

Міністерство освіти і науки України
Національний університет «Києво-Могилянська академія»
Факультет інформатики
Кафедра мультимедійних систем

Кваліфікаційна робота
освітній ступінь – бакалавр

на тему: **«Розробка системи для збору та обробки даних з використанням
Apache Spark»**

Виконав: студент 4-го року навчання,
спеціальності

121 «Інженерія Програмного Забезпечення»

Пінкевич В.М.

Керівник

Старший викладач

Борозенний С.О.

_ травня 2023 р.

Київ – 2023

Календарний план виконання роботи

Тема: Розробка системи для збору та обробки даних з використанням Apache Spark

Календарний план виконання роботи:

№ п/п	Назва етапу кваліфікаційної роботи	Термін виконання етапу	Примітка
1.	Отримання теми роботи	Жовтень 2022 р.	
2.	Отримання завдання на кваліфікаційну роботу	Листопад 2022 р.	
3.	Ознайомлення із літературою та джерелами за темою роботи	Листопад - грудень 2022 р.	
4.	Знайомство із засобами розробки системи	Грудень 2022 р. - січень 2023 р.	
5.	Розробка застосунку для збору даних	Січень – лютий 2023 р.	
6.	Розробка застосунку для обробки даних	Березень - квітень 2023 р.	
7.	Написання текстової частини роботи	Квітень - травень 2023 р.	
8.	Перегляд кваліфікаційної роботи науковим керівником	Квітень - травень 2023 р.	
9.	Внесення змін до роботи згідно з зауваженнями наукового керівника	Квітень - травень 2023 р.	
10.	Захист кваліфікаційної роботи	Червень 2023 р.	

Студент Пінкевич В.М.

Керівник Борозенний С.О.

Зміст

Зміст	1
Анотація.....	3
Використані скорочення	4
Вступ.....	6
Розділ 1. Огляд джерел отримання даних.....	8
1.1 Steam	8
1.2 Steam Spy	11
Розділ 2. Створення архітектури системи.....	13
2.1 Основні компоненти системи	13
2.2 Формат зберігання даних.....	15
2.3 Загальні принципи роботи системи.....	18
Розділ 3. Розробка застосунку для збору даних	20
3.1 Вибір технологій розробки	20
3.2 Steam API.....	21
3.3 Steam Spy API.....	23
3.4 Розробка застосунку.....	24
Розділ 4. Розробка застосунку для обробки даних	32
4.1 Вибір технологій розробки	32
4.2 Структура бази даних.....	34
4.3 Розробка застосунку.....	37
Розділ 5. Аналіз виконаної роботи	47

5.1	Можливості розвитку розробленої системи	47
5.2	Висновки	47
	Список використаних джерел	49
	Додатки	51
	Додаток А. Клас WebClientConfiguration	51
	Додаток Б. Клас SteamService	52
	Додаток В. Клас SteamAppDetailsResponseDeserializer	53
	Додаток Г. Клас SteamDataCsvFileWriter	54
	Додаток Д. Клас SteamAppDetailsProcessor	57
	Додаток Е. Файлова структура проекту застосунку для збору даних	58
	Додаток Є. Файлова структура проекту застосунку для обробки даних	60

Анотація

У роботі описується процес розробки системи для отримання та обробки даних. Розглядаються можливі варіанти реалізації такої системи, їх переваги та недоліки, обрані інструменти розробки, їх особливості та основні можливості. Описуються деталі реалізації системи, можливості розширення та покращення розробленого рішення. Також проводиться аналіз основних проблем, які виникають при розробці таких систем, та основні способи їх вирішення.

Використані скорочення

API (Application Programming Interface) – засоби, які надає система і які можуть бути використані іншими системами для взаємодії із даною системою.

JSON (JavaScript Object Notation) – текстовий формат зберігання та передачі даних, який базується на синтаксисі та форматі об'єктів мови JavaScript.

CSV (Comma Separated Values) – текстовий формат зберігання та передачі даних, який зберігає інформацію у вигляді таблиць що складаються зі стовпців та рядків.

XML (eXtensible Markup Language) – текстовий формат зберігання та передачі даних, який базується на системі тегів.

СКБД (система керування базами даних) – програма або набір програм, які забезпечують контроль над доступом до бази даних і надають можливості для пошуку, додавання, оновлення та видалення даних із бази даних.

SQL (Structured Query Language) – декларативна мова програмування, яка надає можливості для керування базами даних а також інформацією, що зберігається у базах даних.

HTTP (HyperText Transfer Protocol) – мережевий протокол прикладного рівня, який забезпечує комунікацію та передачу даних між пристроями у мережі.

URL (Uniform Resource Locator) – стандартизована адреса унікального ресурсу в мережі.

JVM (Java Virtual Machine) – абстрактна машина, специфікація, яка описує, як має відбуватися виконання програм, написаних мовою Java, або програм, написаних іншими мовами програмування та скомпільованих у байт-код Java. Конкретні реалізації JVM забезпечують виконання байт-коду на різних платформах та кінцевих пристроях.

JDBC (Java Database Connectivity) – специфікація та інтерфейс (API), який надає можливість програмам, написаним мовою Java, взаємодіяти з різними СКБД.

UDF (User Defined Function) – функція, що створюється розробником і може бути зареєстрована та використана у фреймворку Apache Spark так, ніби вона є його частиною.

Вступ

Основою роботи будь-якої системи є дані. Будь-яка система у тому чи іншому вигляді працює із даними. Вона може їх створювати, обробляти, надавати іншим системам, тощо. Часто для повноцінної роботи системі потрібні не лише ті дані, які вона створює, а й сторонні, зібрані іншими системами. Ці дані можуть бути як публічно доступними і не вимагати жодних умов для використання, так і закритими, коли отримання доступу до них має свої особливості. Не менш важливою частиною багатьох систем є також обробка зібраних даних. Часто об'єм цих даних є настільки великим, що для виконання повної їх обробки необхідна велика кількість ресурсів.

Для вивчення особливостей роботи таких систем (а також кожного з етапів роботи з даними) було вирішено розробити власну систему, яка буде отримувати дані із відкритих джерел, проводити їх обробку та зберігати вже оброблені дані.

Об'єктом дослідження є вивчення принципів розробки та функціонування систем збору, обробки та зберігання великої кількості даних.

Предметом дослідження є розробка системи, яка буде взаємодіяти зі сторонніми сервісами через публічні веб API для отримання даних, виконувати очищення, обробку даних та зберігати оброблені дані для їх подальшого використання.

Метою роботи є створення системи, яка буде отримувати дані від сторонніх сервісів Steam та Steam Spy через публічні веб API, проводити їх очистку, трансформацію та зберігати оброблені дані у реляційній базі даних. Через взаємодію зі сторонніми публічними веб API буде досліджено особливості взаємодії із такими API. Використання фреймворку Apache Spark дозволить розглянути основні підходи до роботи із великою кількістю даних та дослідити можливу реалізацію цих підходів на прикладі обраного фреймворку.

Текстова частина роботи складається із п'яти розділів.

У першому розділі виконується огляд сервісів Steam та Steam Spy, API яких будуть використані для отримання даних.

Другий розділ фокусується на розробці архітектури системи. Проводиться огляд декількох можливих варіантів реалізації системи, їх переваг та недоліків, а також вибір фінального рішення.

Третій та четвертий розділи описують безпосередньо розробку обох частин системи: для збору та обробки даних відповідно. Для кожної частини обираються інструменти розробки, описуються їх переваги та недоліки, розглядаються деталі реалізації та проблеми, що виникали в ході розробки, і шляхи їх вирішення. Третій розділ також включає огляд особливостей при роботі зі сторонніми публічними веб API, а четвертий – особливості роботи з великими даними.

У п'ятому розділі підсумовуються результати розробки системи та аналізуються можливості покращення та подальшого розширення розробленої системи.

Розділ 1. Огляд джерел отримання даних

1.1 Steam

Steam, за визначенням своїх розробників – це “найбільше і найкраще місце для гравання, обговорення та створення відеоігор”. [1]

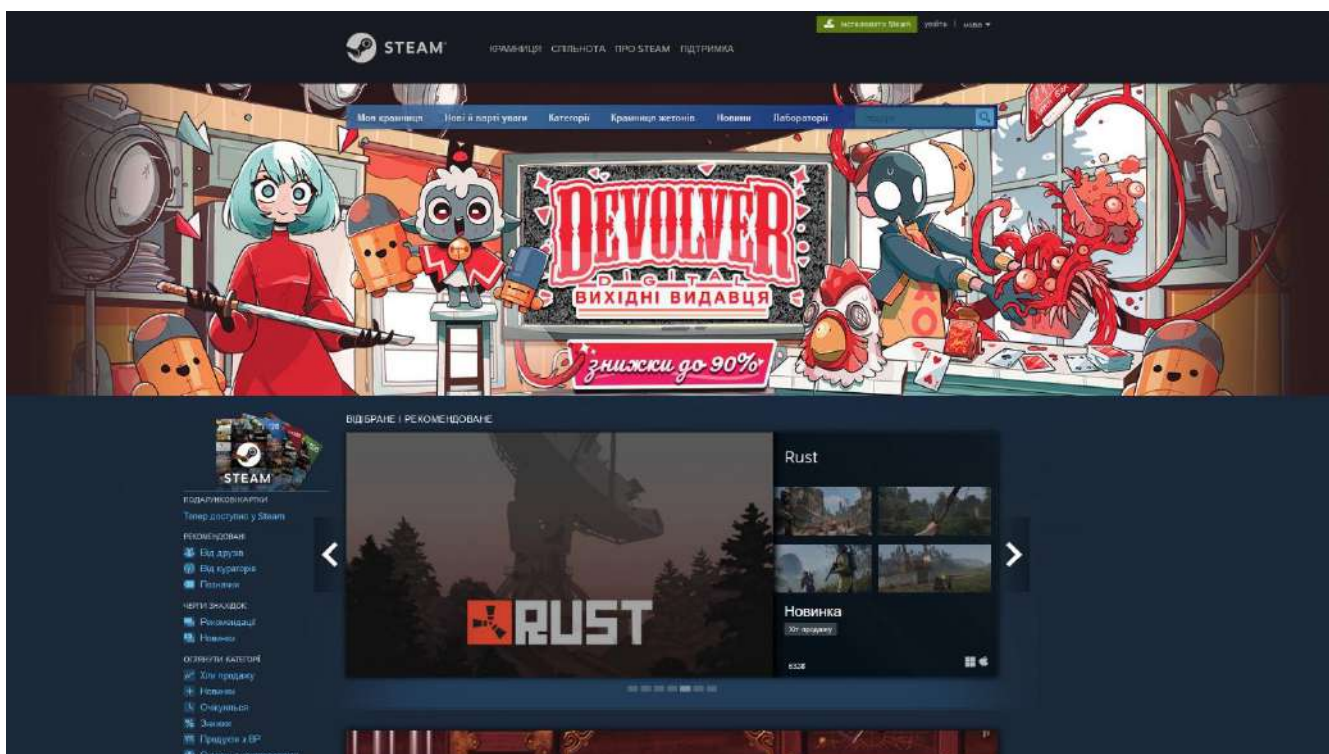


Рисунок 1.1.1. Головна сторінка вебсайту сервісу Steam

Якщо спочатку це був онлайн-сервіс для продажу відеоігор, то тепер це ціла окрема платформа, яка забезпечує велику кількість функцій як споживачам, так і розробникам відеоігор. Основні можливості, які Steam надає користувачам, включають:

- можливість покупки (чи безкоштовного отримання) відеоігор, доповнень до них, саундтреків, програмного забезпечення та їх зберігання у власній бібліотеці користувача;
- наявність детальної інформації про акаунт користувача, статистики користування сервісом (час, проведений у кожній грі, отримані досягнення, нещодавня активність), можливості його персоналізації та

розвитку (наявність системи карток, значків, способів оформлення власного профілю);

- розвинена система взаємодії між користувачами платформи. Вони можуть писати рецензії на ігри, рекомендувати чи не рекомендувати їх іншим користувачам, створювати власні обговорення на окремих форумах, створювати фото та відео за мотивами ігор, модифікації для них, писати інструкції та підказки, тощо. Тобто існують досить широкі можливості не лише для гравця, а й для створення користувацького контенту, який стосується ігор, проте не є їх частиною, задуманою розробниками;
- наявність власної соціальної системи. Користувачі можуть додавати один одного до списку друзів, слідкувати за оновленнями та прогресом в іграх, грати в ігри разом, тощо;
- наявність внутрішнього ринку. Багато ігор надають можливість отримання, покупки та продажу внутрішньоігрових предметів. А Steam має власний ринок, який дозволяє користувачам виставляти предмети на продаж за бажаною ціною, купувати їх чи обмінюватись ними.

Розробникам Steam також надає велику кількість можливостей, серед них:

- повний контроль над власними продуктами у Steam. Розробник може виставляти власні ігри на продаж, влаштовувати акції та розпродажі, піднімати чи зменшувати ціну на гру, керувати додатковими матеріалами (фото, відео, опис, системні вимоги, тощо);
- наявність аналітичних даних щодо власних продуктів. Розробник має можливість спостерігати за кількістю проданих копій гри від часу її додавання, інформацію про кінцевих користувачів, які купили гру, поточну кількість гравців, які грають у гру;

- можливості для взаємодії із гравцями. Розробник може публікувати новини, що стуються гри, проводити прямі трансляції, створювати власний сторонній контент для гри, бачити рецензії та контент, розроблений гравцями. Це дозволяє розробникам контактувати зі своєю аудиторією, отримувати зворотній зв'язок та покращувати якість власних продуктів.

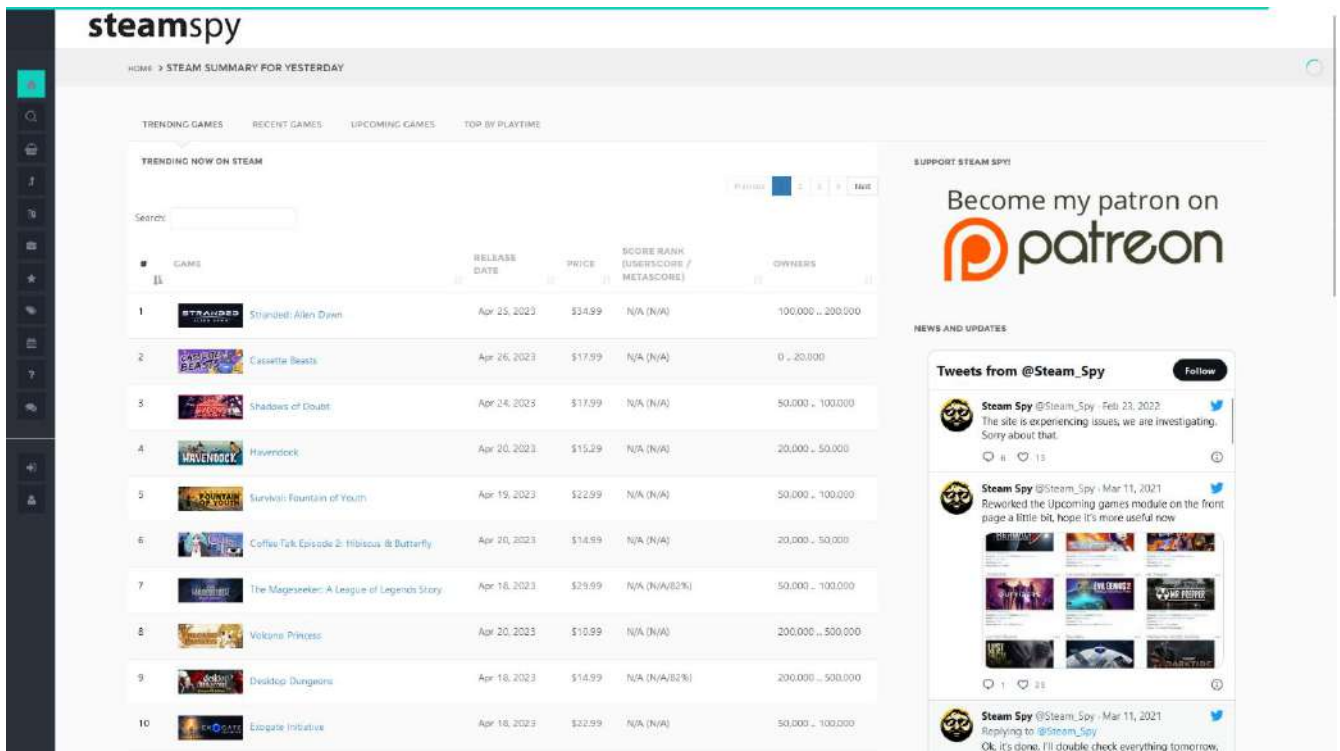
Саме завдяки наявності великої кількості можливостей не лише для покупки ігор, а й для їх обговорення, створення стороннього контенту та взаємодії кінцевих користувачів як між собою, так і напряму із розробниками, Steam є найпопулярнішим сервісом для покупки цифрових копій ігор на сьогоднішній день.

