



UNIVERSITÀ DI PISA



Master's Degree in Artificial Intelligence and Data Engineering

Design and develop of an Application interacting with
NoSQL Databases:

Smart News Aggregator

Instructors:

Prof. Pietro Ducange

Prof. Alessio Schiavo

Student:

Adrien Koumgang

Tegantchouang

Academic Year 2024/2025

Contents

1	Introduction	1
1.1	Overview	1
1.2	Objectives	1
1.3	Structure of the presentation	2
1.4	Repository	3
2	Requirements Analysis	4
2.1	Problem Statement	4
2.2	Stakeholders and Actors	4
2.2.1	Stakeholders	4
2.2.2	Actors	5
2.3	Functional Requirements	5
2.3.1	User Account Management	6
2.3.2	News Aggregation & Display	6
2.3.3	Personalization & Recommendation	6
2.3.4	Admin & Analytics	6
2.3.5	System & Logging	7
2.4	Non-Functional Requirements	7
2.4.1	Performance	7
2.4.2	Scalability	7
2.4.3	Availability	7
2.4.4	Security	8
2.4.5	Maintainability	8
2.4.6	Usability	8
2.4.7	Portability	8

2.5	Use Case Diagrams	8
3	System Design	12
3.1	Architecture Overview	12
3.1.1	Key Components	12
3.2	System Architecture Diagram	13
3.2.1	Client Layer (Frontend)	13
3.2.2	Application Layer (Flask Backend)	14
3.2.3	Data Layer	14
3.2.4	External APIs Layer	14
3.3	UML Class Diagram	15
3.4	UI Wireframes and Mockups	15
3.4.1	Authentication Screens	15
3.4.2	User Screens	16
3.4.3	Admin Screens	17
4	Data Modeling	21
4.1	Databases Technologies Used	21
4.2	MongoDB Data Models	21
4.2.1	Users collection	21
4.2.2	Article log requests collection	23
4.2.3	Articles collection	24
4.2.4	Comments collection	25
4.2.5	User Interactions collection	26
4.2.6	Server Error Logs collection	27
4.3	Redis Data Structures	28
4.3.1	Article Caching	28
4.3.2	User caching	29
5	Data Ingestion and Integration	30
5.1	External API Integration	30
5.2	Data Processing	31
5.2.1	Data Validation	31
5.2.2	Data Normalization	31

5.3	Error Handling and Logging	31
6	Backend Implementation	32
6.1	System Architecture	32
6.1.1	Key Components	32
6.2	API Endpoints	33
6.2.1	Authentication	33
6.2.2	User Management	33
6.2.3	Article Management	35
6.2.4	Administration Management	36
6.3	Security Implementation	37
6.3.1	JWT Authentication	37
6.3.2	Input Validation	37
6.4	Performance Optimizations	38
6.4.1	Caching Strategies	38
6.4.2	Async Task Processing	39
6.5	API Documentation	39
7	Advanced Aggregations and Analytics	41
7.1	MongoDB Aggregation Framework	41
7.1.1	Key Aggregation Concepts	41
7.1.2	Pipeline: All tags	41
7.1.3	Pipeline: Search Articles	42
7.1.4	Pipeline: Article Stats Comment	44
7.1.5	Pipeline User Comments with article	45
7.1.6	Pipeline User Article Interaction by article	47
7.1.7	Pipeline Most Interacted Articles	48
7.1.8	Pipeline User Preferences Most Used Tags	52
7.2	Performance Optimization	53
8	Deployment and Scaling	55
8.1	Containerization with Docker	55
8.1.1	Core Services Setup	55
8.1.2	Backend Setup	57

8.1.3	Frontend Setup	58
8.2	Testing and Monitoring	59
9	Conclusion	61
9.1	Project Achievements	61
9.2	Challenges Overcome	62
9.3	Future Work	62

Chapter 1

Introduction

1.1 Overview

In the digital age, the volume of news articles produced every day is staggering. Readers are increasingly overwhelmed by the amount of available content and often struggle to find reliable and relevant information that matches their personal interests. To address this challenge, we propose the development of a **Smart News Aggregator & Reader Personalization Platform**.

This platform collects news articles from multiple external APIs, stores and processes them using NoSQL databases (MongoDB and Redis), and delivers a personalized reading experience to users. It incorporates intelligent recommendation features, secure user authentication via JWT [1] tokens, and real-time analytics to improve user engagement.

1.2 Objectives

The primary objective of this project is to **design and implement a distributed news aggregation application** of intelligently retrieving, storing, analyzing, and serving large-scale multi-structured data using multiple NoSQL database technologies.

In line with the course requirements for **Large Scale and Multi-Structured Databases** [11], the specific goals of this project include:

- **Data Acquisition & Preprocessing:** Retrieve a real-world, high-volume dataset (50MB) from multiple external news APIs (e.g., MediaStack, newsData, The

Guardian, NYTimes), ensuring **variety** and **velocity**.

- **Distributed NoSQL Architecture:** Use **MongoDB** [6] as the primary **Document Database**, with carefully designed collections, indexes, and aggregation pipelines. Integrate a **Key-Value store (Redis [7])** to optimize caching and access to hot data.
- **System Design and Modeling:** Define functional and non-functional requirements. Design UML class diagrams and user interface mockups. Model the data schema for each NoSQL DB in use.
- **RESTful API Backend:** Build a scalable backend using Flask [9] and Flask-RESTX [3]. Expose secure endpoints with JWT-based authentication. Provide API documentation via Swagger and test interfaces.
- **Advanced Analytics:** Implement at least three real aggregation pipelines in MongoDB for summarizing and analyzing article content and user activity. Provide user-personalized views and filtering using advanced queries.
- **Monitoring, Testing, and Deployment:** Deploy on both AWS and UNIPY virtual clusters [5]. Include performance tests, system logging, and analysis of read/write throughput under different consistency models. Offer a complete presentation and demonstration, including API functionality and analytics insights.

By achieving these goals, the project demonstrates my ability to handle real-world large-scale data applications, integrating theoretical design with practical deployment, and ensuring cross-database consistency and performance.

1.3 Structure of the presentation

The documentation is structured as follows:

- **Chapter 1 - Introduction:** Presents the context, objectives, and structure of the project, outlining the motivations and academic scope.
- **Chapter 2 - Requirements Analysis:** Identifies the functional and non-functional requirements of the application, defines the system actors, and outlines key use cases.

- **Chapter 3 - System Design:** Includes the architectural overview, UML class diagrams, and mockup wireframes of the application's user interface.
- **Chapter 4 - Data Modeling:** Describes the design choices for MongoDB (Document DB) and Redis (Key-Value DB), along with schema definitions, key structures, and CAP theorem considerations.
- **Chapter 5 - Data Ingestion and Integration:** Explains how external APIs are connected, data is fetched and transformed, and stored into the databases. Also includes logging and handling of errors and API rate limits.
- **Chapter 6 - Backend Implementation:** Details the development of RESTful API endpoints using Flask and Flask-RESTX, JWT authentication, Swagger documentation, and the modular blueprint structure.
- **Chapter 7 - Advanced Aggregations and Analytics:** Presents non-trivial aggregation pipelines in MongoDB, user-specific personalization features, and analytics insights extracted from the dataset.
- **Chapter 8 - Deployment and Testing:** Discusses deployment on local and virtualized environments, performance testing scenarios, and consistency benchmarking across NoSQL systems.
- **Chapter 9 - Conclusion:** Summarizes achievements, reflects on encountered challenges, and suggests possible extensions to improve scalability, interactivity, and intelligence.

1.4 Repository

- Backend: Smart News Aggregator API
- Frontend: Smart News Aggregator Frontend React
- Documentation: Smart News Aggregator Documentation

Chapter 2

Requirements Analysis

2.1 Problem Statement

With the exponential growth of online content, users are overwhelmed by the volume of news available across platforms. This leads to difficulty in identifying relevant and trustworthy information. The proposed system aims to aggregate articles from various news APIs and personalize the reading experience based on user behavior and preferences, while maintaining scalability and high performance.

2.2 Stakeholders and Actors

The successful deployment and functioning of the Smart News Aggregator & Reader Personalization Platform depends on multiple stakeholders and actors who interact directly or indirectly with the system. This section outlines their roles and responsibilities.

2.2.1 Stakeholders

- **End Users:** Consume and interact with news content; expect personalized and relevant information.
- **Platform Administrator:** Oversees user activity, ensures data integrity, and manages access or moderation tasks.

- **Project Developers:** Build and maintain the system backend, frontend, and data processing pipelines.
- **External API Providers:** Supply news content (e.g., NYTimes, Guardian, News-Data, CurrentsAPI, MediaStack).
- **Academic Supervisors:** Oversee the project's architecture, correctness, and evaluate its educational objectives.

2.2.2 Actors

- **Visitor (Anonymous user):** Accesses the welcome page and Can register to become a user.
- **Registered User:** Logs into the platform, Personalizes preferences, Views news feed and Interacts with articles (like/comment).
- **Administrator:** Manages users and platform configurations and Monitors logs and analytics.
- **News Aggregator Service:** Fetches articles from external APIs and Normalizes and stores them into MongoDB.
- **External News APIs:** Provide raw article data in JSON format for ingestion by the system.

Each actor has specific actions and permissions that are further detailed in the Use Case and UML Diagrams.

2.3 Functional Requirements

The Smart News Aggregator & Reader Personalization Platform offers a suite of functionalities that serve both user-facing and administrative purposes. The following functional requirements have been identified:

2.3.1 User Account Management

- **Registration:** Users must be able to register with valid email and password.
- **Login/Logout:** Authenticated access via JWT-based login; logout clears session on frontend.
- **Profile Management:** Users can update personal data (e.g., name, preferences).
- **Password Handling:** Secure password storage and update with history tracking.

2.3.2 News Aggregation & Display

- **Fetch Articles from External APIs:** The system retrieves and normalizes article data from third-party providers such as NYTimes, NewsData, CurrentsAPI, etc.
- **Categorized News Display:** Articles are displayed by category (e.g., politics, tech, sports).
- **Search & Filter:** Users can search articles using keywords and filter by category, source, or date.
- **Pagination:** Article lists are paginated to handle large datasets.

