которую представляет тот или иной сотрудник и указать там контакты для сотрудничества. Обычно такие сайты односторонние, под собой не имеют никакой логики. Такие сайты очень удобно парсить. Но чаще всего человек сталкивается с сложными сайтами со сложной логикой, где есть большая база пользователей и контента. Такие сайты просто обязаны обращаться за какой-то логикой на сервер.

Можно заметить, что в разделе 1.2.3 упоминается JS — логическая сторона сайта. Почему не пользоваться этой замечательной компонентой сайта для различной логики? Дело, конечно же, в безопасности. Куда более надежно будет положить какую-то информацию в переменную, а потом ее в будущем отобразить на странице при помощи, например, jQuery. Никто не хочет, чтобы алгоритмы обработки личной информации лежали перед всеми на видном месте, чтобы их можно было запустить и все расшифровать. Или как еще

по-другому можно хранить миллионы миллиардов информации?

Итак, так как предмет исследования как раз идет за нужной информацией на сервер, необходимо понять из каких крупниц информации создаются следующие более сложные запросы на сервер. Для этого необходимо изучить выдаваемую информацию на сайте, запросы, которые исполняются в процессе загрузки страницы. Исследование нужно производить с конца, идя в самое начало. Например, если мне нужно достать нужную мне картинку, то мне нужно найти запрос, который эту картинку запрашивает. Далее нужно найти запрос, который делает предыдущий запрос и так до самого основания.

Таким образом мы можем понять как сайт взаимодействует с сервером при различном получении разных страниц.

1.3.2 Изучение данных

Необходимо сделать несколько оговорок:

- все примеры будут продемонстрированы в браузере Google Chrome, так как в нём имеется богатый инструментарий для программистов, а также поддерживается огромное количество сайтов, ведь именно этот браузер является самым популярным. [5]
- комбинации клавиш, продемонстрированные в примерах, актуальны для операционной системы Windows

Для поиска информации необходимо пользоваться инструментом разработчика. Открывается данное меню разработчика при помощи нажатия клавиши F12. На рисунке 7 продемонстрировано меню.