Steam має декілька різних клієнтів. На даний момент існує вебсайт, доступний через браузер, застосунки для операційних систем Windows, macOS та Linux, додатки для мобільних операційних систем Android та iOS, а також власна операційна система SteamOS (яка базується на операційній системі Debian сімейства Linux), яка уже при встановленні включає клієнт Steam та широко інтегрована з ним. Усі основні функції Steam доступні на будь-якому його клієнті, проте кожен має свої особливі функції. Наприклад, безпосередньо завантажити та грати в ігри зі своєї бібліотеки можна лише через настільні застосунки, а двофакторна аутентифікація доступна лише при використанні мобільних додатків.

Steam було обрано для отримання інформації про ігри саме через те, що він є найбільшим сервісом для продажу ігор, і, відповідно має інформацію про найбільшу кількість ігор, а також через наявність публічного API, який дозволяє отримувати дані без необхідності додаткових дій з боку розробника.

1.2 Steam Spy

Steam Spy – це статистичний сервіс для Steam, який базується на API, який надає компанія Valve (компанія, що розробляє Steam). Він автоматично збирає інформацію із користувацьких профілів Steam, аналізує її та дозволяє переглядати її у простому та зручному форматі. [2]



The screenshot shows the Steam Spy website interface. The main content is a table titled "TRENDING NOW ON STEAM" with columns for Rank, Game, Release Date, Price, Score Rank (UserScore / Metascore), and Owners. The table lists 10 games, with "Stranded: Alien Dawn" at the top.

RANK	GAME	RELEASE DATE	PRICE	SCORE RANK (USERSCORE / METASCORE)	OWNERS
1	STRANDED ALIEN DAWN	Apr 25, 2023	\$34.99	N/A (N/A)	100,000 - 200,000
2	CASSETTE BEASTS	Apr 26, 2023	\$17.99	N/A (N/A)	0 - 20,000
3	SHADOWS OF DOUBT	Apr 24, 2023	\$17.99	N/A (N/A)	50,000 - 100,000
4	HAVERLOCK	Apr 20, 2023	\$15.29	N/A (N/A)	20,000 - 50,000
5	STARVALDI FOUNTAIN OF YOUTH	Apr 19, 2023	\$22.99	N/A (N/A)	50,000 - 100,000
6	COFFEE TALK EPISODE 2: FIBRIOUS & BUTTERFLY	Apr 20, 2023	\$14.99	N/A (N/A)	20,000 - 50,000
7	THE MAGESSEEN: A LEAGUE OF LEGENDS STORY	Apr 18, 2023	\$29.99	N/A (N/A/82%)	50,000 - 100,000
8	YOKUNA PRINCESS	Apr 20, 2023	\$10.99	N/A (N/A)	200,000 - 500,000
9	DESKTOP DUNGEONS	Apr 18, 2023	\$14.99	N/A (N/A/82%)	200,000 - 500,000
10	ELIZABETH HIBISCUS	Apr 18, 2023	\$22.99	N/A (N/A)	50,000 - 100,000

Additional elements visible in the screenshot include a search bar, navigation tabs (TRENDING GAMES, RECENT GAMES, UPCOMING GAMES, TOP BY PLAYTIME), a sidebar with navigation icons, a "SUPPORT STEAM SPY" section with a Patreon link, and a "Tweets from @Steam_Spy" section.

Рисунок 1.2.1. Головна сторінка вебсайту Steam Spy

Основною метою сервісу є отримання доступної інформації через API сервісу Steam та її аналіз, що дозволяє кінцевим користувачам сервісу отримувати найактуальнішу інформацію про наявні у Steam ігри. Крім перегляду звичайної інформації про ігри, яку з легкістю можна отримати у самому Steam, Steam Spy також надає дані про приблизну кількість проданих копій гри, кількість позитивних та негативних оцінок від користувачів, середній проведений час у грі, тощо. Також Steam Spy агрегує дані щодо розподілу ігор за різними критеріями: розробниками, видавцями, доступними мовами, жанрами, тегами, поточними знижками та розпродажами. Більше того, для кожного критерію сервіс надає різноманітну специфічну статистику. Наприклад, для кожного видавця доступна

інформація про загальну кількість його ігор у Steam, загальну кількість проданих копій усіх його ігор, середню ціну гри, середній проведений час у грі.

Як зазначають розробники Steam Spy, дані сервісу, особливо щодо кількості проданих копій ігор, не є точними і надають лише приблизну оцінку. Також у даних можуть виникати деякі аномальні значення, пов'язані із принципами роботи Steam. Наприклад, Steam деколи проводить “безкоштовні вихідні” – період часу, протягом якого деякі ігри можуть бути отримані кінцевими користувачами безкоштовно на цей період часу. Протягом цього періоду Steam позначає кожного кінцевого користувача, який хоча б один раз авторизувався в сервісі, власником цих безкоштовних ігор, а після закінчення цього періоду – повертається до попереднього стану. Дані щодо кількості власників гри Steam надає у публічному API, тому при отриманні даних щодо кількості власників гри у періоди “безкоштовних вихідних” значення може значно відрізнятись від того, яке було отримано поза цими періодами. [2]

Steam Spy було обрано додатковим джерелом інформації про ігри у Steam, оскільки цей сервіс дозволяє отримати дані, які важко знайти використовуючи лише API Steam, або дані, яких там взагалі немає і для отримання яких потрібно виконувати власний аналіз. Важливим фактором також є доступний публічний API, який дозволяє отримувати усю необхідну інформацію.

Розділ 2. Створення архітектури системи

2.1 Основні компоненти системи

Основні задачі, які має виконувати система – це отримання даних від описаних сервісів через публічні API, подальша їх обробка, очищення, трансформація та збереження. На перший погляд, для розробки такої системи достатньо розробити один застосунок, який буде виконувати усі необхідні задачі. Проте, при використанні такого підходу з'являється декілька проблем:

- сильна зв'язаність частин системи, які відповідають за отримання та обробку даних. Частина системи, яка відповідає за отримання даних, не має залежати від тієї частини, яка відповідає за їх обробку і навпаки. При внесенні змін в одну із частин, інша, з великою ймовірністю, також потребуватиме відповідних змін;
- декілька точок відмови. Наприклад, при отриманні даних зі сторонніх публічних API немає жодної впевненості в тому, чи будуть вони працювати в певний момент часу. Також якщо будуть виникати помилки при обробці чи збереженні даних, наслідком яких стане завершення роботи застосунку, частина системи, яка відповідає за отримання даних, також завершить свою роботу, хоча в її роботі не було жодних помилок;
- неефективний розподіл ресурсів. Якщо проаналізувати як має працювати система, то зрозуміло, що при обробці та збереженні даних найбільш потрібними ресурсами є процесор та оперативна пам'ять. Натомість, отримання даних від сторонніх API не потребує великої кількості ресурсів процесора чи оперативної пам'яті, оскільки жодних складних операцій не виконується, адже найбільша кількість часу витрачається на процеси введення-виведення, а саме – на мережеву комунікацію та очікування відповіді від стороннього сервера. Таким чином, при

реалізації усієї системи у вигляді єдиного застосунку досить важко буде реалізувати ефективний розподіл ресурсів між її частинами.

Проаналізувавши принципи роботи системи та описані проблеми, цілком логічним здається розподіл системи на дві частини: частину для отримання інформації від сторонніх API та частину для обробки отриманих даних та їх збереження, причому кожен частину можна реалізувати у вигляді окремого застосунку. Такий розподіл має декілька переваг, серед них:

- незалежні системи отримання та обробки даних. Їх розділення дозволяє вносити зміни до кожної системи окремо, не впливаючи на іншу;
- свобода у виборі технологій. Для кожного застосунку можна обрати набір технологій, який найкраще підійде для вирішення саме його задач;
- вища відмовостійкість. Через незалежність обох застосунків проблеми в роботі одного з них не впливають на роботу іншого, а при завершенні роботи одного з них інший продовжує працювати;
- можливості налаштування та керування ресурсами. Для кожного застосунку можна налаштувати необхідну йому кількість ресурсів, забезпечуючи ефективно їх використання;
- зручність розробки та тестування. Кожен застосунок розгортається окремо, а тому для перевірки правильності роботи одного з них не потрібно розгортати інший, оскільки він не потрібний для тестування.

При розділенні системи на два окремих застосунки виникає потреба у забезпеченні процесу доставки даних від застосунку, який їх отримує, до застосунку, який займається їх обробкою. Для цього можна розглянути декілька варіантів:

- передача даних у вигляді зберігання проміжних результатів у файлах. Найбільш простий варіант, оскільки не вимагає використання жодних сторонніх систем, проте досить гнучкий, оскільки формат для зберігання можна обрати будь-який: JSON, CSV, XML, тощо;
- використання реляційної бази даних для зберігання проміжних результатів. Не дуже вдалий варіант, оскільки отримані дані можуть мати різну структуру, а реляційні бази даних не є найкращим вибором для роботи з такими даними;
- використання брокерів повідомлень для передачі проміжних даних. Такий варіант буде найкращим вибором при розробці системи, яка має отримувати, обробляти та зберігати дані в реальному часі, оскільки дозволяє обробити та зберегти нові дані до системи одразу після їх отримання через брокер повідомлень.

Оскільки система не повинна працювати із даними у реальному часі, то використання брокера повідомлень буде надлишковим, а використання бази даних для зберігання проміжних даних, які після обробки вже не будуть потрібними, є не дуже ефективним та потребує розгортання додаткового застосунку - СКБД. Таким чином, найкращим вибором буде зберігання усіх проміжних даних у файлах.

2.2 Формат зберігання даних

Для зберігання уже оброблених та трансформованих даних необхідно використати певне сховище. Оскільки кінцеві дані можуть бути використані в подальшому для виконання певного аналізу, підрахунку статистики, то найкращим рішенням буде зберігати дані у реляційній базі даних. Основними перевагами такого підходу є:

- зберігання не лише даних, а і зв'язків між ними. Основною особливістю реляційних баз даних є саме можливість будувати зв'язки між даними для подальшого їх використання;
- можливості контролю за вхідними даними. Більшість СКБД підтримують функціонал для створення певних перевірок вхідних даних на основі умов, заданих розробником, що дозволяє уникнути зберігання неправильних даних;
- можливості для редагування та видалення існуючих даних. При збереженні даних у файлах із цим можуть бути проблеми, оскільки зазвичай розробнику потрібно буде створювати власне рішення для правильного оновлення чи видалення даних із файлу. А СКБД надають ці можливості одразу, без необхідності додаткових дій з боку розробника;
- підтримка мови SQL. Керування базою даних через СКБД відбувається за допомогою мови SQL, яка є потужним інструментом для роботи із даними. Крім додавання, редагування та видалення даних, вона має широкі можливості для пошуку необхідних даних та їх аналізу.

Робота із базою даних зазвичай відбувається через СКБД. Саме СКБД є точкою входу в базу даних та відповідає за усі операції, які клієнт хоче виконати. СКБД відрізняються між собою стандартами мови SQL, які вони підтримують, наявністю специфічних можливостей та розширень SQL, таких як власні типи даних чи оператори, продуктивністю, тощо. Для зберігання даних у системі було обрано СКБД MySQL. MySQL – це найпопулярніша реляційна база даних із відкритим вихідним кодом, що розробляється, поширюється та підтримується компанією Oracle. Основними перевагами MySQL є: [3]

- повна підтримка багатьох стандартів мови SQL;
- наявність великої кількості типів даних як фіксованого, так і змінного розміру;

- висока продуктивність та ефективність завдяки можливості використання багатоядерних процесорів або декількох процесорів одночасно (якщо такі доступні на машині, на яку встановлено СКБД);
- підтримка найбільших операційних систем, таких як Windows, macOS та Linux;
- можливість використовувати двигуни як з підтримкою транзакцій, так і без неї;
- підтримка баз даних і таблиць із дуже великою кількістю даних;
- наявність драйверів та клієнтів для багатьох мов програмування;
- широкі можливості для налаштування та легкості розгортання.

Враховуючи усі вищезазначені переваги, використання СКБД MySQL для зберігання даних у системі є цілком доречним рішенням.

Для розгортання СКБД буде використано Docker. Docker – це інструмент для контейнеризації будь-яких застосунків. Він дозволяє розділити програми та інфраструктуру, яка потрібна для їх виконання, забезпечуючи ізоляцію програми від інфраструктури через використання контейнерів. [4] Використання Docker дозволяє уникнути необхідності розгортання СКБД та усіх її залежностей на локальній машині. Щоб уникнути необхідності встановлювати одні й ті ж налаштування для образу СКБД при кожному її запуску буде використано інструмент Docker Compose, який є частиною Docker. Встановлені налаштування образу СКБД через Docker Compose можна побачити на рисунку 2.2.1.

```

docker-compose.yml X
C: > Users > Victor > IdeaProjects > diploma-work > docker-compose.yml
1  version: '3'
2
3  services:
4
5    mysql:
6      container_name: mysql-diploma-work
7      image: mysql:8.0.32
8      restart: always
9      ports:
10     - "3306:3306"
11     environment:
12       MYSQL_DATABASE: diploma-work
13       MYSQL_ROOT_PASSWORD: password
14     volumes:
15       - mysql-data:/var/lib/mysql
16
17 volumes:
18   mysql-data:
19     driver: local
20     driver_opts:
21       type: "none"
22       o: "bind"
23       device: "D:/University/Subjects/Рік 4/Дипломна робота/mysql/mysql-data"
24

```

Рисунок 2.2.1. Налаштування образу СКБД MySQL через Docker Compose

2.3 Загальні принципи роботи системи

Загальний процес роботи системи можна описати покроково:

1. Отримання даних від сторонніх API: Steam API та Steam Spy API.
2. Запис отриманих даних в окремі файли для їх подальшої обробки.
3. Зчитування даних із файлів для їх подальшої обробки.
4. Очищення та трансформація даних.
5. Збереження перетворених даних у базу даних.