2.3.3 Personalization & Recommendation

- **Save Preferences:** Users can set preferred categories, sources, or languages.
- **Personalized Feed:** Articles are filtered and ranked based on user preferences.
- **Like & Save Articles:** Users can like or save articles for future reading.
- **Commenting System:** Users can add comments and replies to articles.

2.3.4 Admin & Analytics

- **Dashboard Access:** Admin can view usage statistics and system logs.
- **User Management:** Admin can deactivate or delete user accounts.

- **API Health Monitoring:** Admin is notified of errors when external APIs fail.
- **Article Quality Control:** Admin can remove duplicated or malformed articles.

2.3.5 System & Logging

- **Authentication Event Logs:** Login/Logout attempts are store for audit.
- **Error Tracking:** Failed requests or errors are stored in MongoDB.

2.4 Non-Functional Requirements

This section outlines the quality attributes and constraints the Smart News Aggregator & Reader Personalization Platform must meet to ensure reliability, security, and scalability.

2.4.1 Performance

The system must handle a high number of concurrent users with minimal latency. Article feed pages must load within 1 second under normal network conditions. External API data must be cached using Redis to reduce response time and load.

2.4.2 Scalability

The backend architecture must support horizontal scaling to accommodate growing numbers of users and API integrations. MongoDB and Redis must be configured to handle large-scale document and key-value datasets efficiently.

2.4.3 Availability

The application must ensure high availability (99.9%) during user access hours. In the case of third-party API failure, the system should provide fallback responses or cached data when available.

2.4.4 Security

User authentication must use JWT with asymmetric encryption (RS256). Passwords must be securely hashed (using bcrypt) and stored with history to prevent reuse. All endpoints must enforce token validation for sensitive operations. Cross-Origin Resource Sharing (CORS) must be properly configured to restrict access to allowed domains.

2.4.5 Maintainability

The system must be modular, separating concerns by feature (auth, articles, users). Logging and exception handling must be centralized to simplify debugging and maintenance. Code must be written following best practices and TypeScript (frontend) and Python (backend) style guides.

2.4.6 Usability

The frontend must offer an intuitive user interface with minimal learning curve. Responsive design must ensure accessibility on desktops, tablets, and mobile devices.

2.4.7 Portability

The application must run on Unix-based systems and be container-ready (Docker-compatible). APIs must be documented using Swagger/OpenAPI to allow easy integration by external systems.

2.5 Use Case Diagrams

These diagrams were created using PlantUML.

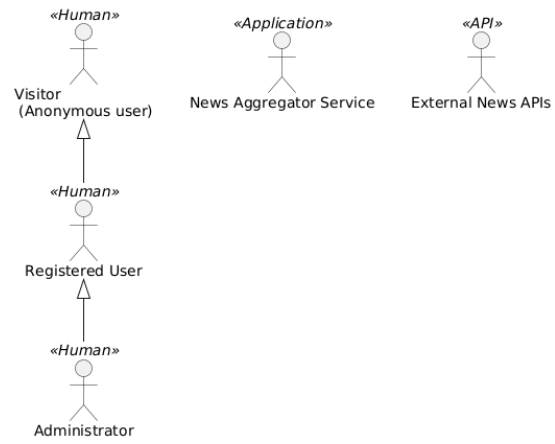


Figure 2.1: Use Case Diagram : Actors

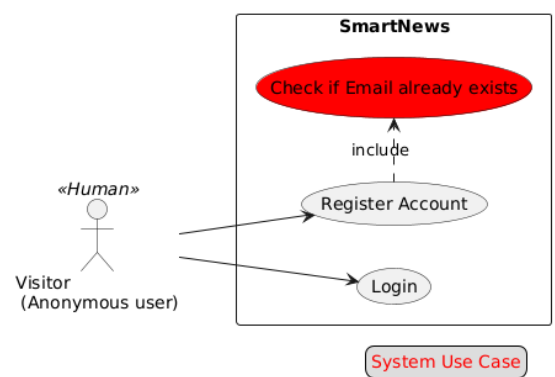


Figure 2.2: Use Case Diagram: Anonymous User

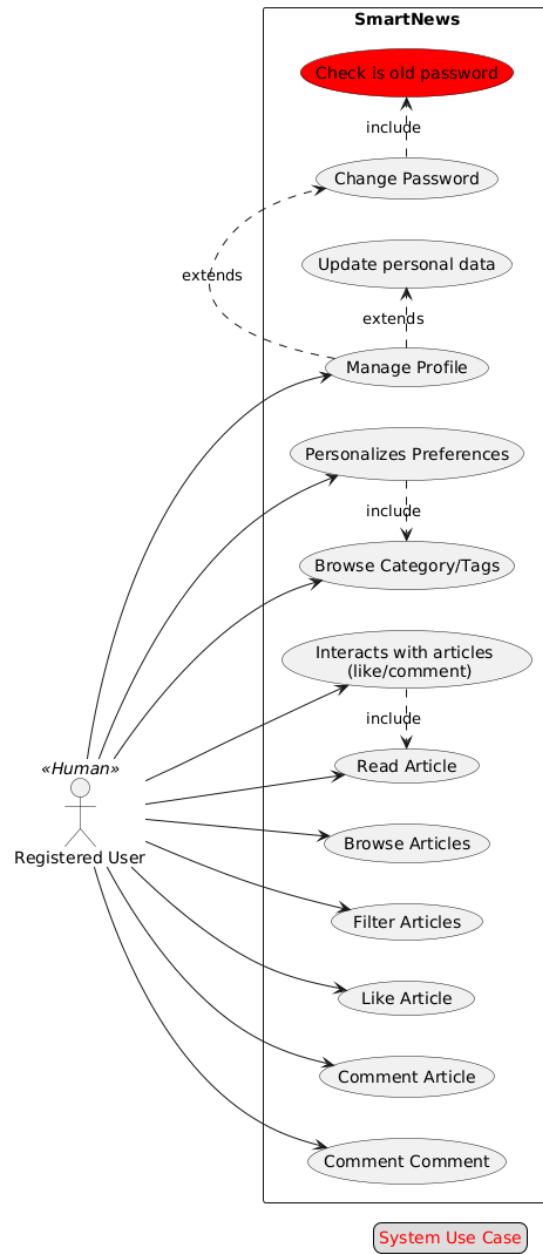


Figure 2.3: Use Case Diagram: Registered User

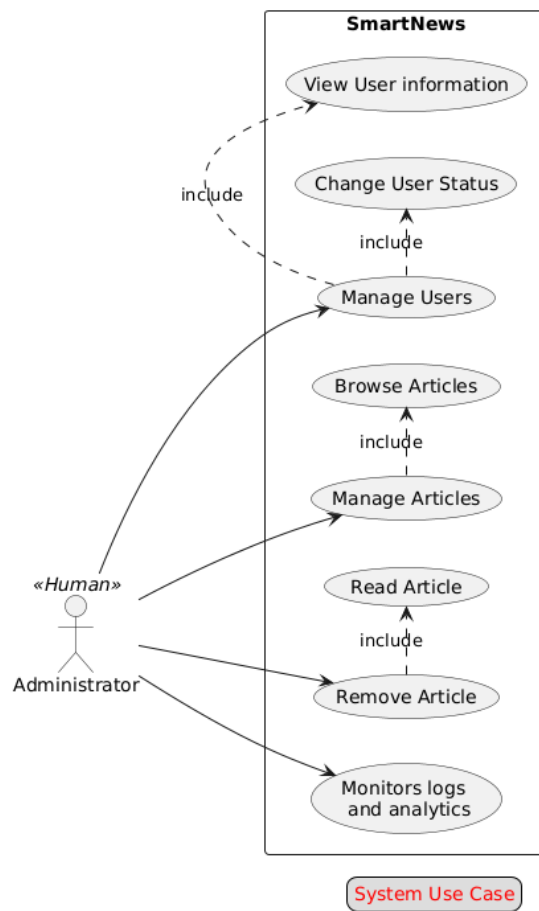


Figure 2.4: Use Case Diagram: Administrator

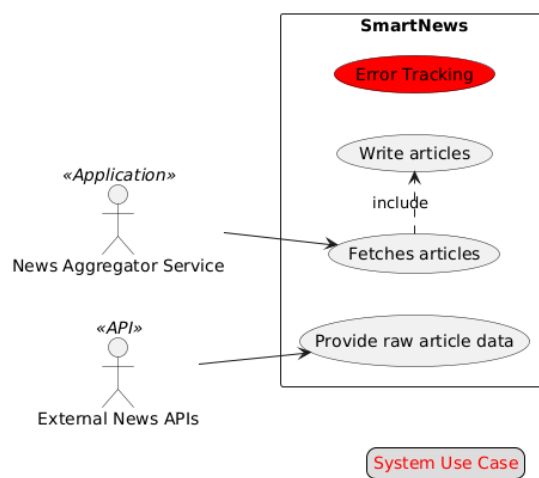


Figure 2.5: Use Case Diagram: System

Chapter 3

System Design

This chapter presents the architectural and structural design of the Smart News Aggregator & Reader Personalization Platform. It includes a high-level system architecture, database modeling strategy, UML class diagrams, and mockups of key user interfaces.

3.1 Architecture Overview

The platform is built using a **modular and layered architecture** to separate concerns and ensure maintainability, performance, and scalability. It integrates multiple technologies:

3.1.1 Key Components

- **Frontend (React [4,12] + TypeScript):** Handles UI, user interaction, API communication and Token-based authentication.
- **Backend (Flask + Flask-RESTX):** RESTful API with Blueprint organization, JWT authentication with RS256 and Logging and error handling.
- **Document Database (MongoDB):** Stores user profiles, articles, comments, and logs. Supports complex queries and aggregation pipelines.
- **Key-Value Store (Redis):** Caches trending articles and recent queries. Session tracking.

- **External APIs:** Article data sources such as CurrentsAPI, NewsData, NYTimes, Guardian.