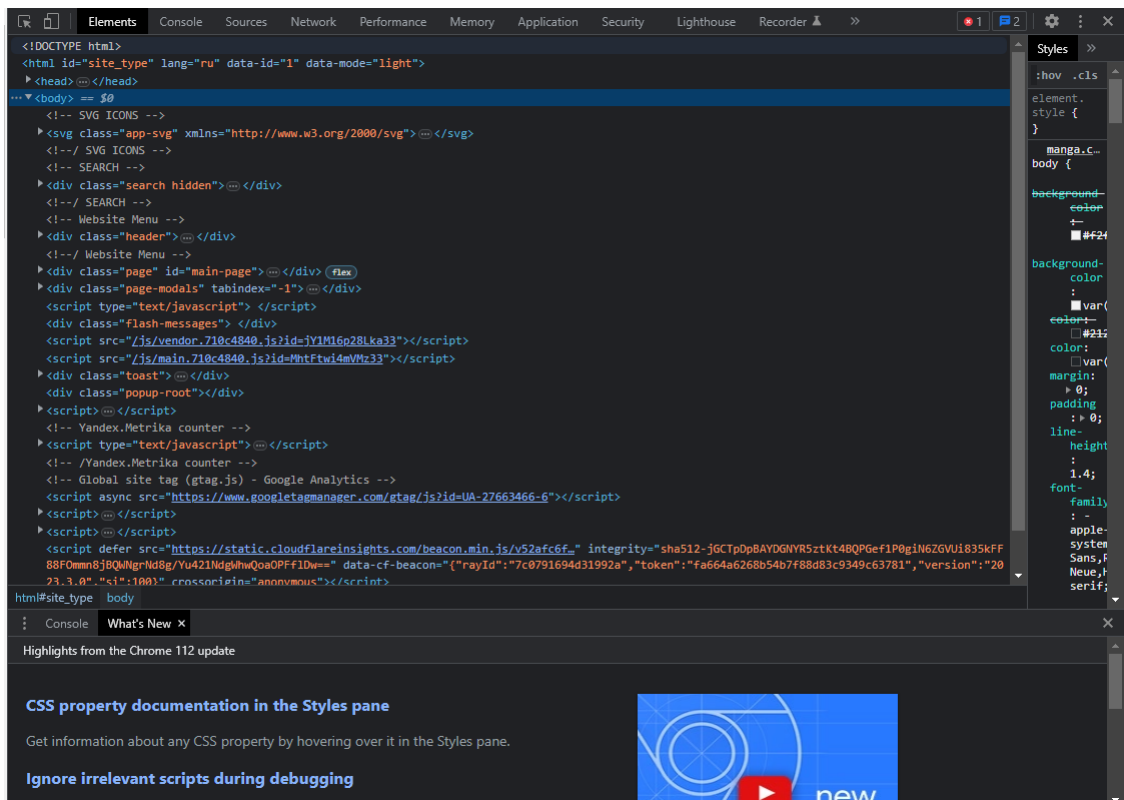


Рисунок 7 – Меню разработчика

У браузера есть свойство, при котором он исполняет все скрипты в тегах языка HTML. Как говорилось ранее, эти скрипты могут содержать запросы на сервер за получением информации на страницу. Также еще нужно знать, что вписывание ссылки в браузер и последующем нажатием подтверждения ввода, по сути своей делает GET запрос на домен с соответствующим роутингом. Зная также тот факт, что на один запрос приходится один ответ, мы можем определенно точно сказать, что программа наша получит сайт с первичной информацией. То есть будет загружены скрипты, которые делают последующие запросы на сервер. Из этого мы можем сделать вывод, что просто так никакая информация из ниоткуда не берется и с большой вероятностью мы при помощи только одного запроса не получим весь сайт. Для тестирования запросов можно пользоваться удобным инструментом postman [6]. На рисунке 8 продемонстрирован пример одного GET запроса на сайт через инструмент postman. На рисунке А.1 приложения изображена команда, которую можно сделать, чтобы посмотреть полностью код HTML.