Оскільки система складається із двох застосунків: для отримання даних та для їх обробки, то, відповідно, перший застосунок буде виконувати кроки 1 та 2, а

другий – кроки 3, 4 та 5. Структуру системи та взаємодію її частин можна побачити на рисунку 2.3.1:

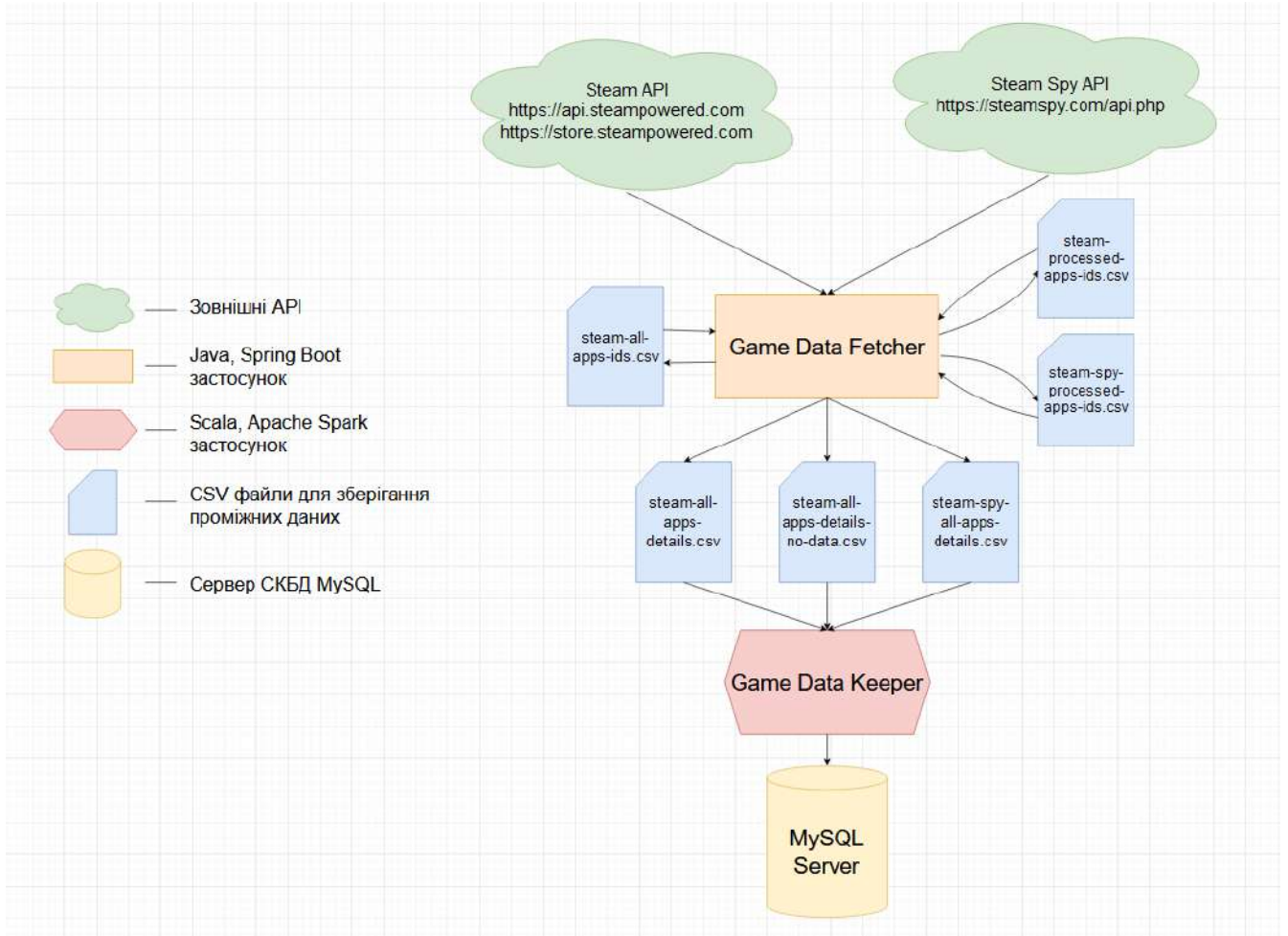


Рисунок 2.3.1. Архітектура системи

Розділ 3. Розробка застосунку для збору даних

3.1 Вибір технологій розробки

Основними завданнями застосунку є отримання даних від сторонніх API у форматі JSON, перетворення їх у формат CSV та збереження у файли. Для їх виконання найкраще буде використати:

- HTTP-клієнт, який дозволяє робити мережеві запити. Він забезпечить отримання даних від сторонніх API;
- засоби для роботи із даними у форматі JSON, оскільки сторонні API надають дані саме у форматі JSON;
- засоби для роботи із даними у форматі CSV - для збереження даних у проміжні CSV файли.

Для реалізації застосунку було обрано мову програмування Java, оскільки вона має велику кількість готових рішень для роботи із форматами JSON та CSV, а також фреймворк Spring та його модуль Spring WebFlux.

Зазвичай Spring та Spring WebFlux використовуються для розробки вебзастосунків, а точніше – реактивних вебзастосунків, оскільки модуль WebFlux реалізує підтримку парадигми реактивного програмування для Spring і базується на бібліотеці Project Reactor. [5] Реактивне програмування є функціональним та декларативним за своєю суттю. Замість опису конкретних кроків, які необхідно виконати для обробки даних послідовно, ця парадигма радше описує певний конвеєр або потік, через який дані проходять. Реактивний підхід не вимагає наявності усіх даних для їх обробки, натомість, він надає можливість для обробки даних в той момент, коли вони стають доступними. [6] Також Spring WebFlux надає зручний високорівневий HTTP-клієнт, який не вимагає ручного контролю за відкриттям та закриттям мережевих з'єднань, ручної обробки помилок, тощо. Водночас, цей клієнт має широкі можливості для налаштування, які дозволяють підібрати найкращу конфігурацію для комунікації із потрібним сервером.

3.2 Steam API

Steam надає досить широкий API, який має можливості, пов'язані із майже всіма функціями сервісу – від отримання інформації про сутність (гру, програму, відео, тощо) до перегляду поточної ситуації на внутрішньому ринку предметів та встановлення двофакторної аутентифікації для акаунту. Не всі можливості є публічними, оскільки деякі з них забезпечують виконання дій, специфічних для акаунту конкретного користувача, а тому вимагають попередньої авторизації. Майже всі доступні можливості Steam API можна знайти на вебресурсах [7] та [8]. Проте, через велику кількість можливостей, не усі вони описані у документації, а тому деякі з них були знайдені та досліджені сторонніми розробниками. Деякі їх знахідки описані на вебресурсі [9].

Для отримання даних про сутності буде використано два ендпоїнти. Вони є публічними, тому для їх використання не потрібно проходити авторизацію чи використовувати ключ для доступу до API.

Перший - GET <http://api.steampowered.com/ISteamApps/GetAppList/v2>. Він дозволяє отримати список усього контенту, який доступний у Steam – ігор, їх доповнень, програм, тощо. Він не має ніяких додаткових параметрів і повертає список сутностей з їх ідентифікаторами та назвою. Отримані ідентифікатори будуть використані для отримання інформації про конкретну сутність.

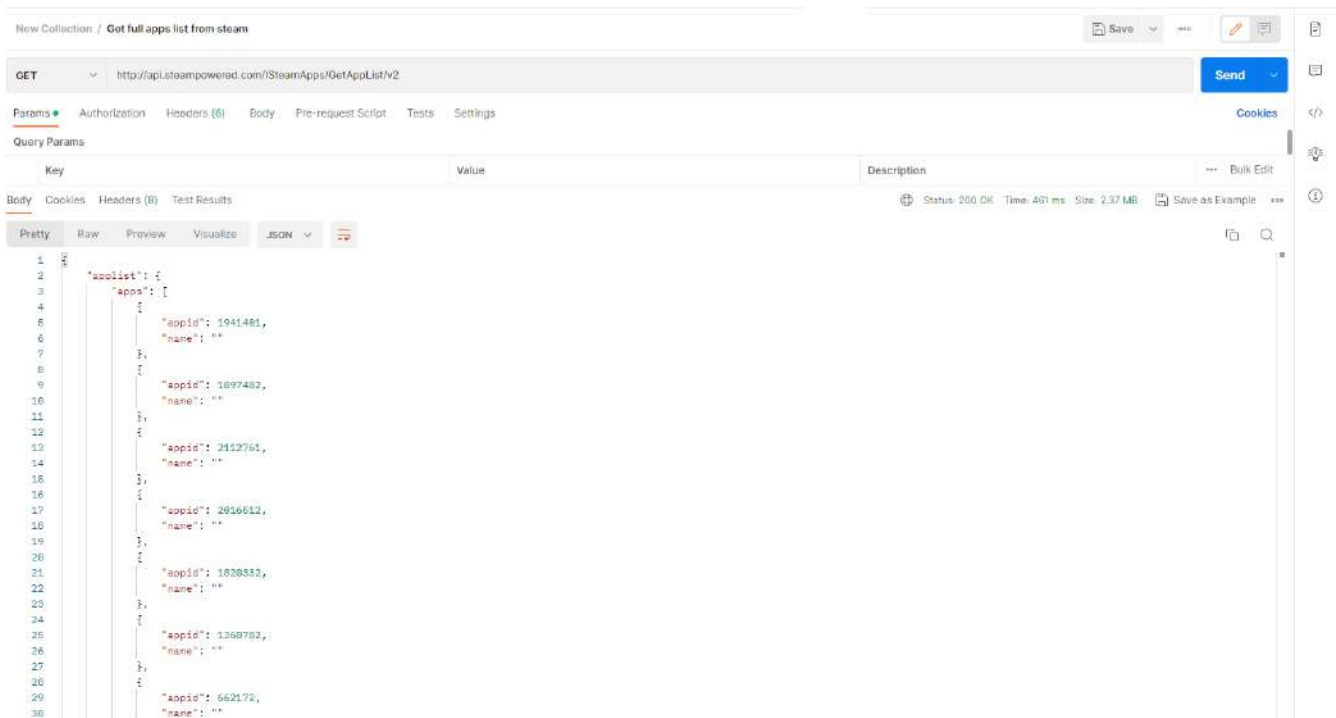


Рисунок 3.2.1. Частина відповіді на запит для отримання списку усього доступного контенту від Steam API

Другий – GET <https://store.steampowered.com/api/appdetails>. Він повертає детальну інформацію про сутність і включає її тип (гра, програма, відео, тощо) назву, опис, ціну, дату виходу, список розробників та видавців, список жанрів та багато іншої інформації. Також для різних типів сутностей дані, що повертаються, можуть мати різну структуру, наприклад, деякі поля будуть присутні або відсутні. Ендпоінт має один обов’язковий HTTP-параметр – “appid”. Йому потрібно надати ідентифікатор сутності, інформацію про яку потрібно отримати. Також є можливість вказати два додаткових параметри – “cc”, який дозволяє вказати код країни, у валюті якої потрібно рахувати ціни, та “l”, який дозволяє вказати назву мови, якою потрібно перекласти текстову інформацію у відповіді.

The screenshot shows a REST client interface with the following details:

- Request:** GET https://store.steampowered.com/api/appdetails?appid=382940&cc=US&l=english
- Query Params:**

Key	Value	Description
appid	382940	
cc	US	
l	english	
- Response Body (JSON):**

```

{
  "382940": {
    "success": true,
    "data": {
      "type": "Game",
      "name": "HITMAN™",
      "steam_appid": 236879,
      "required_age": "17",
      "is_free": false,
      "controller_support": "Full",
      "detailed_description": "Play the beginning of HITMAN for free and become the master assassin. The first location in the game is a secret training facility, where players step into the shoes of Agent 47 for the very first time and must learn what it takes to become an agent for the International Contract Agency. ESSENTIAL COLLECTION: Buy this bundle to save 36% off the three recent Hitman games! The bundle includes HITMAN – The Complete First Season, Hitman Absolution and Hitman Blood Money. Agent 47 is back with a vengeance and vengeance has seldom tasted sweeter. Reviews and Accolades: Recommended Agent 47 is back with a vengeance, and vengeance has seldom tasted sweeter. Eurogamer: Hitman has gone from an enjoyable dalliance to a minor obsession to a bona fide stealth classic. Kotaku: They've managed to prove to the rest of the industry that AAA games done episodically can work, and can work incredibly well. Steam Community: The Hitman experiment has been a success. Hitman-2016-complete-first-season-review/tumblr_contentbuffer036&utm_medium=social&utm_source=twitter&utm_campaign=buffer: Digital Spy: A beautiful puzzle box of a game. The Guardian: The Guardian: The new episodic Hitman is one of Agent 47's greatest adventures to date and is one of the best stealth games to be released in years. Gaming Nexus: A true strength that's suited to IO's vast and nuanced sandboxes. The new money lives up to the Blood Money in this darkly comic, icy cool

```

Рисунок 3.2.2. Частина відповіді на запит для отримання детальної інформації про сутність від Steam API

3.3 Steam Spy API

Steam Spy має значно простіший API, який лише надає інформацію про сутності за різними критеріями. Усі його можливості описані на вебресурсі [10].

Для виконання усіх запитів використовується основний URL - <https://steamspy.com/api.php>. Для усіх запитів обов'язковим HTTP-параметром є "request" – він визначає тип запиту та які ще параметри потрібно використати для отримання даних. Для отримання детальної інформації про сутність необхідно встановити значення "request" = "appdetails", а також для параметру "appid" вказати ідентифікатор сутності. Оскільки Steam Spy повністю базується на API, який надає Steam, то ідентифікатори сутностей повністю співпадають із тими, які можна отримати від Steam API.