3.2 System Architecture Diagram

The architecture is divided into four layers:

- **Client Layer:** Browser-based React frontend
- **Application Layer:** Flask server handling HTTP requests, routing, and validation
- **Data Layer:** MongoDB for structured content, Redis for quick key access
- **External Sources:** Third-party APIs for news ingestion

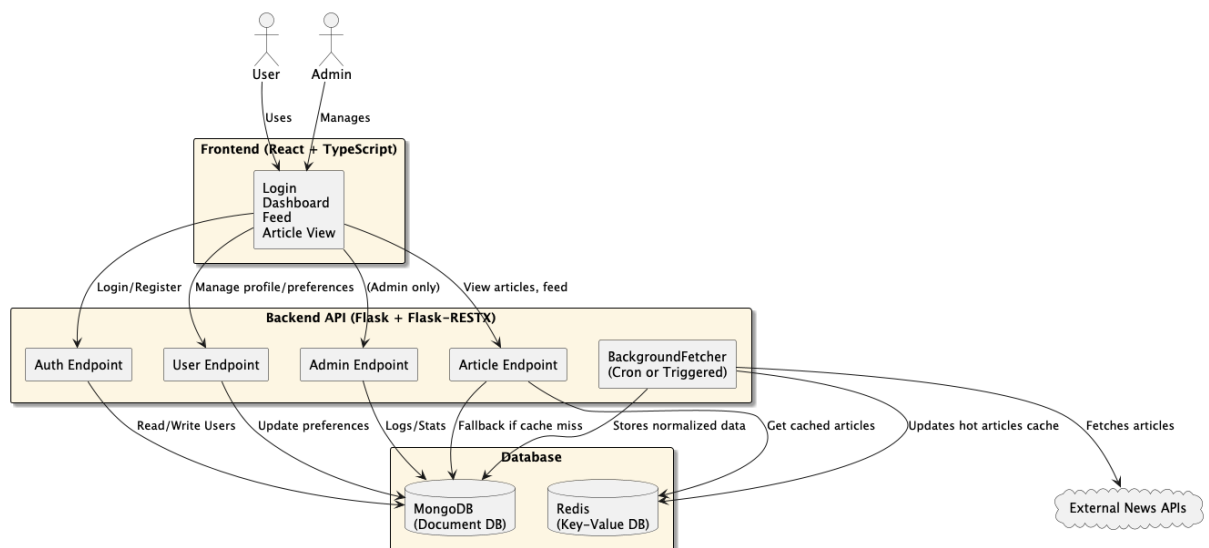


Figure 3.1: System Architecture

3.2.1 Client Layer (Frontend)

At the forefront is the **Client Layer**, which consists of a web-based frontend developed using React and TypeScript. This layer handles the user interface, manages routing, and communicates with the backend API through Axios. JWT tokens are stored in the browser's `localStorage` for session management. Users interact with interfaces such as login, registration, the personalized news feed, article detail views, and, for administrators, a dedicated dashboard.

3.2.2 Application Layer (Flask Backend)

The **Application Layer** is built on Flask with Flask-RESTX to expose RESTful endpoints in a modular architecture. This layer is responsible for authenticating users using asymmetric JWT (RS256), routing client requests to appropriate service handlers, and managing background tasks like scheduled fetching of articles. It includes modules like `auth_endpoint`, `user_endpoint`, `article_endpoint`, and `admin_endpoint`, as well as libraries for handling authentication keys and logging system events into MongoDB.

3.2.3 Data Layer

The **Data Layer** integrates two complementary NoSQL technologies. MongoDB, as a document-oriented database, is used to store structured collections such as users, articles, comments, preferences, and system logs. It supports advanced aggregation pipelines and enables efficient historical analysis. Redis, on the other hand, serves as a key-value store optimized for performance. It caches the latest articles, manages trending topics, and supports rate limiting. Redis ensures low-latency data access and reduces the load on MongoDB for repetitive queries.

3.2.4 External APIs Layer

At the periphery, the **External APIs Layer** includes third-party news services like MediaStack, CurrentsAPI, NewsData, the New York Times API, and the Guardian API. A background fetcher service, either scheduled or triggered, connects to these APIs, retrieves articles, normalizes the data, and saves the cleaned results into MongoDB. This ensures that the content database is continuously updated with fresh news.

The data flow follows a secure and performance-driven path. When a user logs in, they are authenticated via JWT, and the token is returned in the Authorization header. The frontend then retrieves the latest articles by calling the `/article/latest` endpoint. If the content is cached in Redis, it is served directly; otherwise, Flask queries MongoDB. Meanwhile, the background fetcher periodically updates the article repository. Personalized recommendations are generated by combining user preferences with

MongoDB aggregations and Redis-cached data. Administrators can monitor platform activity and user engagement via metrics obtained through aggregation queries stored in MongoDB.

3.3 UML Class Diagram

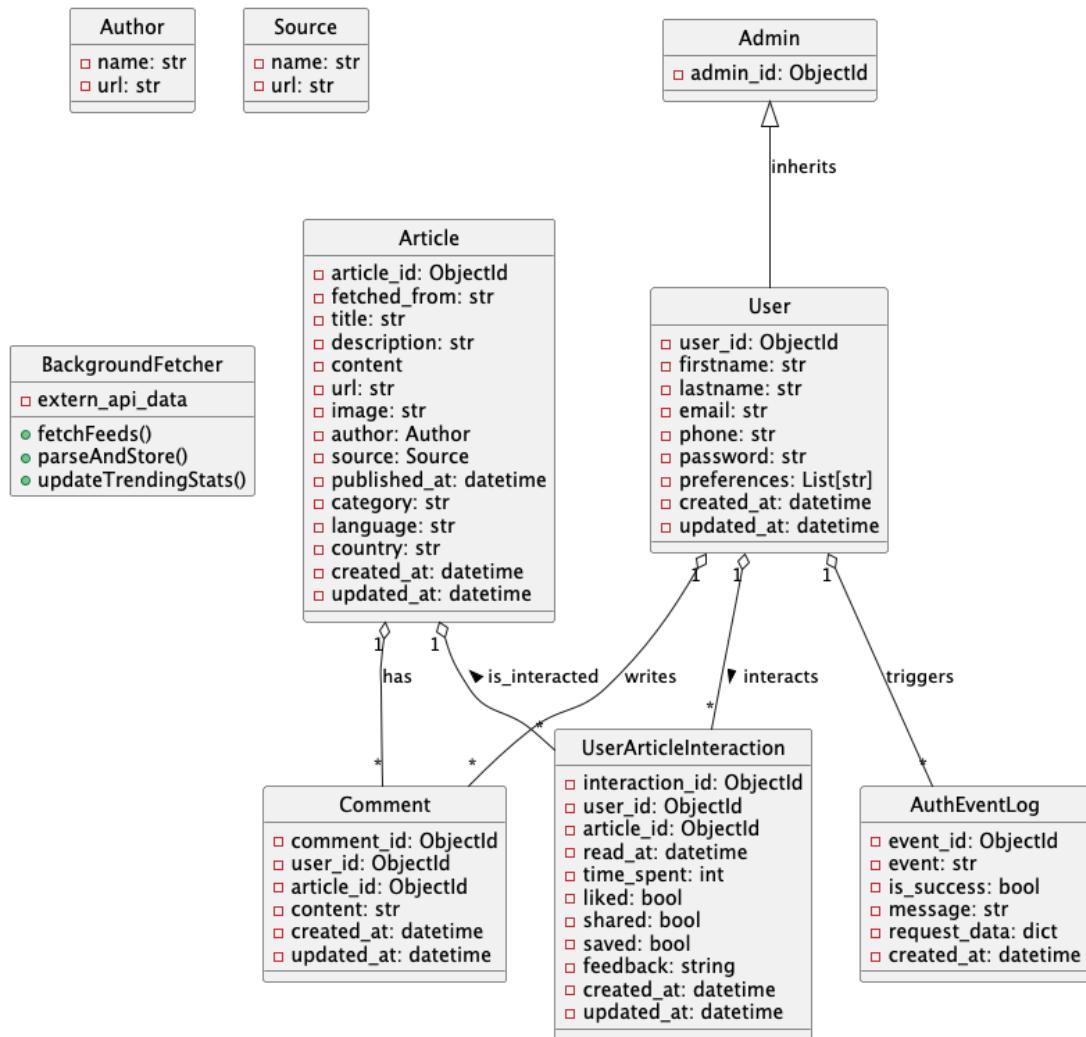
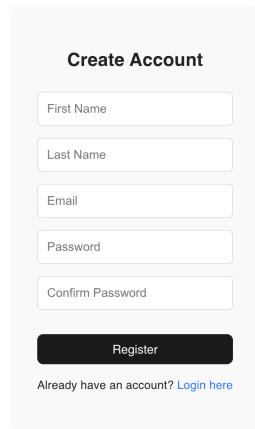


Figure 3.2: Use Case Diagram: System

3.4 UI Wireframes and Mockups

3.4.1 Authentication Screens

For registered new user and login:



Create Account

First Name

Last Name

Email

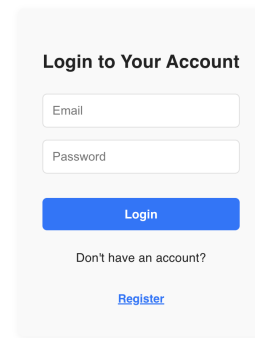
Password

Confirm Password

Register

Already have an account? [Login here](#)

Figure 3.3: Authentication screens: registration page



Login to Your Account

Email

Password

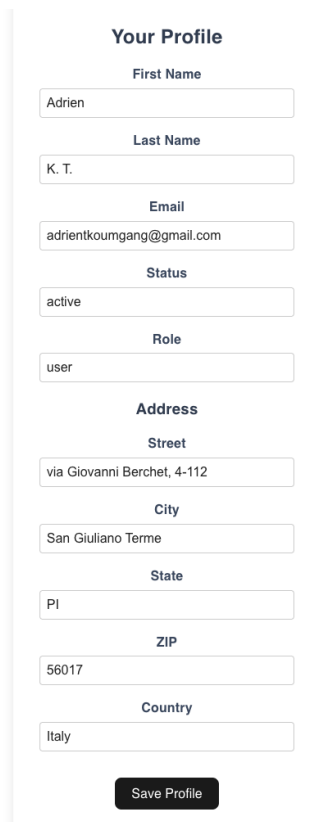
Login

Don't have an account? [Register](#)

Figure 3.4: Authentication screens: login page

3.4.2 User Screens

For see user all user information. It's not possible for user with role "user" to modify status and role.