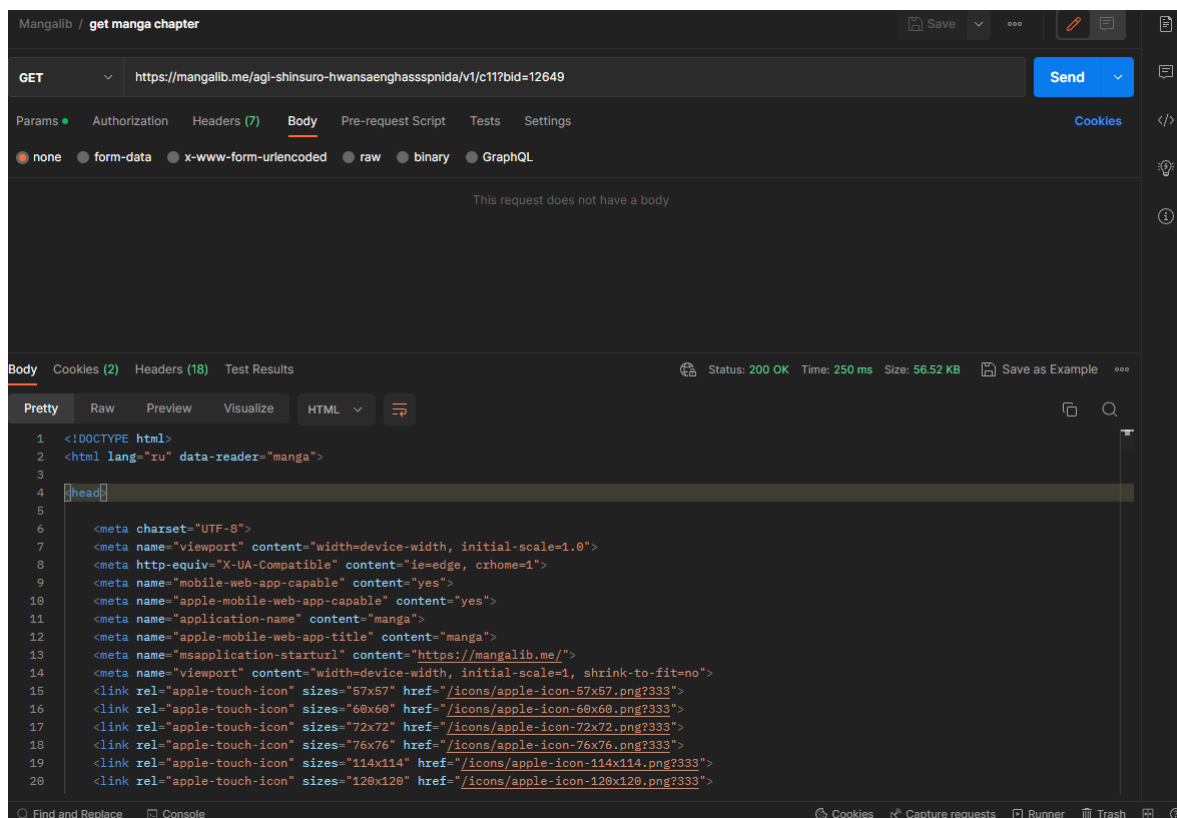


Рисунок 8 – GET запрос

Из кода можно увидеть, что HTML код очень сильно отличается от того, что изображен в меню разработчика. Это все потому, что мы не проделали все те запросы, что делает браузер, когда начинает загрузку страницы. К счастью, нам этого делать не надо. Это лишь доказательство того, что до любой информации мы можем одним или несколькими путями добраться.

Для того, чтобы начать изучение данных, необходимо зайти на сайт и найти самостоятельно ту конечную информацию, которая нужна. В данном случае ведется поиск картинок. После того, как мы нашли нужную картинку, необходимо понять как выглядит ссылка на нее. Исходя из той логики, что разработчик сайта будет составлять ссылку из абсолютно случайного набора букв и цифр – крайне мала, так как это решение бесполезно и очень сильно влияет на быстроту разботки в худшую сторону, мы можем попробовать догадаться по какой логике конструируется ссылка на ресурс. Итак, чтобы посмотреть на ссылку картинки, нам необходимо посмотреть код этой картинки в HTML. На рисунке 9 видно, как выделен тэг картинки, которую мы бы хотели получить.

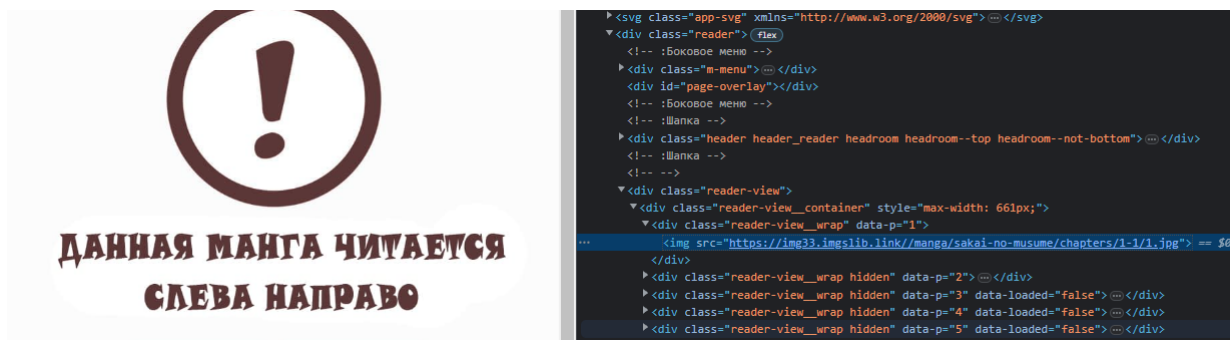


Рисунок 9 – Представление картинки в кода HTML

Посмотрев так на другие картинки, можно увидеть, что у них у всех есть один источник. А ссылка формируется по паттерну из рисунка А.2. Следовательно возможно эту ссылку конфигурировать своими руками в сервисе. Для этого нам нужны:

- название манги
- правильных id главы
- хеш картинки

Теперь необходимо узнать где на странице еще лежит вышеописанная информация. Для этого нам нужно проанализировать запросы. Для этого нам пригодится вкладка Network в меню разработчика. На странице картинки открываем меню разработчика и перезагружаем страницу. Это нужно для того, чтобы во вкладке появились все абсолютно запросы, которые делает браузер для того, чтобы прогрузить полностью страницу. На рисунке 10 продемонстрирован весь список запросов, который делается при загрузке страницы определенного сайта. Очевидно этот список может отличаться от того, что продемонстрировано на рисунке 10.