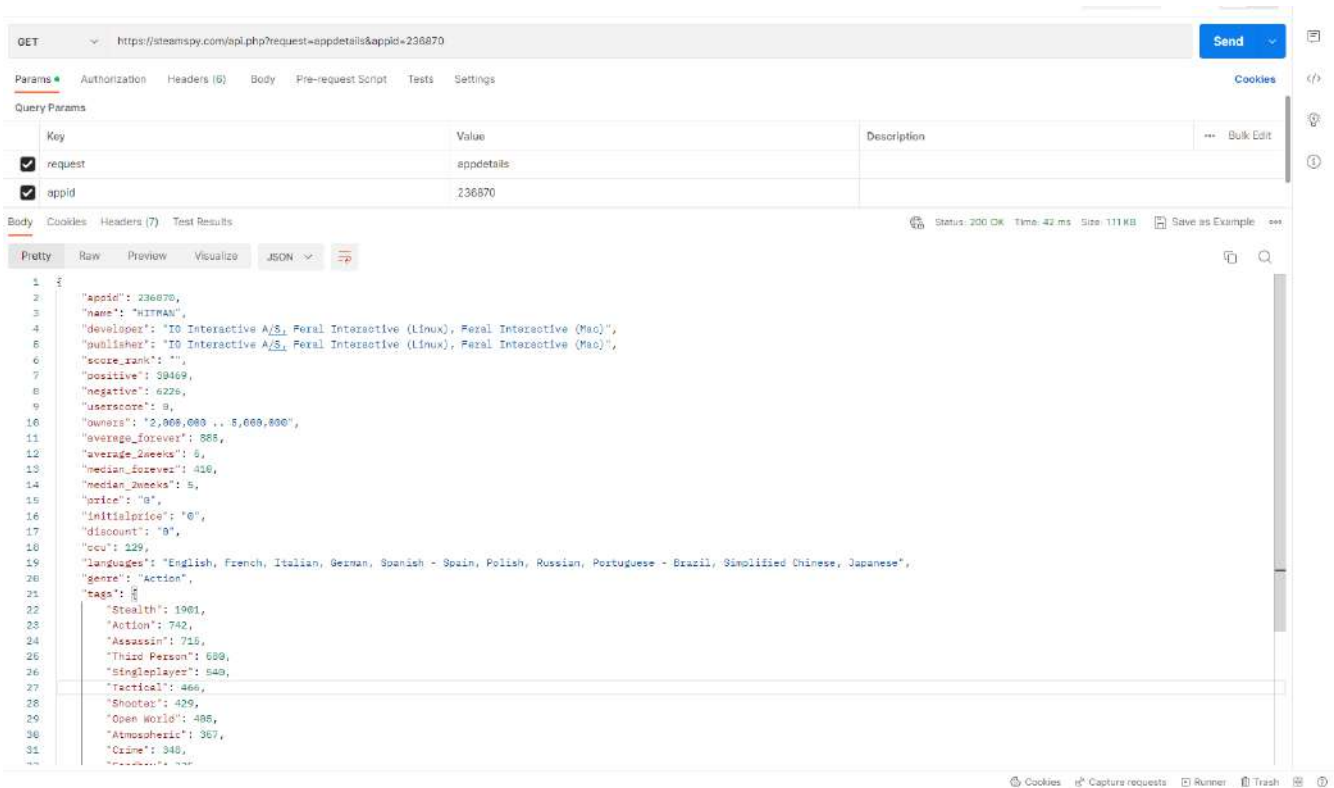


Рисунок 3.3.1. Частина відповіді на запит для отримання детальної інформації про сутність від Steam Spy API

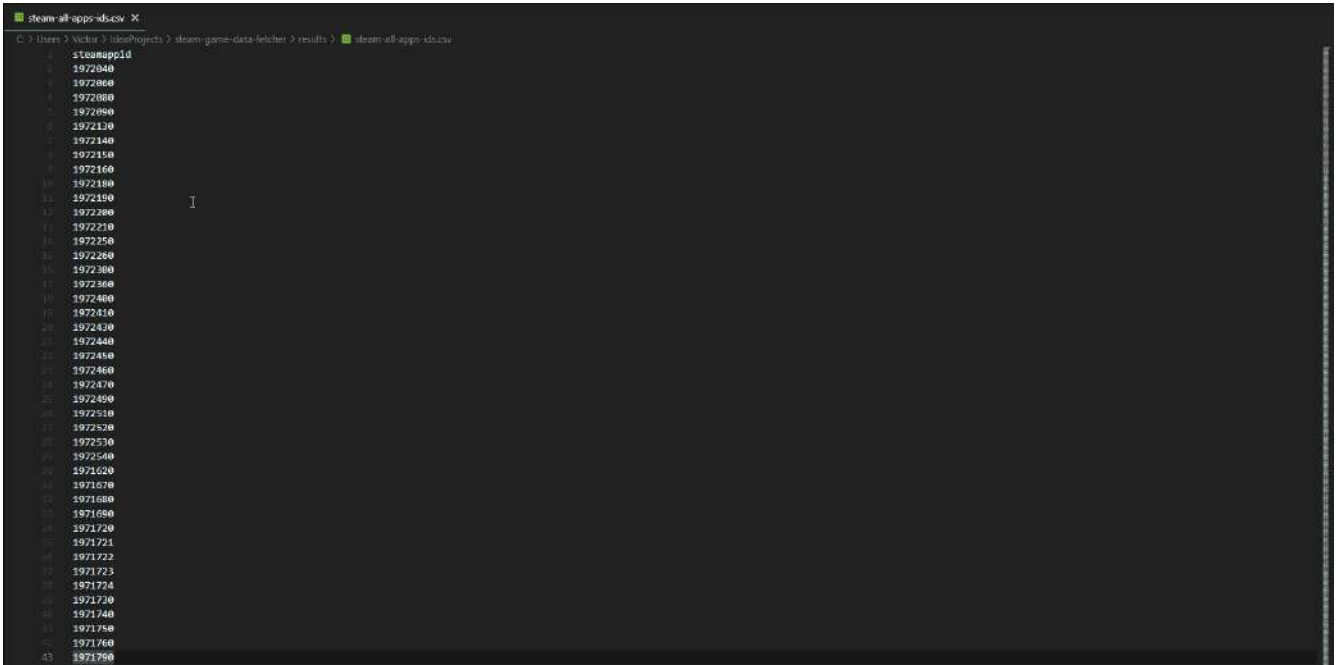
3.4 Розробка застосунку

Принцип роботи застосунку, є досить простим. Його задачами є:

- Отримати список ідентифікаторів усіх сутностей від Steam API.
- Для кожної сутності із цього списку отримати інформацію про неї зі Steam API та Steam Spy API і записати цю інформацію у файли.

Кількість сутностей у Steam оновлюється не дуже часто, а тому немає сенсу кожен раз при необхідності отримати список усіх ідентифікаторів сутностей звертатися до Steam API. Достатньо буде отримати дані один раз і зберегти їх локально. Зберігання цих даних в оперативній пам'яті не є найкращим рішенням, оскільки після завершення роботи застосунку вони зникнуть і доведеться знову звертатися до Steam API. Тому для зберігання усіх ідентифікаторів було створено файл steam-all-apps-ids.csv. При старті застосунку для отримання усіх ідентифікаторів достатньо просто зчитати цей файл. Він має досить просту

структуру – лише один стовпчик із назвою “steamappid”, а в кожному рядку зберігається ідентифікатор сутності. На момент розробки у файлі було 157523 унікальних ідентифікатори.



```

steam-all-apps-ids.csv
C:\Users\Victor\Documents\steam-game-data-feicher\results\steam-all-apps-ids.csv
1,14088991d
2,1972840
3,1972860
4,1972880
5,1972890
6,1972130
7,1972140
8,1972150
9,1972160
10,1972180
11,1972190
12,1972200
13,1972210
14,1972250
15,1972260
16,1972300
17,1972360
18,1972400
19,1972410
20,1972430
21,1972440
22,1972450
23,1972460
24,1972470
25,1972490
26,1972510
27,1972520
28,1972530
29,1972540
30,1971520
31,1971570
32,1971580
33,1971600
34,1971690
35,1971720
36,1971721
37,1971722
38,1971723
39,1971724
40,1971730
41,1971740
42,1971750
43,1971760
44,1971790

```

Рисунок 3.4.1. Частина файлу *steam-all-apps-ids.csv*

Після отримання усіх ідентифікаторів потрібно отримати дані про конкретну сутність на основі її ідентифікатора через Steam API та Steam Spy API. Цей процес дуже схожий для обох API, тому його можна описати на прикладі роботи зі Steam API.

Процес отримання даних про сутності зі Steam API складається із декількох кроків:

1. Отримати усі ідентифікатори сутностей, які ще не були оброблені.
2. Для кожного ідентифікатора зробити запит для отримання даних від Steam API.
3. Перетворити отриманий результат із формату JSON у формат CSV та записати у файл.

Розглянемо кожен крок детальніше.

Для отримання усіх ідентифікаторів сутностей, які ще не були оброблені, потрібно від усіх ідентифікаторів із файлу `steam-all-apps-ids.csv` відняти ті, які вже були оброблені. Вони зберігаються в окремому файлі, який має назву `steam-processed-apps-ids.csv` і має таку ж структуру, як і `steam-all-apps-ids.csv`. Цей файл формується і оновлюється після отримання даних про конкретну сутність від Steam API. Зберігання списку усіх оброблених ідентифікаторів дозволяє при завершенні та відновленні роботи застосунку не починати обробку усіх ідентифікаторів заново. Для роботи зі Steam Spy API існує аналогічний файл зі списком оброблених ідентифікаторів.

Наступним кроком є отримання даних від Steam API. Комунікація зі Steam API відбувається через окремий клас-сервіс `SteamService` і використовує HTTP-клієнт `WebClient`, який надає `Spring WebFlux`. При роботі зі стороннім публічним API ніколи не можна бути впевненим, що API буде стабільним, тому варто перебачити та врахувати сценарії, при яких процес отримання даних може піти не за планом. Основні можливі проблеми:

- тимчасова недоступність API. Наприклад, через велике навантаження API може дуже довго відповідати на клієнтські запити, або взагалі на них не відповідати. Іншим прикладом є проведення технічних робіт, під час яких API може бути повністю недоступним;
- наявність обмеження на кількість запитів від одного клієнта. Майже усі публічні API встановлюють певні обмеження на кількість запитів від одного клієнта щоб уникнути мережових атак, пов'язаних із надсиланням великої кількості запитів для перевантаження сервера.

Ці проблеми були однією з причин використання окремих файлів для зберігання уже оброблених ідентифікаторів. Такі файли дозволяють відновити роботу застосунку при її завершенні через проблеми, спричинені сторонніми API. Також для зменшення впливу цих проблем були додані додаткові налаштування

для екземплярів класу `WebClient`, які можна побачити у *Додатку А*. Було встановлено максимальний час підключення до сервера рівний 30 секундам, максимальний час отримання відповіді рівний 60 секундам та розмір буферу рівний 20 мегабайтам. Останнє налаштування потрібне для того, щоб отримати список усіх ідентифікаторів від Steam API, оскільки їх кількість досить велика (кінцева відповідь у форматі JSON займає декілька мегабайт).

У класі `SteamService` знаходиться основна логіка отримання даних від Steam API, його методи можна побачити у *Додатку Б*. На прикладі цього класу можна побачити, що працювати із класом `WebClient` дуже зручно, оскільки для звернення до стороннього ресурсу достатньо обрати HTTP-метод, вказати URL ресурсу, якщо потрібно – додати HTTP-параметри. Також можна налаштувати логіку відправлення повторних запитів до сервера, якщо на попередній запит сервер повернув помилку. Відправлення повторних повідомлень використовується при отриманні даних про конкретну сутність за її ідентифікатором. Якщо сервер має повернути якісь дані, наприклад, документ у форматі JSON чи XML і його необхідно перетворити на екземпляр певного Java класу, то достатньо викликати спеціальний метод, у якому вказати назву класу. Оскільки Spring має інтеграцію із багатьма бібліотеками для серіалізації та десеріалізації даних у різні формати, він використовує їх для автоматичного перетворення отриманих даних із відповіді сервера в екземпляр бажаного класу. Цей спосіб використовується для перетворення відповіді від Steam API в екземпляр класу `SteamAppDetailsResponse`, який містить усю інформацію, яку можна отримати про сутність.

Для серіалізації та десеріалізації даних у формат JSON використовується бібліотека Jackson. Jackson – це набір інструментів для обробки даних для мови Java та платформи JVM, який надає можливості для парсингу та створення JSON-файлів а також перетворення Java об'єктів у формат JSON та створення їх із формату JSON. [11] Цю бібліотеку також використовує Spring як основну для

роботи із даними у форматі JSON. Класом, який відображає дані про сутності, отримані від Steam API є `SteamAppDetailsResponse`. Зазвичай, для перетворення даних із формату JSON в екземпляр Java класу достатньо визначити поля бажаного класу, надавши їм типи даних, які відповідають типам даних формату JSON, та налаштувати перетворення назв полів у JSON-документі в назви полів Java класу. Але при роботі зі Steam API з'ясувалось, що деякі поля JSON-документу, який повертає Steam API, мають нефіксований тип даних. Наприклад, деякі поля, якщо вони є JSON-об'єктами і мають вкладені поля, повертаються як JSON-об'єкти, а якщо не мають – як JSON-масиви. Таким чином, неможливо задати однозначний тип даних Java, який відповідає одночасно і об'єкту, і масиву, а тому потрібно створювати власний клас, який буде виконувати десеріалізацію JSON-документу в екземпляр класу `SteamAppDetailsResponse`. Для цього було створено клас `SteamAppDetailsResponseDeserializer`, який обробляє такі випадки і забезпечує правильну десеріалізацію. Реалізацію класу можна побачити у *Додатку В*.

Для запису проміжних даних у файли формату CSV використовується бібліотека Apache Commons CSV. Вона була розроблена для створення єдиного та простого інтерфейсу для читання та запису даних у форматі CSV. [12] Для запису даних, отриманих зі Steam API використовується клас `SteamDataCsvFileWriter`, реалізацію якого наведено у *Додатку Г*. Він наслідується від абстрактного класу `AbstractCsvFileWriter`, який містить логіку для налаштування формату CSV файлів та створення допоміжних екземплярів класів, що використовуються для запису у файли. Цю логіку було винесено в абстрактний клас, оскільки для збереження даних, отриманих від Steam Spy API, існує окремий клас `SteamSpyDataCsvFileWriter`, який використовує таку ж логіку для налаштування формату CSV. Використання абстрактного класу у цьому випадку дозволяє запобігти дублюванню коду. Основною задачею класу `SteamDataCsvFileWriter` є

не збирати інформацію із трьох різних файлів, а використовувати лише один, який містить усю необхідну інформацію.

Повний цикл отримання даних від Steam API та їх збереження забезпечує клас `SteamAppDetailsProcessor`, реалізацію якого можна побачити у *Додатку Д*. Він для кожного ідентифікатора отримує дані про сутність за допомогою `SteamService`, а потім записує отримані результати у файли за допомогою `SteamDataCsvFileWriter`. Важливою функцією класу `SteamAppDetailsProcessor` є забезпечення певної паузи між обробкою кожного ідентифікатора. Ця потреба виникає через наявність обмеження у Steam API – він дозволяє зробити максимум 200 запитів кожні 5 хвилин для одного клієнта. Таким чином, якщо робити усі запити послідовно, то для дотримання цього ліміту потрібно робити по одному запиту із мінімальним інтервалом між ними в 1.5 секунди. Одним зі способів обійти це обмеження є використання проксі-серверів. Це рішення дозволить видати запити, які робляться одним клієнтом (застосунком), за запити, відправлені декількома різними клієнтами, що дозволяє обробляти декілька ідентифікаторів одночасно. Проте, використання публічних проксі-серверів не є надійним, оскільки розробник не має ніякого контролю над ними, тому при роботі з ними доведеться витратити зайвий час на обробку можливих помилок, а для розгортання власних проксі-серверів потрібно витратити додатковий час та ресурси, що не буде дуже ефективним у даному випадку, оскільки кількість запитів, які необхідно зробити, не є дуже великою. Тому було обрано простіше рішення із встановленням затримки між обробкою ідентифікаторів у 1.6 секунди. Ця затримка трохи більша, ніж мінімально можлива, але вона дозволяє врахувати неточності при обрахунку часу між запитами, оскільки немає жодної інформації про те, з якою точністю рахує час Steam API.

Для усіх описаних вище класів існують аналогічні, вони забезпечують роботу зі Steam Spy API. Це `SteamSpyService`, `SteamSpyAppDetailsResponse`,

SteamSpyAppDetailsResponseDeserializer, SteamSpyDataCsvFileWriter та SteamAppDetailsProcessor. Завдяки схожому процесу отримання даних від обох API було використано можливість створення абстрактних класів для реалізації спільної логіки в одному місці та уникнення дублювання коду. Файлову структуру проекту застосунку можна побачити у *Додатку E*.