Your Profile

First Name

Adrien

Last Name

K. T.

Email

adrientkoumgang@gmail.com

Status

active

Role

user

Address

Street

via Giovanni Berchet, 4-112

City

San Giuliano Terme

State

PI

ZIP

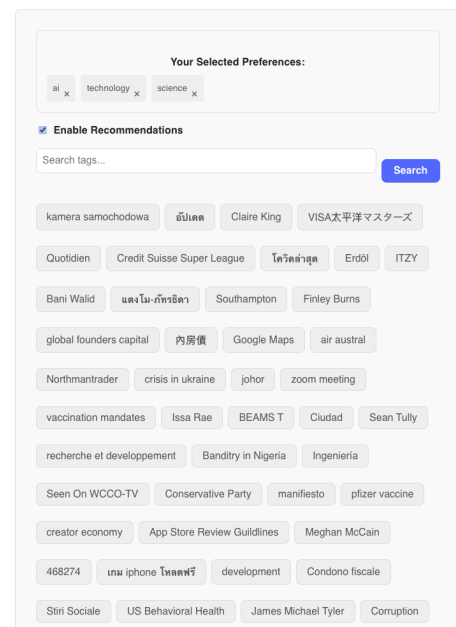
56017

Country

Italy

Save Profile

Figure 3.5: User screens: profile page



Your Selected Preferences:

ai x technology x science x

☒ Enable Recommendations

Search tags... [Search](#)

kamera samochodowa x หนังสื Claire King x VISA太平洋マスターズ

Quotidien x Credit Suisse Super League x โทริค่าสุท Endöl x ITZY

Bani Walid x แดงโม-ภักธิลา Southampton x Finley Burns

global founders capital x 内房債 Google Maps x air austral

Northmantrader x crisis in ukraine johor x zoom meeting

vaccination mandates x Issa Rae x BEAMS T Ciudad x Sean Tully

recherche et developpement x Banditry in Nigeria x Ingenieria

Seen On WCCO-TV x Conservative Party x manifesto x pfizer vaccine

creator economy x App Store Review Guidelines x Meghan McCain

468274 x เกม iphone โดคคพรี development x Condono fiscale

Stiri Sociale x US Behavioral Health x James Michael Tyler x Corruption

Figure 3.6: User screens: settings page

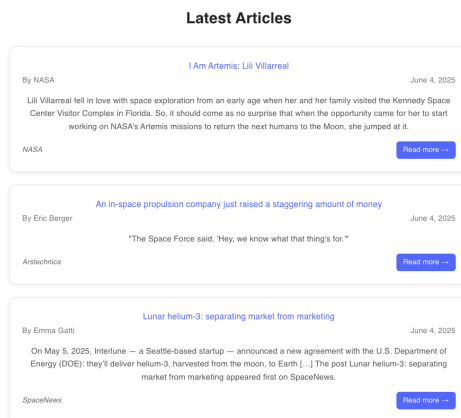


Figure 3.7: User screens: latest page

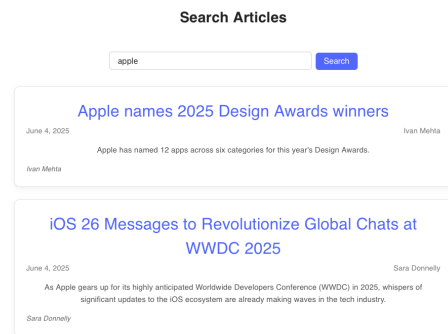


Figure 3.8: User screens: Search page

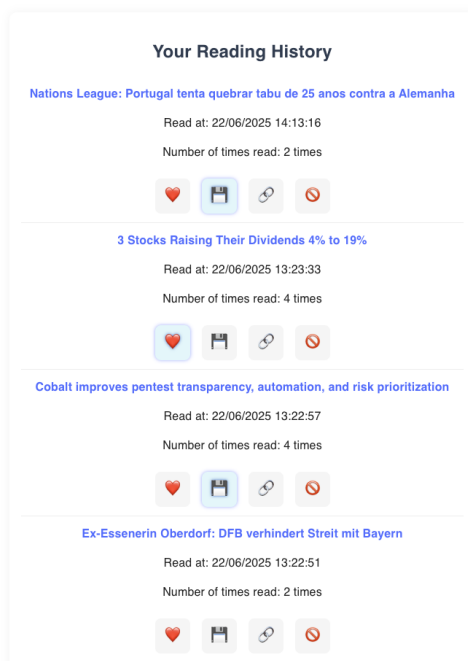


Figure 3.9: User screens: History page

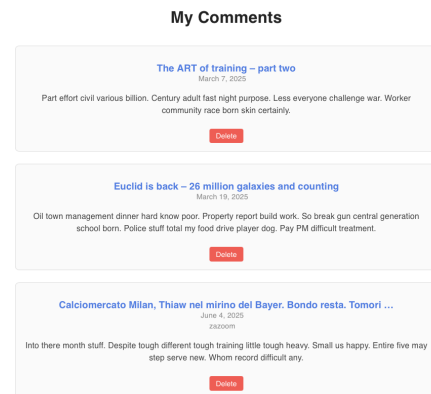


Figure 3.10: User screens: Comments page

3.4.3 Admin Screens

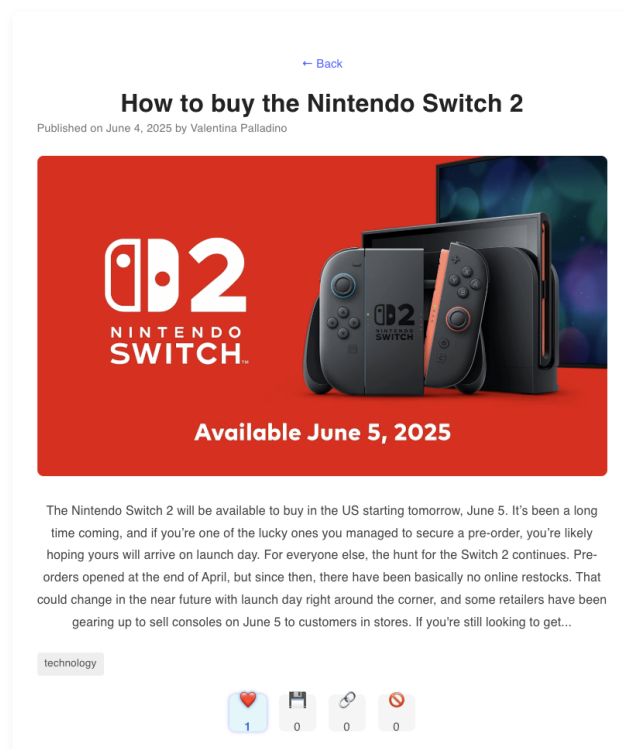


Figure 3.11: User screens: articles details page

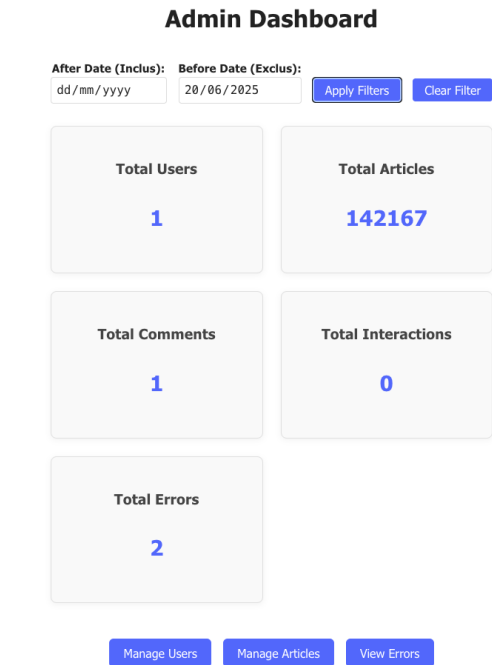


Figure 3.13: Admin screens: admin dash-
board page

Figure 3.12: Admin screens: admin dash-
board page

Admin – Article Management				
API	Title	Author	Published At	Actions
NewsData	Terveys I Suomalaiset saavat käyttöönsä vain puolet uusista lääkkeistä	Helsingin Sanomat	6/3/2025, 7:30:00 AM	Delete
NewsData	Terveys I Suomalaiset saavat käyttöönsä vain puolet uusista lääkkeistä	Helsingin Sanomat	6/3/2025, 7:30:00 AM	Delete
NewsData	Wetter in Baden-Württemberg: Regen und Gewitter sorgen für stürmische Tage	red/dpa/sw	6/3/2025, 7:30:00 AM	Delete
NewsData	ひろゆきはなぜ論破王と呼ばれるのか 彼が振るう「論破力」その力の正体を徹底追及	Aera Dot.	6/3/2025, 7:30:00 AM	Delete
NewsData	Journalisme en Tunisie : Quelle réforme en temps de répression ?	Mahdi Jlassi	6/3/2025, 7:30:00 AM	Delete
NewsData	EU closing in on the 2030 climate and energy targets, according to national plans	EU Reporter Correspondent	6/3/2025, 7:30:00 AM	Delete
NewsData	Nova Agritech Consolidated March 2025 Net Sales at Rs 81.34 crore, up 14.93% Y-o-Y	Moneycontrol	6/3/2025, 7:30:00 AM	Delete
NewsData	Pourquoi planter des palmiers sur le parvis de l'église à La Seyne?	Var-matin	6/3/2025, 7:30:00 AM	Delete
NewsData	Amber Implants Announces Successful One-Year Follow-Up Data in First-in-Human Clinical Trial of VCFix® Spinal System	Pr Newswire Uk	6/3/2025, 7:30:00 AM	Delete
NewsData	교육부·서울교육청 '리막스를 놀봄' 4일부터 합동 실태조사 착수	손현경 (son89@etoday.co.kr)	6/3/2025, 7:30:00 AM	Delete

[Previous](#)Page 1 of 14217[Next](#)

Figure 3.14: Admin screens: admin articles page

Admin – User Management				
Name	Email	Role	Status	Actions
Adrien K. T.	adrientkoumgang@gmail.com	user	active	<button>View</button>
Adrien Koumgang Tegantchouang	a.koumgangtegantc@studenti.unipi.it	admin	active	<button>View</button>

Previous
Page 1 of 1
Next

Figure 3.15: Admin screens: admin users page

Edit User: Adrien K. T.