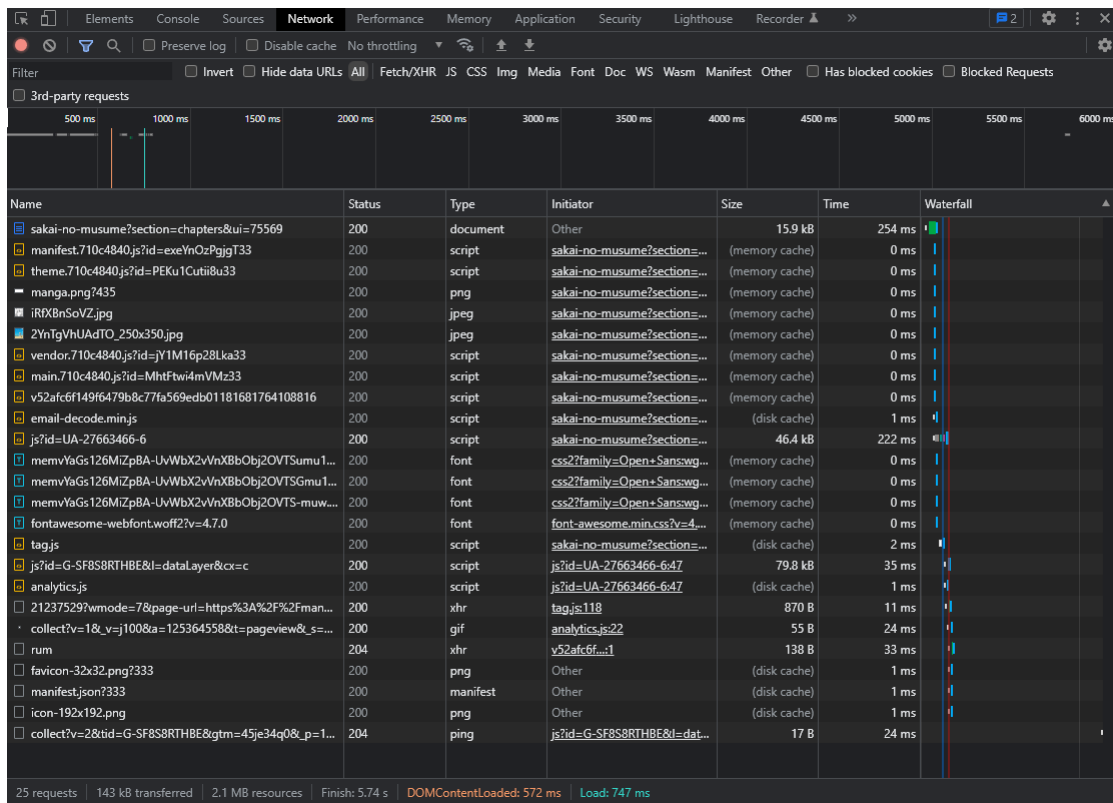


Рисунок 10 – Меню разработчика network вкладка

Далее необходимо найти эту картинку в запросах. В меню есть фильтр по типам файлов. Выбираем там картинки. На рисунке 11 мы видим, что не очень много запросов было на картинки. С легкостью можем найти наш файл (на рисунке выделен). Можно еще также заметить, что разработчики постарались и заранее прогружают следующую страницу, чтобы было проще смотреть контент.

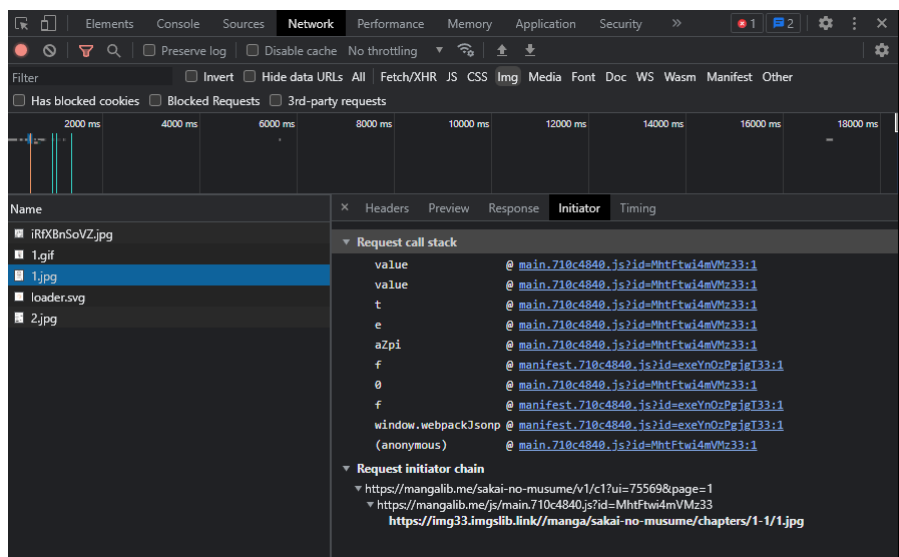


Рисунок 11 – Меню разработчика network вкладка с фильтром по картинкам

Теперь нам необходимо узнать как был получен запрос на эту картинку. Помним, что из-за особенности устройства браузера информация не появляется из ниоткуда. Следовательно мы можем с легкостью узнать откуда пришел запрос. Благо нам не нужно особо много, чтобы раздобыть подобного рода информацию. Меню разработчика предоставляет интерфейс для показа стека вызова. На рисунке 11 можно увидеть вкладку Initiator. Это и есть наш СТЭК ВЫЗОВОВ.