Розділ 4. Розробка застосунку для обробки даних

4.1 Вибір технологій розробки

Основними завданнями застосунку є отримання даних із файлів, що були підготовлені застосунком, який відповідає за збір даних, їх обробка, очищення, трансформація та збереження у базі даних. При виборі технологій також варто врахувати декілька аспектів:

- дані, із якими потрібно працювати, не проходили процес попередньої обробки, тому обрані технології мають забезпечувати роботу із неструктурованими та “сирими” даними;
- обрані технології мають підтримувати обробку великих об’ємів даних, оскільки файл із даними, отриманими від Steam API має досить значний об’єм – 1,02 гігабайта.

Враховуючи вищезазначені вимоги та деталі, найкращим вибором для виконання поставленого завдання є фреймворк Apache Spark.

Apache Spark – це двигун для роботи із великими даними. [13] Він має велику кількість можливостей, серед яких:

- робота як зі структурованими, так і з неструктурованими даними;
- робота зі статичними даними та даними, що отримуються в реальному часі;
- наявність інструментів для виконання аналізу даних;
- наявність інструментів для машинного навчання;
- наявність клієнтів для декількох мов програмування: Scala, Java, Python, R;
- підтримка великої кількості джерел даних (різні файлові формати, реляційні та нереляційні бази даних, брокери повідомлень).

Spark базується на master-worker архітектурі. Основою цієї архітектури є поділ її компонентів на два види – головний (master) та підлеглий (worker). Зазвичай, використовується один головний та декілька підлеглих вузлів. Їх основна відмінність – це функції, які виконує кожен із них.

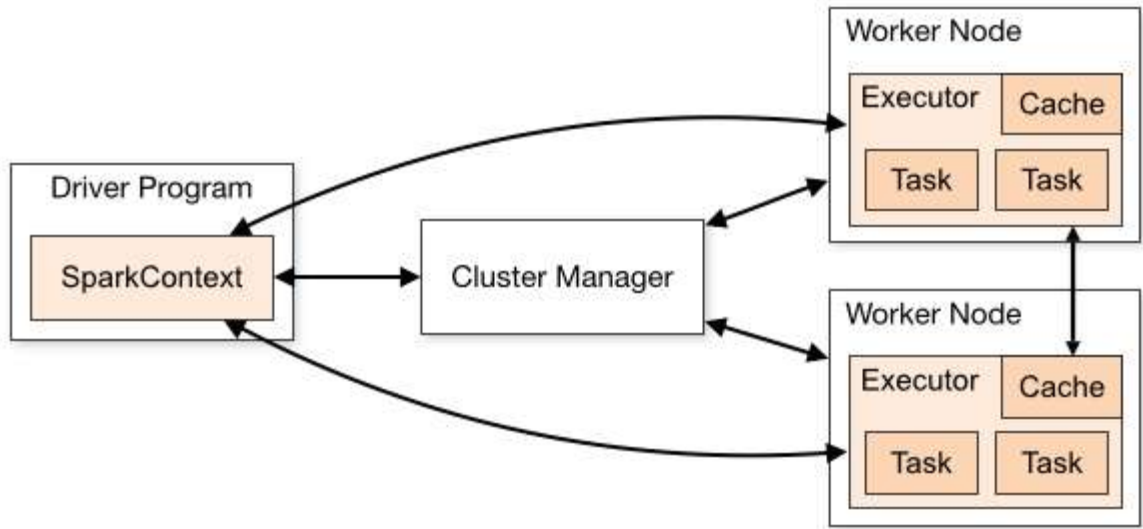


Рисунок 4.1.1. Master-worker архітектура Apache Spark

Головний компонент у Spark – це Driver. Він відповідає за поділ усіх задач, які має виконати застосунок, на етапи, побудову послідовності їх виконання та розподіл задач між підлеглими вузлами. Для представлення задач Spark використовує направлений ациклічний граф, який будується на основі усіх дій, які виконує застосунок. Цей граф також додатково оптимізується для забезпечення якнайбільшої продуктивності, а потім на його основі виокремлюються етапи, необхідні для виконання кожної задачі. Після цього виконується планування та розподіл усіх задач між доступними підлеглими компонентами, які мають назву Executor.

В свою чергу, кожен підлеглий компонент відповідає за безпосереднє виконання задач, отриманих від головного компонента. Саме підлегли вузли виконують читання даних, операції для їх трансформації та їх запис. Підлегли вузли також можуть обмінюватись даними один з одним, якщо виникає така

потреба. Якщо один із підлеглих вузлів не закінчив виконання своїх задач із певних причин (наприклад, фізична машина, на якій відбувалось виконання, стала недоступна), головний вузол виконує перерозподіл його задач між іншими доступними підлеглими вузлами.

Spark надає декілька можливостей для розгортання застосунків. Серед них:

- розгортання усього застосунку в одному процесі. Це найпростіший варіант, оскільки не потребує наявності декількох машин для роботи застосунку, а головний та підлеглі вузли знаходяться в одному процесі операційної системи;
- розгортання застосунку у вигляді декількох процесів на одній фізичній машині. Це один процес для головного вузла та один або більше процесів для підлеглих вузлів;
- розгортання застосунку у вигляді різних процесів на різних фізичних машинах. Для цього варіанту необхідно використати один із менеджерів кластерів, що підтримуються у Spark: власний менеджер кластеру, який надає Spark, Apache Mesos, Hadoop YARN чи Kubernetes.

Для виконання задач роботи було обрано варіант із розгортанням застосунку в одному процесі, оскільки об'єм даних не є настільки великим, що потрібно розгортати окремий кластер для роботи із ними.

У якості мови програмування було обрано Scala, оскільки Spark розроблений саме нею. Це дозволяє використовувати увесь доступний API, який надає Spark, а при необхідності навіть розширювати його можливості для свої потреб.

4.2 Структура бази даних

Для побудови структури бази та визначення необхідних таблиць необхідно розглянути доступні дані, обрати ключову інформацію, знайти зв'язки між

даними, а також визначити дані, які не є корисними і можуть бути проігноровані. Аналізуючи дані, отримані з обох джерел, можна помітити наступні деталі:

- кожна сутність має певний наперед визначений тип, тому можна прослідкувати зв'язок виду “один до багатьох” між типом та сутністю. Таким чином, тип може стати окремою моделлю у предметній області;
- кожна сутність має від 0 до декількох розробників, видавців, категорій та жанрів, тобто між сутністю та кожною із наведених моделей існує зв'язок виду “багато до багатьох”;
- для кожної сутності як Steam, так і Steam Spy надають деяку спільну інформацію (наприклад, назву, розробників та видавців), проте також присутня інформація, специфічна для кожного із сервісів. Тому доречним буде сформулювати одну кореневу сутність, яка буде містити основну інформацію, а також допоміжні, у яких будуть зберігатись дані, специфічні для кожного із сервісів, причому допоміжні сутності будуть пов'язані із кореневою зв'язком виду “один до одного”. Ще одним аргументом на користь такого рішення є те, що про у кожному сервісі може не бути інформації про деякі сутності;
- деяка інформація, яку надають Steam та Steam Spy не є корисною в рамках системи, що розроблюється. Наприклад, Steam для кожної сутності повертає список скріншотів та трейлерів, які користувач бачить на сторінці цієї сутності на вебсайті або у настільному додатку. Така інформація може бути вилучена, оскільки ніяк не використовується.

На основі проведеного аналізу було розроблено схему бази даних, яку можна побачити на рисунку 4.2.1.



Рисунок 4.2.1. Схема бази даних


Кореневу сутність представляє таблиця “application”. Вона пов’язана зв’язком виду “багато до одного” з таблицею “application_type”, яка представляє тип сутності, а також зв’язками виду “багато до багатьох” із таблицями “developer”, “publisher”, “category” та “genre”, які представляють розробника, видавця, категорію та жанр відповідно. Для реалізації зв’язку виду “багато до багатьох” було використано проміжні таблиці. Також вона пов’язана зв’язками виду “один до одного” із таблицями “steam_application_details” та

“steam_spy_application_details”, які представляють дані про сутність, отримані від сервісів Steam та Steam Spy відповідно.

4.3 Розробка застосунку

Основними задачами застосунку є зчитування даних із попередньо підготовлених файлів формату CSV, їх обробка, трансформація та запис у базу даних. Усі ці завдання дозволяє виконати Spark. Для початку роботи зі Spark необхідно ініціалізувати екземпляр класу `SparkSession`, через який відбувається робота із Spark API.

Першим етапом роботи є зчитування даних із файлів формату CSV. Для цього було створено об’єкт `CsvFileReader`, який використовує `SparkSession` для зчитування даних. Цей об’єкт встановлює налаштування формату CSV, які потребує Spark для правильного читання, а також ім’я файлу, із якого необхідно зчитати дані. Результатом зчитування даних є датафрейм.



```

1 package ua.edu.ukma.diploma.pinkevych
2 package reader
3
4 import ...
5
6 object CsvFileReader {
7
8   def read(sparkSession: SparkSession, csvFormat: CsvFormat, fileName: String): DataFrame = {
9     sparkSession.read
10      .option("header", csvFormat.header)
11      .option("multiline", csvFormat.multiline)
12      .option("delimiter", csvFormat.delimiter)
13      .option("escape", csvFormat.escapeCharacter)
14      .csv(fileName)
15   }
16 }
17
18 }
19
20

```

Рисунок 4.3.1. Реалізація об’єкту `CsvFileReader`

Датафрейм – це основна одиниця представлення даних, яку використовує Spark. Він складається зі стовпчиків, кожен із яких має назву і представляє певну властивість, та рядків, кожен з яких представляє одну сутність. Датафрейми дуже схожі на таблиці, які представляють дані у реляційних базах даних. Датафрейми є

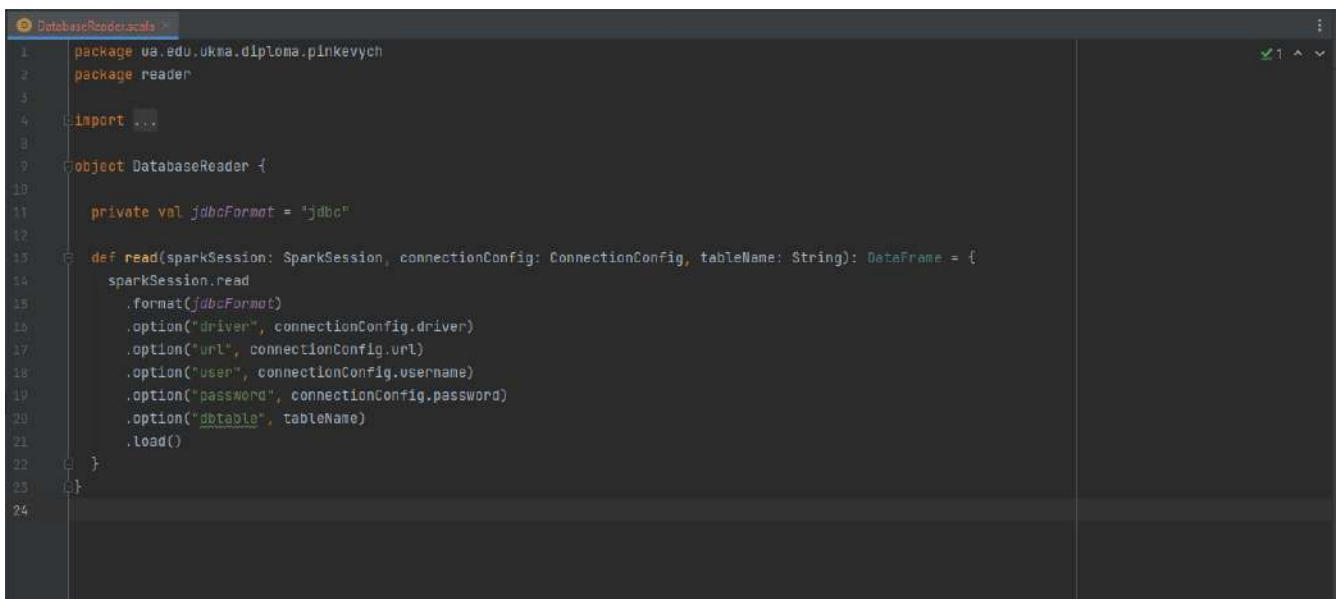
незмінними і будь-яка операція, яка виконується над датафреймом, перетворює його на інший датафрейм, але не змінює початковий датафрейм. Кожен датафрейм має схему, яка визначає, які стовпчики наявні у датафреймі і тип даних для кожного з них.

За допомогою `CsvFileReader` необхідно зчитати 3 файли, дані з яких будуть використані для подальшої трансформації та збереження. Це `steam-all-apps-details.csv`, `steam-all-apps-details-no-data.csv` та `steam-spy-all-apps-details.csv`. Усі отримані після читання датафрейми будуть застосовані в подальшому для отримання необхідних даних. Проте, якщо один датафрейм необхідно використати декілька разів, то усі операції, потрібні для його отримання, будуть виконані заново. Таким чином, якщо один із отриманих вище датафреймів потрібно використати 5 разів, то він буде прочитаний із файлу 5 разів, що є дуже неефективним. Щоб уникнути перерахування датафрейму `Spark` дозволяє зберегти його у процесі виконання на певний час. Датафрейм може бути збережений в оперативній пам'яті машини, в постійній пам'яті, або використовувати обидва типи пам'яті. Після виконання усіх необхідних дій датафрейм може бути видалений із пам'яті. У застосунку усі датафрейми було збережено лише в оперативній пам'яті одразу після зчитування та видалено в кінці його роботи.

Для перетворення та збереження даних у кожную таблицю бази даних було розроблено єдину схему роботи із початковими датафреймами. Розглянемо її на прикладі збереження у базу даних інформації про відношення між сутностями та їх розробниками.