Email

adrientkoumgang@gmail.com

First Name

Adrien

Last Name

K. T.

Role

User

Status

Active

Save Changes
Delete User

Figure 3.16: Admin screens: admin user details page

Admin - Error Logs

- UnsafeException** Delete

Test Error: GET

```
curl -X GET -H 'Host: 127.0.0.1:5000' -H 'Sec-Fetch-Dest: empty' -H 'User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/605.1.15 (KHTML, like Gecko) Version/18.5 Safari/605.1.15' -H 'Accept: application/json' -H 'Referer: http://127.0.0.1:5000/api/docs/' -H 'Sec-Fetch-Site: same-origin' -H 'Sec-Fetch-Mode: cors' -H 'Accept-Language: en-US,en;q=0.9' -H 'Priority: u=3, i' -H 'Accept-Encoding: gzip, deflate' -H 'Connection: keep-alive' 'http://127.0.0.1:5000/api/test/error'
```

► Request Data
- UnsafeException** Delete

Test Error: POST

```
curl -X POST -H 'Host: 127.0.0.1:5000' -H 'Sec-Fetch-Site: same-origin' -H 'Accept: application/json' -H 'Origin: http://127.0.0.1:5000' -H 'Sec-Fetch-Mode: cors' -H 'User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/605.1.15 (KHTML, like Gecko) Version/18.5 Safari/605.1.15' -H 'Referer: http://127.0.0.1:5000/api/docs/' -H 'Sec-Fetch-Dest: empty' -H 'Accept-Language: en-US,en;q=0.9' -H 'Priority: u=3, i' -H 'Accept-Encoding: gzip, deflate' -H 'Connection: keep-alive' 'http://127.0.0.1:5000/api/test/error'
```

► Request Data

Previous
Page 1 of 1
Next

Figure 3.17: Admin screens: admin error logs page

Chapter 4

Data Modeling

In designing the project, I employed a **hybrid data modeling strategy** that combines both **relational** and **NoSQL (document-based and key-value)** approaches to balance data consistency, flexibility, and performance at scale.

4.1 Databases Technologies Used

- **MongoDB (Document Database, version 6.0):** Used for modeling core entities like user, articles, comments, user interactions and logs.
- **Redis (Key-Value Store, version 7.0):** Used for caching.

4.2 MongoDB Data Models

MongoDB is used as the primary **document database** to store structured yet flexible data such as users, articles, comment, interactions and errors server. Its schema-less nature and native support for nested documents make it suitable for dynamic content and high-volume ingestion [2].

MongoDB stores structured documents for core domain entities. Each entity is represented as a collection with JSON-like documents. The key collections are:

4.2.1 Users collection

```
1 {  
2   "id": {
```

```

3      "$oid": "683c758466d99ae51f0508d4"
4    },
5    "createdat": {
6      "$date": "2025-06-01T15:45:08.196Z"
7    },
8    "updatedat": {
9      "$date": "2025-06-21T03:25:04.294Z"
10   },
11   "firstname": "Adrien",
12   "lastname": "K. T.",
13   "email": "adrientkoumgang@gmail.com",
14   "password": "$2b$12$tSDjoy8m4IMBS6T/As5zveHctHxvkUmaslrONegD4kRPpkgPARlJK",
15   "account": {
16     "status": "active",
17     "role": "user"
18   },
19   "passwordhistory": [
20     {
21       "password": "$2b$12$guC9jr36y2MEW50le9S4wOdKbZsVnmRrXAvk5S2IMJ4QSkJSTaG5e"
22       ,
23       "createdat": {
24         "$date": "2025-06-01T15:45:08.196Z"
25       }
26     },
27     {
28       "password": "$2b$12$tSDjoy8m4IMBS6T/As5zveHctHxvkUmaslrONegD4kRPpkgPARlJK"
29       ,
30       "createdat": {
31         "$date": "2025-06-05T20:51:45.379Z"
32       }
33     }
34   ],
35   "preferencesenable": true,
36   "preferences": [
37     "ai",
38     "technology",
39     "science"
40   ],
41   "address": {

```

```

40     "street": "via Giovanni Berchet, 4-112",
41     "city": "San Giuliano Terme",
42     "state": "PI",
43     "zip": "56017",
44     "country": "Italy"
45 }
46 }

```

Indexes:

1. $\{email : 1\}$ (Unique) : Email indexes enables fast login lookups and registration check.
2. $\{preferences : 1\}$ (Multikey) : Preferences index accelerates personalized feed generation.

4.2.2 Article log requests collection

```

1 {
2   "id": {
3     "$oid": "68402a68a26f0389a8c4ed9e"
4   },
5   "createdat": {
6     "$date": "2025-06-04T11:13:44.758Z"
7   },
8   "updatedat": {
9     "$date": "2025-06-04T11:13:44.758Z"
10  },
11  "source": "CurrentsAPI",
12  "url": "https://api.currentsapi.services/v1/latest-news",
13  "request": {
14    "url": "https://api.currentsapi.services/v1/latest-news",
15    "headers": {
16      "Authorization": "97ueDaaD6tgJItwJhzKGqd1FYoXm3xZ6H3-lAxKXloZmRYdm"
17    },
18    "params": {}
19  },
20  "response": {
21    "statusCode": 200,

```

```

22     "returned": 30
23 },
24     "fetchcount": 30
25 }

```

Indexes:

1. $\{created_at : -1, source : 1\}$ (Compound): optimize time-based queries filtered by API source. Monitoring recent API fetch operations by provider.
2. $\{response.status_code : -1\}$: quickly identify failed requests for error monitoring and alerting.

4.2.3 Articles collection

```

1 {
2   "id": {
3     "$oid": "6840d72e34b1dd6ef4c87b05"
4   },
5   "createdat": {
6     "$date": "2025-06-04T23:30:53.992Z"
7   },
8   "updatedat": {
9     "$date": "2025-06-04T23:30:53.992Z"
10  },
11  "externid": "b13b8ed2543f2ca7e1874fc2e63282e0",
12  "externapi": "NewsData",
13  "title": "EU closing in on the 2030 climate and energy targets, according to
14          national plans",
15  "description": "EU member states have significantly closed the gap to
16                achieving the 2030 energy and climate targets, according to the European
                Commission's assessment of the National Energy and Climate Plans (NECPs).
                EU countries have substantially improved their plans following Commission
                recommendations in December 2023. As a result, the EU is closing in
                collectively on a 55% reduction in greenhouse gas (GHG) emissions, as
                committed [...]",
17  "content": "ONLY AVAILABLE IN PAID PLANS",
18  "url": "https://www.eureporter.co/environment/2025/06/03/eu-closing-in-on-the
          -2030-climate-and-energy-targets-according-to-national-plans/",

```

```

17  "author": {
18      "name": "EU Reporter Correspondent"
19  },
20  "source": {
21      "name": "Eureporter Co",
22      "url": "https://www.eureporter.co"
23  },
24  "publishedat": "2025-06-03 07:30:00",
25  "language": "english",
26  "country": "united kingdom",
27  "tags": [
28      "environment",
29      "full-image",
30      "climate-neutral economy",
31      "climate change",
32      "featured"
33  ]
34  }

```

Indexes:

1. $\{external_api : 1, external_id : 1, title : 1\}$ (Compound, Unique) : External ID index prevents duplicates articles.
2. $\{published_at : -1\}$: Date indexes support chronological queries
3. $\{tags : 1, published_at : -1\}$ (Compound): this compound index optimizes feed generation
4. $\{tags : 1\}$: for search all possible tags in all articles
5. $\{title : "text", description : "text"\}$ (Text): Text index enables full-text search. (Creation failure due to invalid characters such as Japanese and other non-Latin characters.)

4.2.4 Comments collection

```

1  {
2      "id": {

```

```

3     "$oid": "684d39437a0814577744e01d"
4   },
5   "createdat": {
6     "$date": "2025-06-14T08:26:06.817Z"
7   },
8   "updatedat": {
9     "$date": "2025-06-14T08:26:06.818Z"
10  },
11  "userid": "683c758466d99ae51f0508d4",
12  "articleid": "6840dbaa2c0129e299fbd9db",
13  "content": "My first comment"
14 }

```

Indexes:

1. $\{article_id : 1, created_at : -1\}$ (Compound): this index optimizes comment threading.
2. $\{user_id : 1\}$: User index supports profile activity views
3. $\{user_id : 1, created_at : -1\}$ (Compound): User index supports profile activity views
4. $\{comment_fk : 1, created_at : -1\}$ (Compound): this index enables efficient nested comment retrieval

4.2.5 User Interactions collection

```

1 {
2   "id": {
3     "$oid": "6855a9d4ecbd85c7cbce08e2"
4   },
5   "articleid": "6840db802c0129e299fbbe29",
6   "userid": "683c758466d99ae51f0508d4",
7   "articletitle": "Apple names 2025 Design Awards winners",
8   "levelinteraction": "article",
9   "liked": true,
10  "readat": {
11    "$date": "2025-06-20T18:35:00.451Z"

```

```

12 },
13 "saved": true,
14 "shared": false,
15 "timespent": 2,
16 "updatedat": {
17   "$date": "2025-06-20T18:35:00.451Z"
18 },
19 "report": false
20 }

```

Indexes:

1. $\{user_id : 1, read_at : -1\}$ (Compound): this index powers reading history
2. $\{article_id : 1\}$: this index supports engagement analytics
3. $\{article_id : 1, updated_at : -1\}$ (Compound): this index powers reading article stats