На самом деле, если мы будем смотреть код каждого запроса, сможем легко закопаться в нем, так как чаще всего подобного рода сайты подключают разные сервисы для генерации логики. Этот сайт не исключение. Здесь добавлена защита от DDOS атак от компании Яндекс, также трекер использования сайта. Через подобного рода скрипты происходит процесс анализа запросов. Поэтому в середине запроса может быть сгенерированный код, который почти невозможно понять без нескольких дней его анализа. Благо в нашем случае его анализировать не надо. Можно начать смотреть код, который находится вокруг места, откуда произошел вызов того или иного запроса. Таким образом натываемся на самый первый вызов в стеке и замечаем, что это наш любимый HTML код. Изучив его, можно увидеть, что есть большой скрипт. В нем есть переменная «window.__DATA__» (при помощи вызова скрипта на рисунке А.1 можно увидеть этот скрипт). Кажется это то, что нам нужно.

Переменная «window.__DATA__» — JSON, в котором явно написана

информация о манге. Исследовав этот объект, обнаруживаем, что там есть информация о названиях каждой страницы. Можно также выявить и сколько страниц в данной конкретной главе. На рисунке 12 можно посмотреть на получившуюся структуру, в которую можно парсить из JSON.

```
1  type pageList struct {  
2      Page int    `json:"p"`  
3      Url  string `json:"u"`  
4  }  
5
```

Рисунок 12 – Структуры содержания информации о главах и картинках

В Url атрибуте хранится тот самый хэш картинки. ID главы можно узнать из атрибутов на той же странице. Название манги совсем просто найти. Оно находится буквально в каждой ссылке данного сайта. Поэтому можно просто скопировать название оттуда.

Примеры выше были для сайта mangalib.me, участвующий в итоговом дипломе. Анализ каждого сайта происходит в индивидуальном порядке, так как везде свои разные алгоритмы и подходы. Задача микросервиса сделать так, чтобы добавление новых сайтов было достаточно простым и требовало исключительно анализа и исполнении всех алгоритмов для автоматического поиска.

1.3.3 Сайты с API

Есть некоторые сайты, разработчики которых подумали о пользователях и о том, что эти пользователи возможно хотели бы автоматизировать процесс поиска и скачивания контента. Одним из таких сайтов является mangadex [7].

Для подобных сайтов нет необходимости заниматься реверс-инженерингом, как было представлено выше. Достаточно изучить документацию [8] и найти хендлеры, которые выполняют поставленную задачу и использовать их. Однако, это не означает, что предварительную обработку не нужно делать. Благо вышесказанный сайт имеет много примеров использования и свой SWAGGER [9] для тестирования.

Тестирование подобного API происходит через принцип «проб и ошибок», так как не всегда случается такое, что появляется на выходе

ожидаемая информация. В некоторых случаях нужно заниматься ее обработкой.

1.4 База данных в микросервисах

Для некоторых микросервисов база данных совсем не нужна. Такие микросервисы обычно не используют никакую информацию, например, сервисы по конвертированию данных. Там запускается какой-то разовый алгоритм и все. Часто сохранять информацию в один csv файл. Однако сервисы, которые в основе своей работают с данными пользователями, просто необходим. Собранную информацию потом возможно использовать для анализа, чтобы, например, делать на ее основе модели искусственного интеллекта.

Одной из основных причин использования баз данных в микросервисах является возможность распределения данных и обработки запросов между несколькими сервисами. Каждый сервис может иметь свою собственную базу данных или использовать общую базу данных, что позволяет сократить нагрузку на каждый сервис и повысить производительность всей системы.

Кроме того, базы данных обеспечивают целостность данных и защиту от потери информации. Они позволяют хранить информацию в структурированном виде, что упрощает ее поиск и обработку. Базы данных также обеспечивают безопасность данных, позволяя контролировать доступ к ним и устанавливать права доступа для каждого сервиса.

Базы данных также позволяют решать проблемы масштабирования и управления ресурсами. Благодаря возможности горизонтального масштабирования баз данных, можно легко расширять хранение данных и увеличивать производительность системы при необходимости. Кроме того, базы данных позволяют эффективно управлять ресурсами, например, памятью и дисковым пространством, что повышает эффективность работы всей системы.