У датафреймах, отриманих із файлів, розробники сутності зберігаються у стовпчику `CSV` файлу, який представляє собою стрічку у форматі `JSON`, що представляє `JSON`-масив із назвами компаній розробників. У базі даних розробники зберігаються у вигляді таблиці, яка має 2 поля – ідентифікатор та назву компанії. Для зберігання відношень між сутністю та її розробниками у базі

даних існує таблиця, яка представляє це відношення у вигляді таблиці із 2 полями: ідентифікатором сутності та ідентифікатором розробника. Тому для зберігання відношень потрібно для кожної сутності із датафрейму додати ідентифікатори усіх її розробників. Для цього треба зчитати таблицю з розробниками в датафрейм. Для цього було створено об'єкт `DatabaseReader`, який дуже схожий на `CsvFileReader`. Він використовує `SparkSession`, встановлює необхідні налаштування для під'єднання до бази даних, назву таблиці, яку необхідно прочитати, та повертає датафрейм із даними з таблиці. Реалізацію об'єкту `DatabaseReader` наведено на рисунку 4.3.2.



```

1 package ua.edu.ukna.diptona.pinkevych
2 package reader
3
4 import ...
5
6 object DatabaseReader {
7
8     private val jdbcFormat = "jdbc"
9
10    def read(sparkSession: SparkSession, connectionConfig: ConnectionConfig, tableName: String): DataFrame = {
11        sparkSession.read
12            .format(jdbcFormat)
13            .option("driver", connectionConfig.driver)
14            .option("url", connectionConfig.url)
15            .option("user", connectionConfig.username)
16            .option("password", connectionConfig.password)
17            .option("database", tableName)
18            .load()
19    }
20 }
21
22 }
23
24

```

Рисунок 4.3.2. Реалізація об'єкту `DatabaseReader`

Після отримання усіх необхідних даних можна почати їх перетворення. За виконання перетворень над датафреймами сутностей та розробників відповідає об'єкт `SteamApplicationDeveloperTransformer`, реалізацію якого можна побачити на рисунку 4.3.3. Першим кроком є підготовка даних про сутність, а саме – отримання усіх імен її розробників у форматі масиву стрічок із JSON-масиву. Для цього використовується окрема визначена користувачем функція та вищезгадана бібліотека `Jackson`, оскільки вона також доступна і для мови `Scala`.

```

1 package us.edu.ukma.diploma.pinkevych
2 package transformer.developer
3
4 import ...
5
13 object SteamApplicationDeveloperTransformer {
14
15     def transform(
16         steamAppsDetailsDF: DataFrame,
17         developersDF: DataFrame,
18         broadcast: Broadcast[JsonMapperProvider]
19     ): DataFrame = {
20
21         val developerArrayFromJson = udf((rawDeveloperNamesString: String) =>
22             TransformationUtil.stringArrayFromJson(rawDeveloperNamesString, broadcast.value.getJsonMapper)
23         )
24
25         val steamApplicationDevelopersDF = steamAppsDetailsDF
26             .select(
27                 col( colName = "steam_app_id"),
28                 explode(developerArrayFromJson(col( colName = "developers"))).as( alias = "developer_name")
29             )
30
31         steamApplicationDevelopersDF
32             .as( alias = "steam_application_developers")
33             .join(
34                 developersDF.as( alias = "developers"),
35                 steamApplicationDevelopersDF("developer_name") === developersDF("developer_name")
36             )
37             .select(
38                 col( colName = "steam_application_developers.steam_app_id").as( alias = "steam_application_id"),
39                 col( colName = "developers.id").as( alias = "developer_id")
40             )
41             .dropDuplicates()
42     }
43 }
44

```

Рисунок 4.3.3. Реалізація об'єкту *SteamApplicationDeveloperTransformer*

Визначені користувачем функції (UDF) – це функції, яка дозволяють виконувати операції над стовпцями у датафреймах, але не є частиною фреймворку Spark, оскільки визначені сторонніми розробниками. Функції можуть приймати від 0 до 22 аргументів, але повертають лише одне значення, яке зберігається у стовпчику датафрейму. [14] Ця концепція є дуже потужною, оскільки дозволяє легко розширювати можливості фреймворку та додавати необхідний функціонал. Функції можуть бути написані мовою, якою виконується розробка, таким чином дозволяючи використовувати будь-які засоби мови, включно зі сторонніми бібліотеками та фреймворками. Ці функції можуть бути викликані як напряму, так і через Spark SQL. Для їх використання достатньо попередньо зареєструвати

бажані функції через `SparkSession` або використовуючи спеціальні функції, які `Spark` надає у клієнтах для кожної підтримуваної мови.

Також для перетворення JSON-масиву у масив стрічок необхідно використати об'єкт класу `JsonMapper`. Якщо створювати його напряду при виклику функції для перетворення, то новий об'єкт буде створюватись для кожного виклику, що є дуже неефективним. Найкращим рішенням буде передавати його параметром у функцію. Проте це не вийде зробити просто створивши об'єкт, оскільки усі дії над датафреймами виконуються окремими підлеглими вузлами, кожен із яких має свою пам'ять, кеш та є повністю незалежним від інших підлеглих вузлів. Проте `Spark` має механізм, який дозволяє передавати спільні дані від головного вузла усім підлеглим – трансльовані змінні.

Трансльовані змінні (`Broadcast variables`) - це змінні, які доступні усім вузлам кластера лише для читання і призначені для ефективного зберігання та кешування спільних даних кожним вузлом. Без використання трансльованих змінних `Spark` буде передавати необхідні для виконання задачі дані кожен раз перед виконанням нової задачі. Використання трансльованих змінних дозволяє передати дані лише один раз і зберігати їх у кеші вузла протягом усієї роботи. [15]

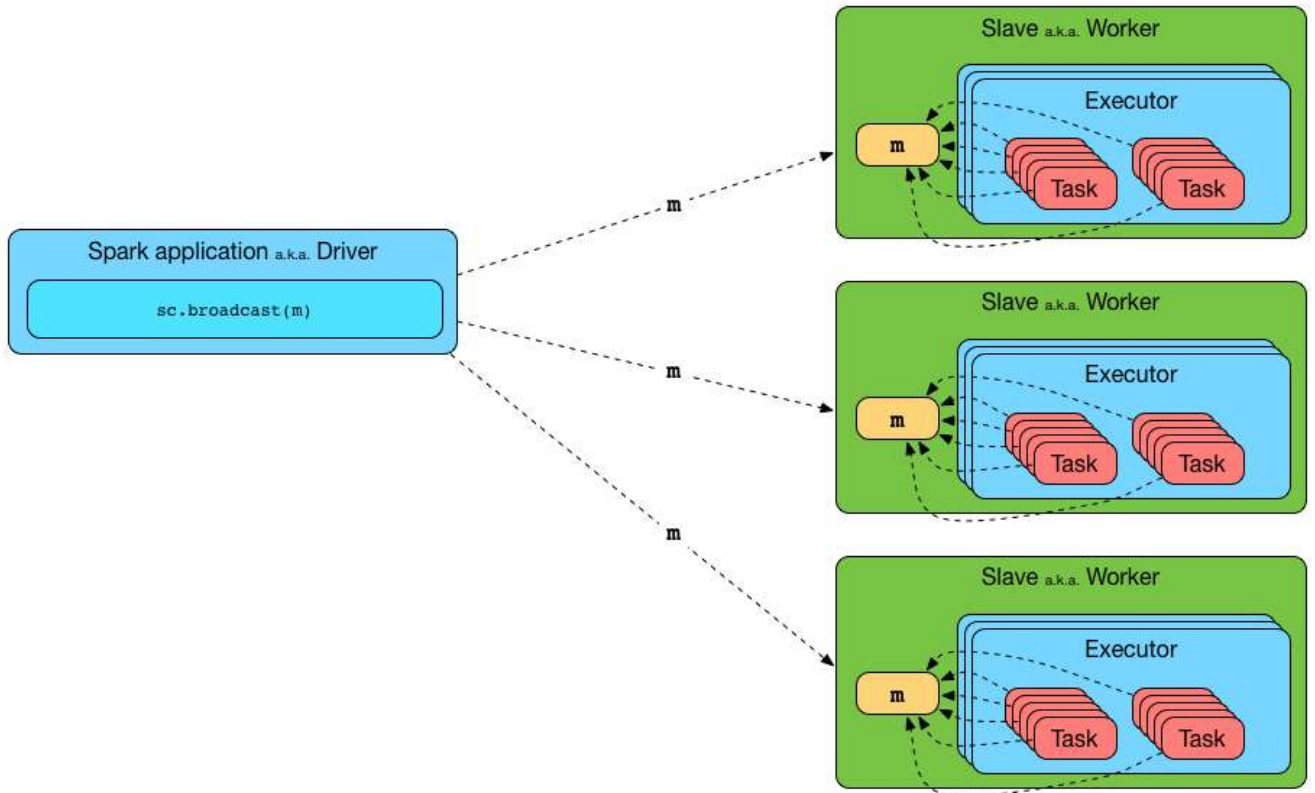


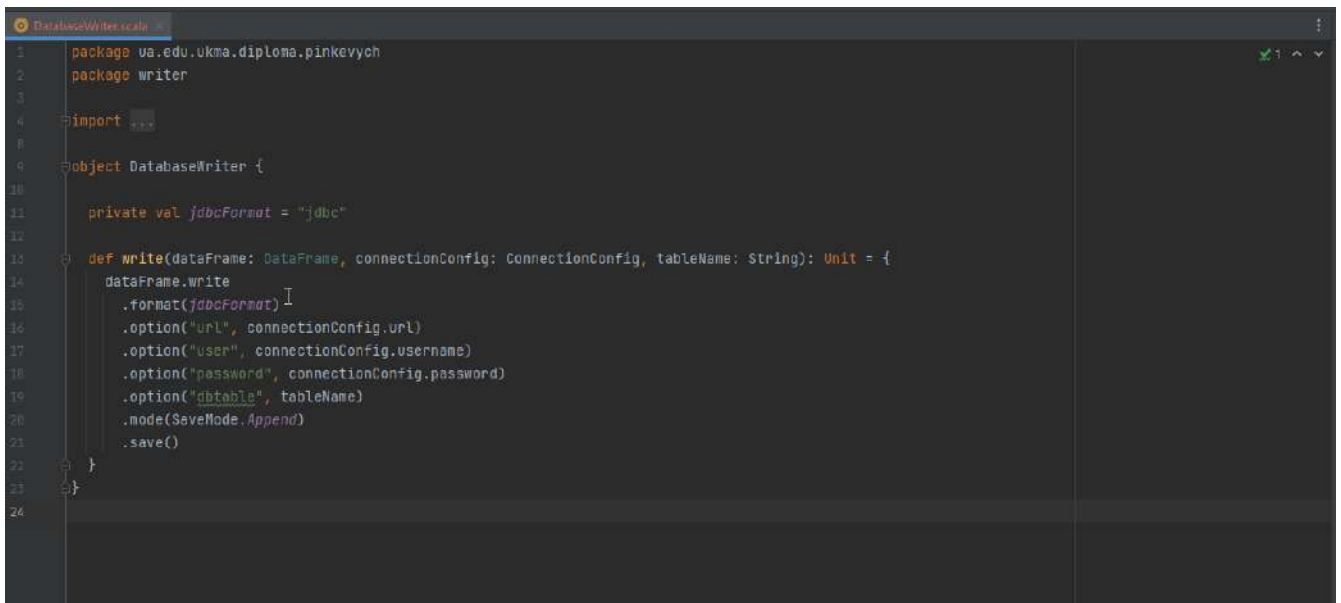
Рисунок 4.3.4. Принцип роботи трансльованих змінних в Spark

За допомогою трансльованих змінних у застосунку кожен підлеглий вузол отримує екземпляр класу `JsonMapperProvider`, який створює екземпляр класу `JsonMapper` із необхідними налаштуваннями. Таким чином, кожен підлеглий вузол має власний екземпляр класу `JsonMapper`.

Після підготовки даних про сутності потрібно пов'язати їх із даними про розробників. Для поєднання даних між датафреймами Spark надає можливість виконувати з'єднання датафреймів, які за своєю суттю є аналогічними до з'єднань таблиць у реляційних базах даних. Spark підтримує усі ті ж види з'єднань, що і реляційні бази даних: внутрішні, ліві, праві та повні. Для з'єднання датафреймів потрібно обрати два датафрейми, які потрібно з'єднати, та задати умову, при виконанні якої рядки датафреймів будуть об'єднані. Результатом об'єднання є новий датафрейм, який містить усі стовпчики із обох датафреймів, а також нові рядки, значення в яких залежать від типу з'єднання.

Для того щоб пов'язати ідентифікатори сутностей та ідентифікатори їх розробників достатньо виконати внутрішнє з'єднання на основі рівності імен розробників. Результатом внутрішнього з'єднання є датафрейм, який містить лише ті рядки, для яких виконується умова з'єднання, тобто ті, для яких було знайдено ідентифікатори розробників за їх іменами.

Після формування кінцевого датафрейму його потрібно зберегти у таблицю бази даних. Spark надає зручний API для запису датафреймів у будь-яке сховище: файл, базу даних, тощо. Для збереження датафрейму достатньо вказати формат та відповідні налаштування. Оскільки в застосунку усі датафрейми зберігаються у базу даних, було створено об'єкт `DatabaseWriter`, який для заданого датафрейму встановлює відповідні налаштування для запису у базу даних та зберігає його. Реалізацію цього об'єкту можна побачити на рисунку 4.3.5.



```

1 package ua.edu.ukma.diploma.pinkevych
2 package writer
3
4 import scala
5
6
7
8 object DatabaseWriter {
9
10
11   private val jdbcFormat = "jdbc"
12
13
14   def write(dataFrame: DataFrame, connectionConfig: ConnectionConfig, tableName: String): Unit = {
15     dataFrame.write
16       .format(jdbcFormat)
17       .option("url", connectionConfig.url)
18       .option("user", connectionConfig.username)
19       .option("password", connectionConfig.password)
20       .option("dbtable", tableName)
21       .mode(SaveMode.Append)
22       .save()
23   }
24 }

```

Рисунок 4.3.5. Реалізація об'єкту `DatabaseWriter`

Spark дозволяє як зберігати датафрейми в існуючі таблиці бази даних, так і створювати нові на основі схеми датафрейму. Проте підхід зі створенням таблиць на основі схеми датафрейму не підійде, якщо потрібно створити таблицю із певними обмеженнями, чітко вказати тип даних стовпчика чи задати первинні та зовнішні ключі, створити індекси. В таких випадках найкращим рішенням є

створення таблиць заздалегідь і налаштування Spark для запису лише в уже існуючі таблиці. Щоб уникнути необхідності створювати таблиці у базі даних вручну перед кожним запуском застосунку було створено об'єкт `SchemaInitializer`, який на початку роботи застосунку створює необхідні таблиці у базі даних, якщо вони відсутні. Таблиці створюються з використанням SQL-скриптів, що зберігаються у зовнішніх файлах, а взаємодія із базою даних відбувається через JDBC – API для мови Java та платформи JVM, який забезпечує доступ до баз даних через використання відповідних драйверів. Оскільки Spark працює на платформі JVM, то він також використовує JDBC для роботи із базами даних. Реалізацію об'єкту `SchemaInitializer` можна побачити на рисунку 4.3.6.

```

12 object SchemaInitializer {
13
14   def initializeSchema(schemaConfig: SchemaConfig, connectionConfig: ConnectionConfig): Unit = {
15     val files: Seq[String] = Seq(
16       schemaConfig.applicationTypeFile,
17       schemaConfig.applicationFile,
18       schemaConfig.steamApplicationDetailsFile,
19       schemaConfig.steamSpyApplicationDetailsFile,
20       schemaConfig.categoryFile,
21       schemaConfig.applicationCategoryFile,
22       schemaConfig.genreFile,
23       schemaConfig.applicationGenreFile,
24       schemaConfig.developerFile,
25       schemaConfig.applicationDeveloperFile,
26       schemaConfig.publisherFile,
27       schemaConfig.applicationPublisherFile
28     )
29     val queries: Seq[String] = files.map(readQueryFromFile)
30     queries.foreach(println)
31
32     val connection: Connection =
33       DriverManager.getConnection(connectionConfig.url, connectionConfig.username, connectionConfig.password)
34     val statement: Statement = connection.createStatement()
35     queries.foreach(statement.execute)
36     statement.close()
37     connection.close()
38   }
39
40   private def readQueryFromFile(fileName: String): String = {
41     val filePath = getClass.getClassLoader.getResource(fileName).getPath
42     val source = Source.fromFile(filePath)
43     val query = source.getLines().mkString
44     source.close()
45     query
46   }
47 }
48