4.2.6 Server Error Logs collection

```

1 {
2   "id": {
3     "$oid": "684241a852621fd670c36a1c"
4   },
5   "createdat": {
6     "$date": "2025-06-06T01:17:05.512Z"
7   },
8   "updatedat": {
9     "$date": "2025-06-06T01:17:05.512Z"
10  },
11  "requestdata": {
12    "url": "http://127.0.0.1:5000/api/test/error",
13    "method": "POST",
14    "body": {},
15    "args": {},
16    "headers": {
17      "Host": "127.0.0.1:5000",
18      "Sec-Fetch-Site": "same-origin",
19      "Accept": "application/json",

```



```

20     "Origin": "http://127.0.0.1:5000",
21     "Sec-Fetch-Mode": "cors",
22     "User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit
        /605.1.15 (KHTML, like Gecko) Version/18.5 Safari/605.1.15",
23     "Referer": "http://127.0.0.1:5000/api/docs/",
24     "Sec-Fetch-Dest": "empty",
25     "Content-Length": "0",
26     "Accept-Language": "en-US,en;q=0.9",
27     "Priority": "u=3, i",
28     "Accept-Encoding": "gzip, deflate",
29     "Connection": "keep-alive"
30 },
31     "form": {}
32 },
33     "curl": "curl -X POST -H 'Host: 127.0.0.1:5000' -H 'Sec-Fetch-Site: same-
        origin' -H 'Accept: application/json' -H 'Origin: http://127.0.0.1:5000' -
        H 'Sec-Fetch-Mode: cors' -H 'User-Agent: Mozilla/5.0 (Macintosh; Intel Mac
        OS X 10_15_7) AppleWebKit/605.1.15 (KHTML, like Gecko) Version/18.5
        Safari/605.1.15' -H 'Referer: http://127.0.0.1:5000/api/docs/' -H 'Sec-
        Fetch-Dest: empty' -H 'Accept-Language: en-US,en;q=0.9' -H 'Priority: u=3,
        i' -H 'Accept-Encoding: gzip, deflate' -H 'Connection: keep-alive' 'http
        ://127.0.0.1:5000/api/test/error'",
34     "exceptionname": "UnsafeException",
35     "exceptionmessage": "Test Error: POST"
36 }

```

Indexes:

1. `{created_at}` : this index support log rotation

4.3 Redis Data Structures

4.3.1 Article Caching

When an article appears in a user's News Feed, there's a probability that the user will want to read it, so the article is fully cached for quick access. This with a ttl of 10 minutes.

4.3.2 User caching

The user comes fully cached for fast access especially for his reading preferences. This with an infinite ttl.

Chapter 5

Data Ingestion and Integration

The **Smart News Aggregator** relies on real-time and scheduled data ingestion from multiple external news APIs. This chapter details the data pipeline architecture, API integration strategies, preprocessing steps, and error handling mechanisms to ensure a consistent and reliable flow of news articles into the system.

5.1 External API Integration

The system integrates with the following news providers:

API	Description	Rate Limits	Data Format
MediaStack	News articles from 7,500+ sources	500 requests/month	JSON
CurrentsApi	-	x requests/month	JSON
Gnews	-	x requests/month	JSON
MarketAux	-	x requests/month	JSON
NYTimes	Premium news content	4000 requests/month	JSON
News Api	-	x requests/month	JSON
News Data	Global news coverage	100 requests/month	JSON
Space Flight News Api	-	x requests/month	JSON
The Guardian	High-quality journalism	5,000 requests/month	JSON

Table 5.1: External Api Integration

Each API uses API keys stored securely in environment variables. Scheduled cron job triggers API calls every monday at 00:00. Requests are made via requests. Re-

sponse Handling: On success, Normalize and store in MongoDB. On failure, Log error, retry with backoff (max 3 attempts).

5.2 Data Processing

5.2.1 Data Validation

- **Duplicate Detection:** checks *externalid* and mongodb index $\{external_api : 1, external_id : 1, title : 1\}$ (Compound, Unique)
- **Schema Validation:** Ensures required field (*title*, *url published_at*)

5.2.2 Data Normalization

Standardized Schema

5.3 Error Handling and Logging

For error type is **rate limit exceeded**, retry the next week. For **network failure**, exponential backoff (max 3 retries). And for **malformed data**, skips invalid entries.

Chapter 6

Backend Implementation

This backend of the **Smart News Aggregator** is built using **Flask** (Python) with **Flask-RESTX** for API development.

6.1 System Architecture

6.1.1 Key Components

Component	Technology	Purpose
API Server	Flask [3, 9] + Gunicor	HTTP request handling
Auth	JWT(RS256) [1]	User authentication
Caching	Redis	Session/store hot data
Asybc Tasks	Flask Cron	Background jobs (API fetching)
Docs	Swagger UI [10]	Interactive API documentation

6.2 API Endpoints

6.2.1 Authentication

Endpoint	Method	Description
<i>/auth/register</i>	POST	User registration
<i>/auth/login</i>	POST	JWT token generation
<i>/auth/login – alt</i>	POST	JWT token generation without password validation
<i>/auth/change_password</i>	POST	Update user password
<i>/auth/me</i>	GET	Get current user information like email, name, status, role

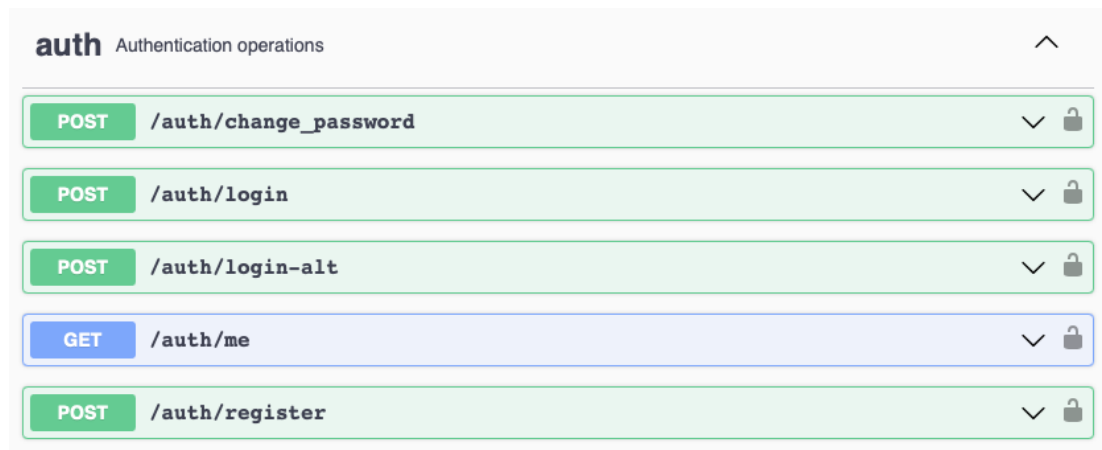


Figure 6.1: API Endpoint: Authentication endpoint

6.2.2 User Management

Endpoint	Method	Description
<i>/user/article/preference</i>	POST	Add article preference tags for current user
<i>/user/article/preference</i>	GET	Get article preference tags for current user
<i>/user/me</i>	POST	Update current user information
<i>/user/me</i>	GET	Get current user information

user User related operations		^
POST	/user/article/preference	✓ 🔒
GET	/user/article/preference	✓ 🔒
POST	/user/me	✓ 🔒
GET	/user/me	✓ 🔒

Figure 6.2: API Endpoint: User endpoint

6.2.3 Article Management

Endpoint	Method	Description
<i>/article/comment/me</i>	GET	List of comment make by current user
<i>/article/history</i>	GET	List of read articles
<i>/article/latest</i>	GET	List of all articles order by recent published
<i>/article/search</i>	GET	List of article where title or description content query
<i>/article/tags</i>	POST	-
<i>/article/tags</i>	GET	-
<i>/article/{article_id}</i>	GET	Get all article data
<i>/article/{article_id}/comment</i>	POST	Add comment on article
<i>/article/{article_id}/comment</i>	GET	List of comment for this specific article
<i>/article/{article_id}/comment/{comment_id}</i>	GET	Get all comment data
<i>/article/{article_id}/comment/{comment_id}</i>	DELETE	Delete specific comment
<i>/article/{article_id}/comment/{comment_id}/interaction</i>	POST	add interaction in this comment
<i>/article/{article_id}/comment/{comment_id}/interaction</i>	GET	Get interaction information in this comment
<i>/article/{article_id}/interaction</i>	POST	Add interaction in this article
<i>/article/{article_id}/interaction</i>	GET	Get article interaction information
<i>/article/{article_id}/summary</i>	GET	Get summable details of this article

article Article endpoint		^
GET	/article/comment/me	✓ 🔒
GET	/article/history	✓ 🔒
GET	/article/latest	✓ 🔒
GET	/article/search	✓ 🔒
POST	/article/tags	✓ 🔒
GET	/article/tags	✓ 🔒
GET	/article/{article_id}	✓ 🔒
POST	/article/{article_id}/comment	✓ 🔒
GET	/article/{article_id}/comment	✓ 🔒
GET	/article/{article_id}/comment/{comment_id}	✓ 🔒
DELETE	/article/{article_id}/comment/{comment_id}	✓ 🔒
POST	/article/{article_id}/comment/{comment_id}/interaction	✓ 🔒
GET	/article/{article_id}/comment/{comment_id}/interaction	✓ 🔒
POST	/article/{article_id}/interaction	✓ 🔒
GET	/article/{article_id}/interaction	✓ 🔒
GET	/article/{article_id}/summary	✓ 🔒

Figure 6.3: API Endpoint: Articles endpoint

6.2.4 Administration Management

Endpoint	Method	Description
/admin/article/{article_id}	DELETE	Delete Article
/admin/articles	GET	List of all articles
/admin/dashboard/errors	GET	List of log of errors occur in server during execution
/admin/dashboard/errors/{server_error_log_id}	DELETE	delete server error log
/admin/users	GET	List of all users
/admin/user/{user_id}	GET	Get all user data
/admin/user/{user_id}	PUT	Add user
/admin/user/{user_id}	POST	Update user
/admin/user/{user_id}	DELETE	Delete User
/admin/	GET	-

admin Administration endpoint			^
GET	/admin/article/{article_id}		✓ 🔒
DELETE	/admin/article/{article_id}		✓ 🔒
GET	/admin/articles		✓ 🔒
GET	/admin/dashboard/activity		✓ 🔒
GET	/admin/dashboard/errors		✓ 🔒
GET	/admin/dashboard/errors/{server_error_log_id}		✓ 🔒
DELETE	/admin/dashboard/errors/{server_error_log_id}		✓ 🔒
GET	/admin/dashboard/summary		✓ 🔒
GET	/admin/dashboard/top-articles		✓ 🔒
GET	/admin/dashboard/top-tags		✓ 🔒
GET	/admin/reload-config		✓ 🔒
GET	/admin/user/{user_id}		✓ 🔒
PUT	/admin/user/{user_id}		✓ 🔒
DELETE	/admin/user/{user_id}		✓ 🔒
GET	/admin/users		✓ 🔒