Использование баз данных в микросервисах также упрощает разработку и тестирование сервисов. Каждый сервис может использовать свою собственную базу данных для тестирования и разработки, что позволяет избежать конфликтов между сервисами и обеспечить более быстрое и эффективное тестирование.

Также, базы данных помогают поддерживать консистентность данных.

Очень важно эту особенность учитывать и в сервисе по поиску картинок. На ее основе можно сохранять ранее полученную информацию и не беспокоить сайт при каждом новом запросе. Эта особенность очень важна в тех случаях, если сайт имеет защиту против DDOS атак. Если сервисом будет пользоваться даже сотни людей, то у сайта может возникнуть вопрос о появлении стольких запросов с одного IP. Именно поэтому намного лучше и проще будет сохранять всю информацию в базе данных.

Приведу пример, есть популярная манга, которая так или иначе пользуется спросом. При помощи поиска этой же манги и конкретной главы возможно сразу из памяти БД достать информацию, лишняя раз не нагружая сайт запросами. Таким образом такое поведение очень хорошо скажется на:

- а) работоспособности приложения. Не придется каждый раз ждать ответа от конкретного сайта.
- б) внезапно не будет заблокирован наш микросервис за большое количество запросов на сайт
- в) не придется в микросервисе обходить защиту от роботов (captcha)

Перед тем, как использовать БД, нужно также подумать о том что будет храниться там. Очень важно это сделать на этапе конструирования. Иначе может получиться так, что эффекта от такого мощного инструмента не будет.

Также при помощи БД возможно сделать подписку на обновление контента. С ее помощью можно организовать рассылку, достав большой объем нужных данных. С файлом csv это было бы очень трудно сделать.

Также преимуществом работы с БД является транзакции. При помощи транзакции можно изменять чувствительные данные, не боясь при этом гонки (race condition).

1.5 Чистая архитектура

Clean Architecture - это популярный подход к разработке программного обеспечения, который был предложен Робертом Мартином (также известным как Uncle Bob). Этот подход помогает разработчикам создавать гибкие, расширяемые и легко тестируемые приложения.

В Golang Clean Architecture часто реализуется с помощью пакетов и интерфейсов. Пакеты используются для организации кода в логически связанные блоки, а интерфейсы - для определения контрактов между различными компонентами приложения. На рисунке 13 представлена схема

такого взаимодействия пакетов.

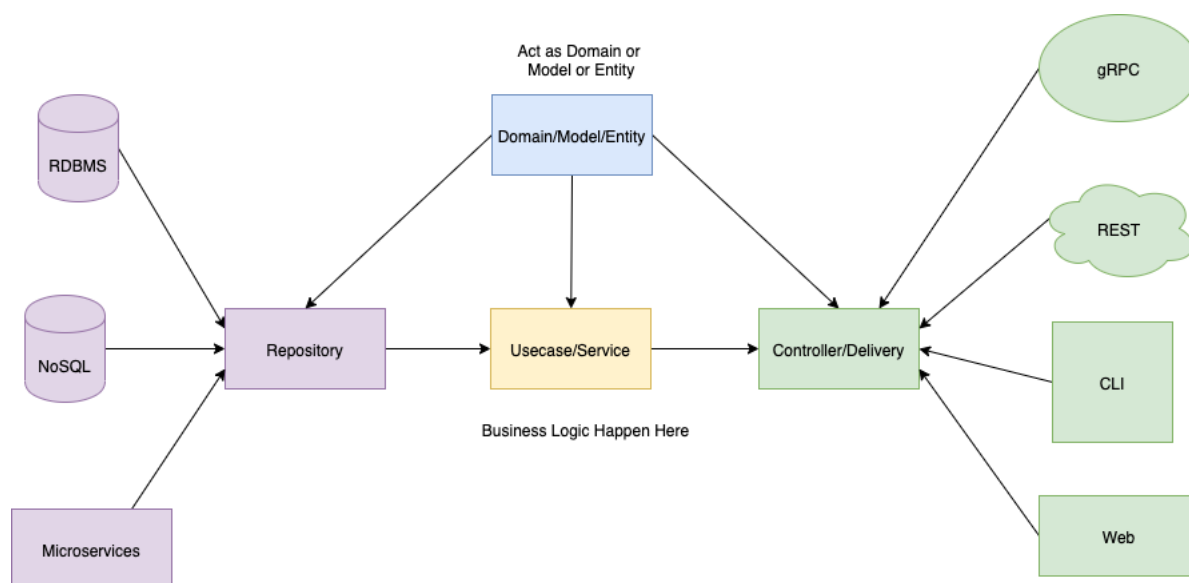


Рисунок 13 – Схема чистой архитектуры

В центре Clean Architecture находится принцип разделения зависимостей (Dependency Inversion Principle), который гласит, что зависимости должны быть направлены от абстракций к конкретным реализациям, а не наоборот. Это позволяет разработчикам легко заменять реализации компонентов приложения, не затрагивая другие части системы. Так, например, область, отвечающая за сохранение данных не должна содержать никакую логику. Она должна только выдавать или сохранять данные.