```

Рисунок 4.3.6. Реалізація об'єкту `SchemaInitializer`

Для об'єднання процесів отримання, обробки та збереження даних про сутності та їх розробників було створено об'єкт `ApplicationDeveloperKeeper`, реалізацію якого наведено на рисунку 4.3.7. Він отримує початкові датафрейми і на їх основі, використовуючи усі описані вище об'єкти, виконує трансформацію та збереження даних. У реалізації об'єкту можна побачити використання датафрейму `steamSpyAppsDetailsWithNoDataInSteamDF`. Він містить дані про усі сутності, про які немає даних у сервісі Steam, але є у сервісі Steam Spy. Його використання дозволяє отримати ключові дані про сутність незалежно від сервісу, з якого вони були отримані. Цей датафрейм також використовується для отримання даних про видавців та жанри сутностей.

```

ApplicationDeveloperKeeper.scala
1 package ua.edu.ukma.diploma.pinkevych
2 package keeper
3
4 import ...
5
6 object ApplicationDeveloperKeeper {
7
8   def keep(
9     sparkSession: SparkSession,
10    steamAppsDetailsDF: DataFrame,
11    steamSpyAppsDetailsWithNoDataInSteamDF: DataFrame,
12    applicationConfig: ApplicationConfig,
13    broadcast: Broadcast[JsonMapperProvider]
14  ): Unit = {
15    val developersDF =
16      DatabaseReader.read(
17        sparkSession,
18        applicationConfig.databaseConfig.connectionConfig,
19        applicationConfig.databaseConfig.tables.developer
20      )
21    developersDF.persist(StorageLevel.MEMORY_ONLY)
22
23    val transformedSteamApplicationDevelopersDF =
24      SteamApplicationDeveloperTransformer.transform(steamAppsDetailsDF, developersDF, broadcast)
25    val transformedSteamSpyApplicationDevelopersDF =
26      SteamSpyApplicationDeveloperTransformer.transform(steamSpyAppsDetailsWithNoDataInSteamDF, developersDF)
27    val transformedDevelopersDF =
28      transformedSteamApplicationDevelopersDF.unionByName(transformedSteamSpyApplicationDevelopersDF)
29
30    DatabaseWriter.write(
31      transformedDevelopersDF,
32      applicationConfig.databaseConfig.connectionConfig,
33      applicationConfig.databaseConfig.tables.applicationDeveloper
34    )
35
36    developersDF.unpersist()
37  }
38 }
39

```

Рисунок 4.3.7. Реалізація об'єкту `ApplicationDeveloperKeeper`

Описаний вище процес отримання даних для їх збереження в одній з таблиць бази даних є спільним для усіх таблиць. Для кожної з них існують об'єкти із суфіксом Keeper та Transformer, а також використовуються об'єкти DatabaseReader та DatabaseWriter для читання та запису даних у базу даних відповідно. Такий підхід забезпечує легке розуміння зв'язків між об'єктами в застосунку і зручність роботи із ним, оскільки завжди зрозуміло, де потрібно зробити зміни, якщо виникає така потреба. Файлову структуру проекту застосунку можна побачити у *Додатку Є*.

Розділ 5. Аналіз виконаної роботи

5.1 Можливості розвитку розробленої системи

Розроблена система уже складається із двох незалежних частин, а тому може стати прототипом для більш функціональної системи, що буде отримувати, обробляти та надавати дані, отримані від сервісів Steam та Steam Spy, у реальному часі. Це вимагатиме досить частих звертань до сторонніх API для перевірки оновлення даних, більшої потужності для обробки більшої кількості даних, а також кращого рішення для комунікації між частинами системи, ніж проміжні файли. Для виконання цих вимог можна застосувати:

- Проксі-сервери, які дозволять звертатися до сторонніх сервісів паралельно та уникнути обмежень на кількість запитів. Це значно пришвидшить процес отримання даних від сторонніх API.
- Брокери повідомлень, що забезпечать комунікацію між частинами системи в режимі реального часу. Із одним файлом одночасно може взаємодіяти лише один процес операційної системи, тому на їх основі не вийде побудувати систему для комунікації між процесами в реальному часі.
- Створення повноцінного Spark-кластеру, що дозволить задіяти більше ресурсів для обробки більшої кількості даних.

Впровадження наведених вище рішень дозволить розробити систему, яка дозволить отримувати найбільш актуальні дані від сторонніх API у реальному часі.

5.2 Висновки

В ході виконання роботи було розроблено систему, яка виконує повний процес отримання, обробки та зберігання даних, отриманих зі сторонніх джерел. Було розроблено реляційну модель предметної області на основі доступних даних.

Розроблена реляційна модель надає основну інформацію для опису сутностей предметної області і може бути розширена при наявності нових джерел інформації.

При розробці застосунку, який відповідає за отримання даних, були розглянуті основні проблеми, які можуть виникнути при взаємодії зі сторонніми публічними API та основні способи їх вирішення. Також було описано особливості роботи із неструктурованими даними, отриманими зі сторонніх джерел. Важливою частиною розробки було використання сторонніх інструментів для комунікації зі сторонніми API та перетворення даних між різними форматами представлення. Використання таких інструментів дозволяє уникнути необхідності власноруч створювати рішення для стандартних проблем і зосередитись на виконанні задач, що стоять перед застосунком.

При розробці застосунку для обробки та збереження даних було розглянуто основні можливості фреймворку Apache Spark. Вони були застосовані для вирішення задач обробки великої кількості неструктурованих даних, отриманих із різних джерел. Використання Apache Spark забезпечило зручність розробки, велику кількість інструментів для виконання обробки даних, можливість для гнучкого розширення стандартного функціоналу Spark а також високу ефективність роботи кінцевого застосунку.

Було розглянуто можливості подальшого розвитку та розширення розробленої системи, які дозволять забезпечити її роботу із даними, отриманими у реальному часі.

Дані, отримані в результаті роботи системи, можуть бути використані для проведення аналізу стану ринку відеогор на момент їх збору.

Список використаних джерел

1. Офіційний сайт сервісу Steam [Електронний ресурс] – Режим доступу до ресурсу: <https://store.steampowered.com/about>.
2. Офіційний сайт сервісу Steam Spy [Електронний ресурс] – Режим доступу до ресурсу: <https://steamspy.com/about>.
3. Документація СКБД MySQL [Електронний ресурс] – Режим доступу до ресурсу: <https://dev.mysql.com/doc/refman/8.0/en/>.
4. Документація Docker [Електронний ресурс] – Режим доступу до ресурсу: <https://docs.docker.com/>.
5. Документація Spring WebFlux [Електронний ресурс] – Режим доступу до ресурсу: <https://docs.spring.io/spring-framework/docs/current/reference/html/web-reactive.html>.
6. Walls C. Spring in Action / Crain Walls., 2018. – 520 с.
7. Документація Steamworks API [Електронний ресурс] – Режим доступу до ресурсу: <https://partner.steamgames.com/doc/webapi>.
8. Документація Steam API [Електронний ресурс] – Режим доступу до ресурсу: <https://steamapi.xpaw.me/>.
9. Official Team Fortress Wiki [Електронний ресурс] – Режим доступу до ресурсу: <https://wiki.teamfortress.com/wiki/User:RJackson/StorefrontAPI>.
10. Документація Steam Spy API [Електронний ресурс] – Режим доступу до ресурсу: <http://steamspy.com/api.php>.
11. Домашня сторінка бібліотеки Jackson на GitHub [Електронний ресурс] – Режим доступу до ресурсу: <https://github.com/FasterXML/Jackson>.
12. Документація Apache Commons CSV [Електронний ресурс] – Режим доступу до ресурсу: <https://commons.apache.org/proper/commons-csv/index.html>.
13. Документація Apache Spark [Електронний ресурс] – Режим доступу до ресурсу: <https://spark.apache.org/docs/latest/>.

14. Perrin J. Spark in Action, Second Edition / Jean-Georges Perrin., 2020. – 576 с.
15. Spark by Examples - Spark Broadcast Variables [Электронный ресурс] – Режим доступа до ресурсу: <https://sparkbyexamples.com/spark/spark-broadcast-variables/>.

Додатки

Додаток А. Клас WebClientConfiguration

```

1 package ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.configuration;
2
3 import io.netty.channel.ChannelOption;
4 import io.netty.handler.timeout.ReadTimeoutHandler;
5 import java.util.concurrent.TimeUnit;
6 import lombok.RequiredArgsConstructor;
7 import org.springframework.context.annotation.Bean;
8 import org.springframework.context.annotation.Configuration;
9 import org.springframework.http.client.reactive.ReactorClientHttpConnector;
10 import org.springframework.web.reactive.function.client.ExchangeStrategies;
11 import org.springframework.web.reactive.function.client.WebClient;
12 import reactor.netty.http.client.HttpClient;
13
14 @Configuration
15 @RequiredArgsConstructor
16 public class WebClientConfiguration {
17
18     3 usages
19     private final WebClientProperties webClientProperties;
20
21     @Bean
22     public WebClient steamWebClient() {
23         return WebClient.builder()
24             .clientConnector(new ReactorClientHttpConnector(createHttpClient()))
25             .exchangeStrategies(
26                 ExchangeStrategies.builder()
27                     .codecs(
28                         codecs ->
29                             codecs
30                                 .defaultCodecs()
31                                 .maxInMemorySize(webClientProperties.getClient().getMaxBufferSize())
32                             .build()
33                     )
34             .build();
35     }
36
37     @Bean
38     public WebClient steamSpyWebClient() {
39         return WebClient.builder()
40             .clientConnector(new ReactorClientHttpConnector(createHttpClient()))
41             .build();
42     }
43
44     2 usages
45     @private HttpClient createHttpClient() {
46         return HttpClient.create()
47             .option(
48                 ChannelOption.CONNECT_TIMEOUT_MILLIS,
49                 Long.valueOf(webClientProperties.getClient().getConnectTimeout().toMillis()).intValue()
50             )
51             .doOnConnected(
52                 connection ->
53                     connection.addHandlerFirst(
54                         new ReadTimeoutHandler(
55                             webClientProperties.getClient().getReadTimeout().toMillis(),
56                             TimeUnit.MILLISECONDS));
57             )
58     }
59 }

```

Додаток Б. Клас SteamService

```
SteamService.java
1 package ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.service.steam;
2
3 import lombok.RequiredArgsConstructor;
4 import org.springframework.stereotype.Service;
5 import org.springframework.web.reactive.function.client.WebClient;
6 import reactor.core.publisher.Mono;
7 import reactor.util.retry.Retry;
8 import ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.configuration.WebClientProperties;
9 import ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.service.steam.response.SteamAllAppsResponse;
10 import ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.service.steam.response.SteamAppDetailsResponse;
11
12 @Service
13 @RequiredArgsConstructor
14 public class SteamService {
15
16     private final SteamProperties steamProperties;
17     private final WebClientProperties webClientProperties;
18     private final WebClient steamWebClient;
19
20     private static final String STEAM_APP_DETAILS_ID_PARAM = "appid";
21     private static final String STEAM_APP_DETAILS_LANGUAGE_PARAM = "l";
22     private static final String STEAM_APP_DETAILS_LANGUAGE_ENGLISH = "english";
23     private static final String STEAM_APP_DETAILS_CURRENCY_PARAM = "cc";
24     private static final String STEAM_APP_DETAILS_CURRENCY_DOLLAR = "US";
25
26     public Mono<SteamAllAppsResponse> getAllSteamApps() {
27         return steamWebClient
28             .get()
29             .uri(steamProperties.getAllAppsUrl())
30             .retrieve()
31             .bodyToMono(SteamAllAppsResponse.class);
32     }
33
34     public Mono<SteamAppDetailsResponse> getSteamAppDetails(Long steamAppId) {
35         return steamWebClient
36             .get()
37             .uri(
38                 steamProperties.getAppDetailsUrl(),
39                 uriBuilder ->
40                     uriBuilder
41                         .queryParam(STEAM_APP_DETAILS_ID_PARAM, steamAppId)
42                         .queryParam(
43                             STEAM_APP_DETAILS_LANGUAGE_PARAM, STEAM_APP_DETAILS_LANGUAGE_ENGLISH)
44                         .queryParam(STEAM_APP_DETAILS_CURRENCY_PARAM, STEAM_APP_DETAILS_CURRENCY_DOLLAR)
45                         .build()
46             )
47             .retrieve()
48             .bodyToMono(SteamAppDetailsResponse.class)
49             .retryWhen(
50                 Retry.fixedDelay(
51                     webClientProperties.getRetry().getAmount(),
52                     webClientProperties.getRetry().getDelay());
53             );
54     }
55 }
```


Додаток Г. Клас SteamDataCsvFileWriter

```
1 package ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.csv.writer;
2
3 import com.fasterxml.jackson.core.JsonProcessingException;
4 import com.fasterxml.jackson.databind.ObjectMapper;
5 import java.io.*;
6 import java.util.Objects;
7 import java.util.Optional;
8 import org.apache.commons.csv.CSVPrinter;
9 import org.apache.commons.lang3.tuple.Pair;
10 import ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.csv.SteamDataProperties;
11 import ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.service.steam.response.SteamAppDetailsResponse;
12
13
14
15
16
17
18
19
20
21 @Usage
22 public class SteamDataCsvFileWriter extends AbstractCsvFileWriter<SteamAppDetailsResponse> {
23
24     @Usage
25     private final CSVPrinter steamAppDetailsCsvPrinter;
26
27     @Usage
28     private final CSVPrinter steamAppDetailsNotParsedCsvPrinter;
29
30     @Usage
31     private final CSVPrinter steamAppDetailsNoDataCsvPrinter;
32
33     @Usage
34     private final CSVPrinter steamProcessedAppsIdsCsvPrinter;
35
36     @Usage
37     private final ObjectMapper objectMapper;
38
39
40     @Usage
41     public SteamDataCsvFileWriter(final SteamDataProperties steamDataProperties) {
42         this.steamAppDetailsCsvPrinter = createCSVPrinter(steamDataProperties.getSteamAllAppsDetails());
43         this.steamAppDetailsNotParsedCsvPrinter =
44             createCSVPrinter(steamDataProperties.getSteamAllAppsDetailsNotParsed());
45         this.steamAppDetailsNoDataCsvPrinter =
46             createCSVPrinter(steamDataProperties.getSteamAllAppsDetailsNoData());
47         this.steamProcessedAppsIdsCsvPrinter =
48             createCSVPrinter(steamDataProperties.getSteamProcessedAppsIds());
49
50         this.objectMapper = new ObjectMapper();
51     }
52
53 }
```

```
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```

```

72 Optional.ofNullable(appDetails.getPcRequirements()) Optional<PcRequirements>
73     .map(SteamAppDetailsResponse.PcRequirements::getMinimum) Optional<String>
74     .orElse( other: null),
75 Optional.ofNullable(appDetails.getPcRequirements()) Optional<PcRequirements>
76     .map(SteamAppDetailsResponse.PcRequirements::getRecommended) Optional<String>
77     .orElse( other: null),
78 Optional.ofNullable(appDetails.getMacRequirements()) Optional<PcRequirements>
79     .map(SteamAppDetailsResponse.PcRequirements::getMinimum) Optional<String>
80     .orElse( other: null),
81 Optional.ofNullable(appDetails.getMacRequirements()) Optional<PcRequirements>
82     .map(SteamAppDetailsResponse.PcRequirements::getRecommended) Optional<String>
83     .orElse( other: null),
84 Optional.ofNullable(appDetails.getLinuxRequirements()) Optional<PcRequirements>
85     .map(SteamAppDetailsResponse.PcRequirements::getMinimum) Optional<String>
86     .orElse( other: null),
87 Optional.ofNullable(appDetails.getLinuxRequirements()) Optional<PcRequirements>
88     .map(SteamAppDetailsResponse.PcRequirements::getRecommended) Optional<String>
89     .orElse( other: null),
90 toJsonArray(appDetails.getDevelopers()),
91 toJsonArray(appDetails.getPublishers()),
92 Optional.ofNullable(appDetails.getPriceOverview()) Optional<PriceOverview>
93     .map(SteamAppDetailsResponse.PriceOverview::getCurrency) Optional<String>
94     .orElse( other: null),
95 Optional.ofNullable(appDetails.getPriceOverview()) Optional<PriceOverview>
96     .map(SteamAppDetailsResponse.PriceOverview::getInitialPrice) Optional<Long>
97     .orElse( other: null),
98 Optional.ofNullable(appDetails.getPriceOverview()) Optional<PriceOverview>
99     .map(SteamAppDetailsResponse.PriceOverview::getFinalPrice) Optional<Long>
100     .orElse( other: null),
101 Optional.ofNullable(appDetails.getPriceOverview()) Optional<PriceOverview>
102     .map(SteamAppDetailsResponse.PriceOverview::getDiscountPercent) Optional<Integer>
103     .orElse( other: null),
104 Optional.ofNullable(appDetails.getPriceOverview()) Optional<PriceOverview>
105     .map(SteamAppDetailsResponse.PriceOverview::getInitialFormatted) Optional<String>
106     .orElse( other: null),
107 Optional.ofNullable(appDetails.getPriceOverview()) Optional<PriceOverview>
108     .map(SteamAppDetailsResponse.PriceOverview::getFinalFormatted) Optional<String>
109     .orElse( other: null),
110 toJsonArray(appDetails.getPackages()), // array
111 toJsonArray(appDetails.getPackageGroups()), // array

```