Figure 6.4: API Endpoint: Admin endpoint

6.3 Security Implementation

6.3.1 JWT Authentication

Algorithm RS256 (asymmetric) for token and **validation middleware**:

```

1 def tokenrequired(f):
2     @wraps(f)
3     def decoratedargs, **kwargs):
4         token = None
5         Extract token from
6         Authorization header
7         auth_header = request.headers.get('Authorization', '')
8         if auth_header.startswith("Bearer "):
9             token = auth_header.split(" ")[1]
10            if not token:
11                return 'error' : 'Token is missing!', 401
12            try:
13                user = TokenManager.decode_token(token=token)
14            except:
15                if 'admin' in request.url and user.role != 'admin':
16                    raise UnauthorizedException('You are not authorized to perform this operation')
17            g.user = user
18            except TokenException as e:
19                return 'error' : str(e), 401
20            return f(*args, **kwargs)
21        return decorated

```

Listing 6.1: Decoration 'Token Required'

6.3.2 Input Validation

All input data is validated via flask-rest's marshal service, which ensures the integrity of input data.

```

1 @nsarticle.route('<string:articleid>\interaction')
2 @nsarticle.param('articleid', 'The article ID')
3 class ArticleInteractionResource(Resource):
4     @tokenrequired
5     @nsarticle.expect(ArticleInteractionType.to_model(namespace=nsarticle))
6     @nsarticle.marshmallow(Model.get_messenger_response_model(namespace=nsarticle))
7     def post(self, articleid):
8         usertoken: UserToken = g.user
9
10        data = request.get_json()
11        interaction = ArticleInteractionType(data)
12        result = UserArticleInteractionModel.update_interaction(user_token, interaction =
13            interaction, article_id = articleid)
14        if result:
15            return "success": True, "message": "Interaction updated"
16        return "success": False, "message": "Interaction not found"

```

Listing 6.2: Input validation with marshal

6.4 Performance Optimizations

6.4.1 Caching Strategies

Cache Key	Data	TTL
<i>user : user_id</i>	Full user json	-
<i>article : article_id</i>	-	-
<i>article : last : count</i>	all last articles count	-
<i>article : last : user_id : count</i>	all last articles count with user preferences	-
<i>article : last : page : limit</i>	full list last articles json by page and limit	-
<i>article : all : count</i>	all articles count	-
<i>article : all : page : limit</i>	full list articles json by page and limit	-
<i>comment : all : filter : page : limit</i>	full list comment by article, page and limit	-

6.4.2 Async Task Processing

When a user accesses a list of items via services such as latest, history, each item in the list is cached.

```
1 class ArticleModel(ArticleSummaryModel):
2     @classmethod
3     def cachearticles(cls, usertoken: UserToken, articles: list):
4         for article in articles:
5             = cls.get(usertoken, str(article.articleid))
6
7     @classmethod
8     def cachearticles(cls, usertoken: UserToken, articles: list):
9         thread = Thread(target=cls.cachearticles, args=(usertoken,
10             articles,))
11         thread.daemon = True
12         thread.start()
```

Listing 6.3: Async Cache Article

6.5 API Documentation

My api's documentation is auto-generated using flask-restx namespaces and swagger accessible via swagger ui via endpoint /api/docs.

```
1 class ArticleSummaryModel(MongoDBBaseModel):
2     @staticmethod
3     def tomodel(namespace: Namespace):
4         return namespace.model('ArticleSummaryModel', {
5             'articleid': fields.String(required=False),
6             'externid': fields.String(required=False),
7             'externapi': fields.String(required=True),
8             'title': fields.String(required=True),
9             'description': fields.String(required=True),
10            'author':
11                fields.Nested(ArticleSourceModel.tomodel(namespace)),
12            'source':
13                fields.Nested(ArticleSourceModel.tomodel(namespace)),
14            'imageurl': fields.String(required=True),
15            'publishedat': fields.String(required=True),
```

```

14         'tags': fields.List(fields.String, description="List of tags"),
15         'currentuserinteraction':
16             fields.Nested(ArticleInteractionStatus.tomodel(namespace)),
17         'totaluserinteraction':
18             fields.Nested(ArticleInteractionStats.tomodel(namespace)),
19     })
20
21 @staticmethod
22 def tomodellist(namespace: Namespace):
23     return namespace.model('ArticleSummaryModelList', {
24         'articles':
25             fields.List(fields.Nested(ArticleSummaryModel.tomodel(namespace))),
26         'total': fields.Integer,
27         'page': fields.Integer,
28         'limit': fields.Integer,
29         'pageCount': fields.Integer,
30     })

```

Listing 6.4: Api Docuementation Article Summary

Chapter 7

Advanced Aggregations and Analytics

7.1 MongoDB Aggregation Framework

7.1.1 Key Aggregation Concepts

- **Pipeline Stages:** Filter (*\$match*), Group (*\$group*), Sort (*\$sort*), Project (*\$project*)
- **Operators:** *\$sum*, *\$avg*, *\$max*, *\$arrayElementAt*, *\$cond*
- **Performance:** Uses indexes, optimized for large datasets

7.1.2 Pipeline: All tags

```
1 class ArticleModel(ArticleSummaryModel):
2     @classmethod
3     def getalltags(cls, usertoken, search: str = None):
4         tags = cls.getalltags()
5         if tags:
6             return tags
7
8         apilogger = ApiLogger(f"[MONGODB] [ARTICLE TAGS] [GET ALL] :
9                               search = {search}")
10        if search:
11            pipeline = [
12                {"$unwind": "$tags"},
13                {"$match": {"tags": {"$regex": f" {search}", "$options":
14                               "i"}}}],
```

```

13         {"$group": {"id": None, "matchedTags": {"$addToSet":
14             "$tags"}}},
15     ]
16     else:
17         pipeline = [
18             {"$unwind": "$tags"},
19             {"$group": {"id": None, "matchedTags": {"$addToSet":
20                 "$tags"}}},
21             {"$project": {"id": 0, "matchedTags": 1}}
22         ]
23     with MONGOQUERYTIME.time():
24         data = cls.collection().aggregate(pipeline)
25
26     result = list(data)
27     tags = result[0]['matchedTags'] if result else []
28
29     apilogger.printlog()
30
31     cls.cachealltags(tags)
32
33     return tags

```

Listing 7.1: Pipeline All Tags

7.1.3 Pipeline: Search Articles

```

1 class ArticleSearchModel(DataBaseModel):
2     @classmethod
3     def searcharticles(cls, usertoken: UserToken, query: str, page: int =
4         1, limit: int = 10):
5         if not query:
6             return ArticleModel.lastarticles(usertoken, page=page,
7                 limit=limit)
8
9         apilogger = ApiLogger(f"[MONGODB] [ARTICLE LATEST] [GET] :
10             query={query}, page={page} and limit={limit}")

```

```

9      pipeline = [
10          {
11              "$match": {
12                  "$text": {
13                      "$search": query,
14                      # "$language": "english"
15                  }
16              }
17          },
18          {
19              "$project": {
20                  'articleid': '$id',
21                  'externapi': 1,
22                  'externid': 1,
23                  "title": 1,
24                  "description": 1,
25                  'source': 1,
26                  'author': 1,
27                  "score": {"$meta": "textScore"},
28                  "publishedat": 1
29              }
30          },
31          {
32              "$sort": {"score": -1, "publishedat": -1} # Relevance +
33                  recency
34          },
35          {
36              "$skip": (page - 1) * limit
37          },
38          {
39              "$limit": limit
40          }
41      ]
42
43      with MONGOQUERYTIME.time():
44          results = ArticleModel.collection().aggregate(pipeline)
45
46      if results is None:
47          apilogger.printerror("Error occurred during article search")

```



```

47         return []
48
49     apilogger.printlog()
50     return [cls*result) for result in list(results)]

```

Listing 7.2: Pipeline search articles

7.1.4 Pipeline: Article Stats Comment

```

1 class ArticleCommentStats(DataBaseModel):
2     @classmethod
3     def getstats(cls, articleid: str, commentid: str = None):
4         statslist = cls.getstats(articleid, commentid)
5         if statslist is not None:
6             return statslist
7
8         apilogger = ApiLogger(f"[MONGODB] [ARTICLE] [MOST COMMENT] :
9                                article={articleid} and comment={commentid}")
10
11         pipeline = [
12             {
13                 '$addFields': {
14                     'articleObjectId': { '$toObjectId': '$articleid' }
15                 }, {
16                     '$group': {
17                         'id': '$articleObjectId',
18                         'commentcount': { '$sum': 1 }
19                     }
20                 }, {
21                     '$sort': {
22                         'commentcount': -1
23                     }
24                 }, {
25                     '$limit': 10
26                 }, {
27                     '$lookup': {
28                         'from': 'articles',
29                         'localField': 'id',

```

```

30         'foreignField': 'id',
31         'as': 'article'
32     }
33 }, {
34     '$unwind': '$article'
35 }, {
36     '$project': {
37         'articleid': '$id',
38         'externapi': '$article.externapi',
39         'title': '$article.title',
40         'source': '$article.source',
41         'author': '$article.author',
42         'publishedat': '$article.publishedat',
43         'commentcount': 1,
44         'id': 0
45     }
46 }
47 ]
48
49 with MONGOQUERYTIME.time():
50     stats = CommentModel.collection().aggregate(pipeline)
51 if stats is None:
52     apilogger.printerror("Error during retrieving statistics")
53     return []
54 apilogger.printlog()
55 statslist = [cls*stat for stat in
               list(stats)]cls.cache(statslist,article_id,comment_id)returnstatslist