Стоит оговориться, что не стоит обязательно точно придерживаться этому концепту. Иногда в этом нет особой надобности. Но такое разделение лучше делать всегда и сохранять его в том или ином виде.

Одним из основных преимуществ Clean Architecture в Golang является возможность легко тестировать каждый слой независимо от других. Это достигается путем использования интерфейсов для определения контрактов между слоями. Это позволяет разработчикам легко заменять реализации компонентов приложения на макеты (mocks) для тестирования.

Кроме того, Golang Clean Architecture обеспечивает высокую гибкость и расширяемость приложения. Благодаря разделению зависимостей и использованию интерфейсов, программисты могут легко добавлять новые функции и изменять существующие, не затрагивая другие части системы.

В Golang Clean Architecture используется многослойная архитектура.

1.5.1 Представление (Presentation Layer)

Этот слой отвечает за взаимодействие с пользователем и обработку входящих запросов. Он содержит контроллеры, представления и другие компоненты, которые отображают данные для пользователя. В вышеизложенном сервисе расположены хендлеры, которые отвечают за обработку запросов по конкретному роутингу.

1.5.2 Использование (Use Case Layer)

Этот слой содержит бизнес-логику приложения. Он определяет, как приложение должно взаимодействовать с данными и какие операции должны быть выполнены. Он также содержит интерфейсы для доступа к данным. Например, в этом слое можно организовать логику парсинга картинок, организацию их по подгруппам и сбор в pdf. Затем данные отправлять в следующий слой.

1.5.3 Хранилище (Storage Layer)

Этот слой отвечает за хранение данных и доступ к ним. Он содержит реализации интерфейсов, определенных в слое использования. Эти реализации могут использовать различные источники данных, такие как базы данных, файлы или внешние API.

1.5.4 Инфраструктура (Infrastructure Layer)

Этот слой содержит код, который не относится к бизнес-логике приложения, но необходим для его работы. Это может быть код для обработки ошибок, логирования, аутентификации и т.д. Этот слой также содержит реализации интерфейсов, определенных в слое использования, для доступа к внешним сервисам и системам.

На картинке Б.3 демонстрируется пакетное представление проекта.

2 ПРАКТИЧЕСКАЯ ЧАСТЬ

2.1 Структура приложения

Приложение состоит из основной серверной части — API и клиентских частей. Клиентским сервисом считается любой сервис, например бот в мессенджере, способный отправлять запросы на серверную часть. Серверная часть выполнена в соответствии с концепцией чистой архитектуры, в соответствии с тем, что описано в разделе 1.5. На рисунке Б.3 изображена схема проекта. Каждый пакет отвечает за свою определенную часть.

Пакет «repository» отвечает за взаимодействие между БД. В «cmd» находится main.go файл, который запускает сервер. «App» содержит в себе реализацию хендлер функций. Пакет «config» используется по всему проекту. В нем описаны функции, помогающие доставать различные секреты приложения, такие как пароли, из переменных окружения. В этом же пакете происходит парсинг тех из yaml файла. В пакете «gen» содержится сгенерированные структуры библиотекой jet из раздела 1.1.3. В «internal/pgk» содержатся все пакеты, которые так или иначе выполняют какую-то логику приложения. Своего рода usecase слой. Пакет «parse» отвечает за логику парсинга описанного для конкретного сайта. Так можно увидеть там один из объектов исследования — сайт mangalib «Pdf_creator», как понятно из названия, занимается описанием логики сбора картинок в один pdf файл. Пакет «mailer» необходим для того, чтобы отправлять файлы через через какие-либо возможные способы передачи, например, отправки тех через электронную почту. В папке «migrations» описан SQL код всех миграций БД.