```

112 Optional.ofNullable(appDetails.getPlatforms()) Optional<Platforms>
113     .map(SteamAppDetailsResponse.Platforms::getWindows) Optional<Boolean>
114     .orElse( other: null),
115 Optional.ofNullable(appDetails.getPlatforms()) Optional<Platforms>
116     .map(SteamAppDetailsResponse.Platforms::getMac) Optional<Boolean>
117     .orElse( other: null),
118 Optional.ofNullable(appDetails.getPlatforms()) Optional<Platforms>
119     .map(SteamAppDetailsResponse.Platforms::getLinux) Optional<Boolean>
120     .orElse( other: null),
121 Optional.ofNullable(appDetails.getMetacritic()) Optional<Metacritic>
122     .map(SteamAppDetailsResponse.Metacritic::getScore) Optional<Integer>
123     .orElse( other: null),
124 Optional.ofNullable(appDetails.getMetacritic()) Optional<Metacritic>
125     .map(SteamAppDetailsResponse.Metacritic::getURL) Optional<String>
126     .orElse( other: null),
127 toJsonArray(appDetails.getCategories()), // array
128 toJsonArray(appDetails.getGenres()), // array
129 toJsonArray(appDetails.getScreenshots()), // array
130 Optional.ofNullable(appDetails.getRecommendations()) Optional<Recommendations>
131     .map(SteamAppDetailsResponse.Recommendations::getTotal) Optional<Long>
132     .orElse( other: null),
133 Optional.ofNullable(appDetails.getReleaseDate()) Optional<ReleaseDate>
134     .map(SteamAppDetailsResponse.ReleaseDate::getComingSoon) Optional<Boolean>
135     .orElse( other: null),
136 Optional.ofNullable(appDetails.getReleaseDate()) Optional<ReleaseDate>
137     .map(SteamAppDetailsResponse.ReleaseDate::getDate) Optional<String>
138     .orElse( other: null),
139 Optional.ofNullable(appDetails.getSupportInfo()) Optional<SupportInfo>
140     .map(SteamAppDetailsResponse.SupportInfo::getURL) Optional<String>
141     .orElse( other: null),
142 Optional.ofNullable(appDetails.getSupportInfo()) Optional<SupportInfo>
143     .map(SteamAppDetailsResponse.SupportInfo::getEmail) Optional<String>
144     .orElse( other: null),
145 appDetails.getBackground(),
146 appDetails.getBackgroundRaw(),
147 toJsonArray(
148     Optional.ofNullable(appDetails.getContentDescriptors()) Optional<ContentDescriptors>
149         .map(SteamAppDetailsResponse.ContentDescriptors::getIds) Optional<Long[]>
150         .orElse( other: null), // array
151 Optional.ofNullable(appDetails.getContentDescriptors()) Optional<ContentDescriptors>
152     .map(SteamAppDetailsResponse.ContentDescriptors::getNotes) Optional<String>

```

```

153         .orElse( other: null));
154     } catch (IOException e) {
155         throw new RuntimeException(e);
156     }
157 }
158
159 @Usage
160 private String toJsonArray(Object[] array) {
161     return Optional.ofNullable(array).map(this::writeArrayAsJsonString).orElse( other: null);
162 }
163
164 @Usage
165 private String writeArrayAsJsonString(Object[] array) {
166     try {
167         return objectMapper.writeValueAsString(array);
168     } catch (JsonProcessingException e) {
169         throw new RuntimeException(e);
170     }
171 }
172
173 @Usage
174 private void writeNoDataSteamAppDetails(Pair<Long, SteamAppDetailsResponse> pair) {
175     System.out.printf(
176         "Saving steamAppDetails with no data for steamAppId=[%s]...\n", pair.getLeft());
177     try {
178         steamAppDetailsNoDataCsvPrinter.printRecord(pair.getLeft());
179     } catch (IOException e) {
180         throw new RuntimeException(e);
181     }
182 }

```

```

183 @Usage
184 private void writeNotParsedSteamAppDetails(Pair<Long, SteamAppDetailsResponse> pair) {
185     System.out.printf("Saving not parsed steamAppDetails for steamAppId=[%s]...\n", pair.getLeft());
186     try {
187         steamAppDetailsNotParsedCsvPrinter.printRecord(
188             pair.getLeft(), pair.getRight().getResponseData().getErrorData().getMessage());
189     } catch (IOException e) {
190         throw new RuntimeException(e);
191     }
192 }
193
194 @Usage
195 private void writeProcessedSteamAppId(Pair<Long, SteamAppDetailsResponse> pair) {
196     System.out.printf("Saving processed steamAppId=[%s]...\n", pair.getLeft());
197     try {
198         steamProcessedAppsIdsCsvPrinter.printRecord(pair.getLeft());
199     } catch (IOException e) {
200         throw new RuntimeException(e);
201     }
202 }

```

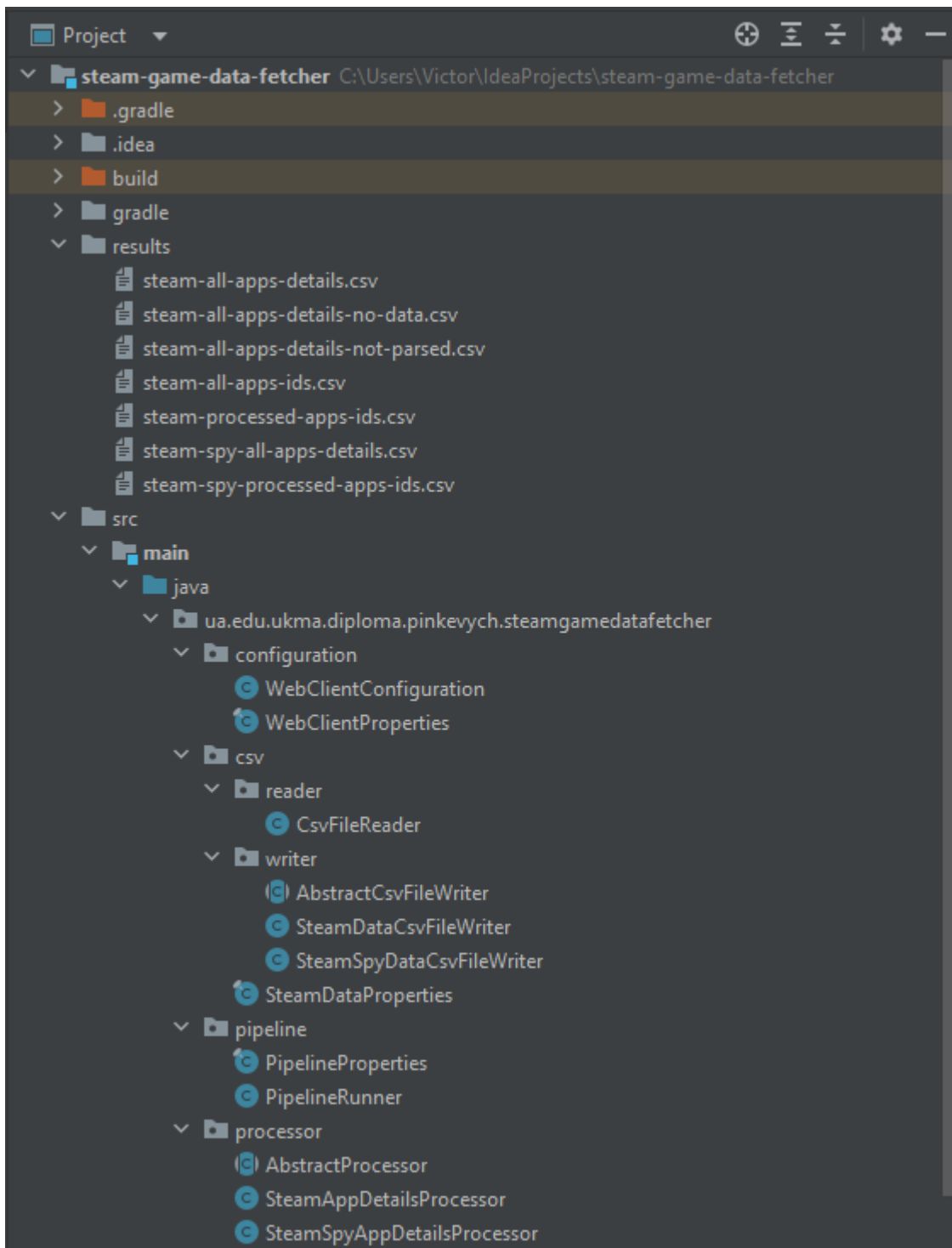
Додаток Д. Клас SteamAppDetailsProcessor

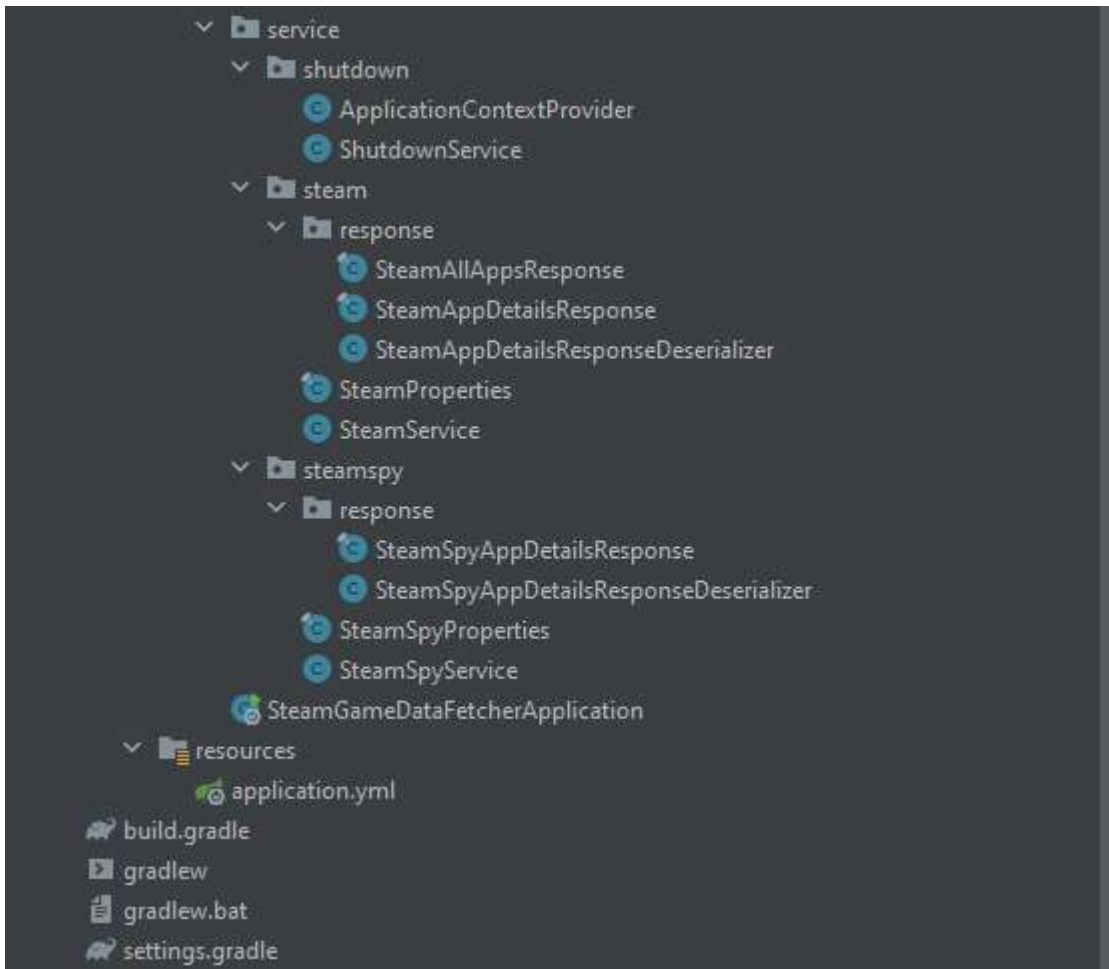
```

1 package ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.processor;
2
3 import java.time.Duration;
4 import org.apache.commons.lang3.tuple.Pair;
5 import org.reactivestreams.Publisher;
6 import org.springframework.stereotype.Component;
7 import reactor.core.publisher.Flux;
8 import ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.csv.SteamDataProperties;
9 import ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.csv.writer.AbstractCsvFileWriter;
10 import ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.csv.writer.SteamDataCsvFileWriter;
11 import ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.service.steam.SteamService;
12 import ua.edu.ukma.diploma.pinkevych.steamgamedatafetcher.service.steam.response.SteamAppDetailsResponse;
13
14 @Component
15 public class SteamAppDetailsProcessor extends AbstractProcessor<SteamAppDetailsResponse> {
16
17     1 usage
18     private static final Long STEAM_API_REQUEST_DELAY_MILLS = 1000L;
19
20     2 usages
21     private final SteamService steamService;
22
23     2 usages
24     private final AbstractCsvFileWriter<SteamAppDetailsResponse> csvFileWriter;
25
26     public SteamAppDetailsProcessor(
27         SteamService steamService, SteamDataProperties steamDataProperties) {
28         this.steamService = steamService;
29         this.csvFileWriter = new SteamDataCsvFileWriter(steamDataProperties);
30     }
31
32     1 usage
33     @Override
34     protected Publisher<Pair<Long, SteamAppDetailsResponse>> doProcess(
35         Flux<Long> appIds, AbstractCsvFileWriter<SteamAppDetailsResponse> csvFileWriter) {
36         return appIds
37             .delayElements(Duration.ofMillis(STEAM_API_REQUEST_DELAY_MILLS)) Flux<Long>
38             .flatMap(
39                 appId ->
40                 steamService
41                     .getSteamAppDetails(appId)
42                     .map(steamAppDetailsResponse -> Pair.of(appId, steamAppDetailsResponse)) Flux<Pair<Long, SteamAppDetailsResponse>>
43             ).doOnNext(csvFileWriter::write);
44     }
45
46     1 usage
47     @Override
48     protected AbstractCsvFileWriter<SteamAppDetailsResponse> getCsvFileWriter() { return csvFileWriter; }
49 }

```

Додаток Е. Файлова структура проекту застосунку для збору даних





Додаток Є. Файлова структура проекту застосунку для обробки даних