Внешние пакеты (те, что находятся не в папке «internal») используются для того, чтобы к ним возможно было получить доступ внешним клиентским сервисам. Так в «api» пакете описаны структуры request и response, отвечающие за боди запроса и ответа соответственно. Также в модуле описан роутинг API. Пакет «pgk» необходим для того, чтобы импортировать тот в клиентский сервис, так как там предполагается описание sdk. В этом пакете описаны кастомные ошибки, которые можно в последствии обработать.

ЗАКЛЮЧЕНИЕ

В результате выполненной НИР был сделан API способный предоставлять быстро и просто офлайн доступ к любой манге для дальнейшего прочтения в удобных мессенджерах или электронной книге. Оценка полноты решений показала, что данное направление позволяет настроить формат чтения под свои потребности и сохранять любимые серии для повторного чтения без необходимости загрузки каждый раз заново. Также была организована база данных для сбора информации о пользователях, скаченной манги и отслеживание ее обновления.

Рекомендуется использовать результаты НИР для продвижения электронных книг в качестве основного способа чтения манги, а также для разработки новых функций и возможностей, которые могут улучшить процесс чтения.

Оценка технико-экономической эффективности внедрения показала, что использование электронных книг для чтения манги может сэкономить время и деньги, которые ранее были затрачены на поиск и загрузку серий в интернете.

Сравнение выполненной НИР с лучшими достижениями в этой области показало, что использование электронных книг для чтения манги является инновационным и перспективным направлением, которое может улучшить процесс чтения и повысить удовлетворенность пользователей.

Репозиторий с результатами работы [**result-cite**].

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. *GitHub*. Pdfcpu. — 2023. — URL: <https://github.com/pdfcpu/pdfcpu> (дата обращения 28.04.2023).
2. *GitHub*. Goquery. — 2023. — URL: <https://github.com/PuerkitoBio/goquery> (дата обращения 28.04.2023).
3. *JQuery*. documentations. — 2023. — URL: <https://jquery.com/> (дата обращения 28.04.2023).
4. *GitHub*. Go-jet. — 2023. — URL: <https://github.com/go-jet/jet> (дата обращения 29.04.2023).
5. *Similarweb*. Доля рынка топ браузеров. — 2023. — URL: <https://www.similarweb.com/ru/browsers/> (дата обращения 01.05.2023).
6. *Postman*. Getting started. — 2023. — URL: <https://www.postman.com/> (дата обращения 01.05.2023).
7. *Mangadex*. Главная страница. — 2023. — URL: <https://mangadex.org/> (дата обращения 01.05.2023).
8. *Mangadex*. API documentation. — 2023. — URL: <https://api.mangadex.org/docs/> (дата обращения 01.05.2023).
9. *Youtube*. Что такое Swagger и OpenAPI за 3 минуты. — 2023. — URL: <https://youtu.be/aaFDBgPdXw4> (дата обращения 01.05.2023).

ПРИЛОЖЕНИЕ А

Приложение к изучению данных

На рисунке А.1 демонстрируется bash скрипт, который можно исполнить для получения HTML кода анализируемой страницы сайта.

```
1 #!/bin/bash
2 curl --location \
3     'https://mangalib.me/agi-shinsuro-hwansaenghassspnida/v1/c11'
```

Рисунок А.1 – Пример HTML кода получения страницы сайта

На рисунке А.2 демонстрируется шаблон, по которому формируется ссылка на картинку сайта mangalib.me

```
1     https://img33.imgslib.link//manga/{manga-name}/chapters/{
2     chapter-hash}/{image-hash}
```

Рисунок А.2 – Шаблон конфигурируемой ссылки для конкретной картинки

ПРИЛОЖЕНИЕ Б

Приложение к чистой архитектуре

На рисунке Б.3 демонстрируется файловая схема проекта.

```
1 +---api
2 |   +---domain
3 |   +---middleware
4 |   \---router
5 +---cmd
6 +---docs
7 +---internal
8 |   +---app
9 |   +---config
10 |   |   \---flags
11 |   +---gen
12 |   |   \---manga_parser
13 |   |       \---public
14 |   |           +---enum
15 |   |           +---model
16 |   |           \---table
17 |   +---pkg
18 |   |   +---err_controller
19 |   |   +---mailer
20 |   |   +---parse
21 |   |   |   +---http_requester
22 |   |   |   \---mangalib
23 |   |   |       \---mock
24 |   |   +---pdf_creator
25 |   |   |   +---mock
26 |   |   |   \---pdf_merger
27 |   |   \---utils
28 |   \---repository
29 +---migrations
30 \---pkg
31     \---custom_errors
```

Рисунок Б.3 – Схема пакетов проекта