```

Listing 7.3: Pipeline Article Stats Comment

7.1.5 Pipeline User Comments with article

```

1 class CommentDetailsModel(CommentModel):
2     @classmethod
3     def getusercommentswitharticle(cls, usertoken: UserToken, userid:
4         str = None, page: int = 1, limit: int = 10):
5         extrafilter = {}
6
7         if userid is None:

```

```

7         userid = str(usertoken.userid)
8
9     pipeline = [
10         {
11             "$match": {
12                 "userid": userid
13             }
14         },
15         {
16             "$sort": {
17                 "createdat": -1
18             }
19         },
20         {
21             "$skip": (page - 1) * limit
22         },
23         {
24             "$limit": limit
25         },
26         {
27             "$addFields": {
28                 "articleidobj": {
29                     "$toObjectId": "$articleid" # Convert string to
30                                             ObjectId
31                 }
32             },
33             {
34                 "$lookup": {
35                     "from": "articles",
36                     "localField": "articleidobj",
37                     "foreignField": "id",
38                     "as": "article"
39                 }
40             },
41             {
42                 "$unwind": "$article"
43             },
44             {

```

```

45         "$project": {
46             "id": 1,
47             "userid": 1,
48             "author": 1,
49             "articleid": 1,
50             "commentfk": 1,
51             "content": 1,
52             "createdat": 1,
53             "updatedat": 1,
54             "articleinfo": {
55                 "externapi": "$article.externapi",
56                 "title": "$article.title",
57                 "description": "$article.description",
58                 "author": "$article.author",
59                 "source": "$article.source",
60                 "publishedat": "$article.publishedat"
61             }
62         }
63     }
64 ]
65
66     apilogger = ApiLogger(f"[MONGODB] [COMMENT] [GET BY USER] :
67         userid={userid}, page={page} and limit={limit}")
68
69     with MONGOQUERYTIME.time():
70         results = cls.collection().aggregate(pipeline)
71
72     apilogger.printlog()
73
74     if results:
75         return [cls*result) for result in results]return []

```

Listing 7.4: Pipeline User Comments with article

7.1.6 Pipeline User Article Interaction by article

```

1 class ArticleInteractionStats(DataBaseModel):
2     @classmethod
3     def getstats(cls, articleid: str, commentid: str = None):

```

```

4     apilogger = ApiLogger(f"[MONGODB] [USER ARTICLE INTERACTION] [GET
      STAT] : article={articleid} and comment={commentid}")
5
6     match = {"articleid": articleid}    ({"commentid": commentid} if
      commentid else {})
7
8     pipeline = [
9         {"$match": match},
10        {
11            "$group": {
12                "id": "$articleid",
13                "liked": {"$sum": {"$cond": [{"$liked", 1, 0}]}},
14                "saved": {"$sum": {"$cond": [{"$saved", 1, 0}]}},
15                "shared": {"$sum": {"$cond": [{"$shared", 1, 0}]}},
16                "report": {"$sum": {"$cond": [{"$report", 1, 0}]}},
17            }
18        }
19    ]
20
21    with MONGOQUERYTIME.time():
22        stats = cls.collection().aggregate(pipeline)
23    if stats is None:
24        apilogger.printerror("Error during retrieving statistics")
25        return ArticleInteractionStats()
26    apilogger.printlog()
27    statslist = list(stats)
28    if statslist:
29        return ArticleInteractionStats(
30            liked=statslist[0]["liked"],
31            saved=statslist[0]["saved"],
32            shared=statslist[0]["shared"],
33            report=statslist[0]["report"],
34        )
35    return ArticleInteractionStats()

```

Listing 7.5: Pipeline User Article Interaction by article

7.1.7 Pipeline Most Interacted Articles

```

1 class ArticleInteractionDashboard(DataBaseModel):

```

```

2 @classmethod
3 def getmostinteractedarticles(cls, datecheck = None):
4     apilogger = ApiLogger(f"[MONGODB] [USER ARTICLE INTERACTION]
5                             [DASHBOARD] [MOST INTERACTED ARTICLES] ")
6
7     if datecheck:
8         pipeline = [
9             {
10                 '$match': {
11                     'updatedat': {'$gte': datecheck}
12                 }
13             }, {
14                 '$addFields': {
15                     'articleObjectId': {'$toObjectId': '$articleid'}
16                 }
17             }, {
18                 '$group': {
19                     'id': '$articleObjectId',
20                     'readcount': {'$sum': {'$cond': [{'$ifNull':
21                                             ['$readat', False]}, 1, 0]}},
22                     'likecount': {'$sum': {'$cond': ['$liked', 1, 0]}},
23                     'savecount': {'$sum': {'$cond': ['$saved', 1, 0]}},
24                     'sharecount': {'$sum': {'$cond': ['$shared', 1,
25                                                         0]}}
26                 }
27             }, {
28                 '$addFields': {
29                     'totalinteractions': {
30                         '$add': ['$readcount', '$likecount',
31                                 '$savecount', '$sharecount']
32                     }
33                 }
34             }, {
35                 '$sort': {
36                     'totalinteractions': -1
37                 }
38             }, {
39                 '$limit': 5
40             }, {

```

```

37         '$lookup': {
38             'from': 'articles',
39             'localField': 'id',
40             'foreignField': 'id',
41             'as': 'article'
42         }
43     }, {
44         '$unwind': '$article'
45     }, {
46         '$project': {
47             'articleid': '$id',
48             'externapi': '$article.externapi',
49             'title': '$article.title',
50             'publishedat': '$article.publishedat',
51             'author': '$article.author',
52             'source': '$article.source',
53             'totalinteractions': 1,
54             'readcount': 1,
55             'likecount': 1,
56             'savecount': 1,
57             'sharecount': 1,
58             'id': 0
59         }
60     }
61 ]
62 else:
63     pipeline = [
64         {
65             '$addFields': {
66                 'articleObjectId': {'$toObjectId': '$articleid'}
67             }
68         }, {
69             '$group': {
70                 'id': '$articleObjectId',
71                 '$sum': {'$cond': [{'$ifNull':
72                     ['$readat', False]}, 1, 0]}},
73                 'likecount': {'$sum': {'$cond': ['$liked', 1, 0]}},
74                 'savecount': {'$sum': {'$cond': ['$saved', 1, 0]}},

```

```

74         'sharecount': {'$sum': {'$cond': ['$shared', 1,
75                                     0]}}
76     }, {
77         '$addFields': {
78             'totalinteractions': {'$add': ['$readcount',
79                                             '$likecount', '$savecount', '$sharecount']}
79         }
80     }, {
81         '$sort': {
82             'totalinteractions': -1
83         }
84     }, {
85         '$limit': 5
86     }, {
87         '$lookup': {
88             'from': 'articles',
89             'localField': 'id',
90             'foreignField': 'id',
91             'as': 'article'
92         }
93     }, {
94         '$unwind': '$article'
95     }, {
96         '$project': {
97             'articleid': '$id',
98             'externapi': '$article.externapi',
99             'title': '$article.title',
100            'publishedat': '$article.publishedat',
101            'author': '$article.author',
102            'source': '$article.source',
103            'totalinteractions': 1,
104            'readcount': 1,
105            'likecount': 1,
106            'savecount': 1,
107            'sharecount': 1,
108            'id': 0
109        }
110    }

```



```

111         ]
112
113     with MONGOQUERYTIME.time():
114         stats =
115             UserArticleInteractionModel.collection().aggregate(pipeline)
116     if stats is None:
117         apilogger.printerror("Error during retrieving statistics")
118
119     statlist = list(stats)
120     apilogger.printlog()
121     # print(statlist)
122
123     for stat in statlist:
124         stat['articleid'] = str(stat['articleid'])
125
126     return [cls*data) for data in statlist]

```

Listing 7.6: Pipeline Most Interacted Articles

7.1.8 Pipeline User Preferences Most Used Tags

```

1 class UserPreferencesDashboard(DataBaseModel):
2     @classmethod
3     def getmosttags(cls, limit: int = 5):
4         apilogger = ApiLogger(f"[MONGODB] [USER TAGS] [DASHBOARD] [MOST
5             TAGS IN PREFERENCES] : limit={limit}")
6
7         pipeline = [
8             {
9                 '$match': {
10                     'preferences': {'$exists': True, '$ne': []}
11                 }
12             }, {
13                 '$unwind': '$preferences'
14             }, {
15                 '$group': {
16                     'id': '$preferences',
17                     'count': {'$sum': 1}
18                 }
19             }
20         ]

```

```

18         }, {
19             '$sort': {
20                 'count': -1
21             }
22         }, {
23             '$limit': limit
24         }, {
25             '$project': {
26                 'tag': '$id',
27                 'count': 1,
28                 'id': 0
29             }
30         }
31     ]
32
33     with MONGOQUERYTIME.time():
34         stats = cls.collection().aggregate(pipeline)
35     if stats is None:
36         apilogger.printerror("Error during retrieving statistics")
37
38     statlist = list(stats)
39     apilogger.printlog()
40
41     return [cls*data) for data in statlist]

```

Listing 7.7: Pipeline User Preferences Most Used Tags

7.2 Performance Optimization

To optimize these aggregation operations, different indexes have been drawn and the aggregation results cached for an average of 1 hour.

- **All Tags**7.1: $\{tags : 1\}$ (1 hour)
- **Search Articles**7.2: $\{title : "text"\}$, $\{description : "text"\}$ and $\{title : "text", description : "text"\}$
- **Article Stats Comment**7.3: $\{article_id : 1\}$ (10 minutes)

- **User Comments with article**7.4: $\{user_id : 1\}$ and $\{user_id : 1, created_at : -1\}$
- **User Article Interaction by article**7.5: $\{article_id : 1\}$
- **Most Interacted Articles**7.6: $\{article_id : 1\}$ and $\{article_id : 1, updated_at : -1\}$
- **User Preferences Most Used Tags**7.7: $\{preferences : 1\}$

Chapter 8

Deployment and Scaling

This chapter covers containerized deployment for the Smart News Aggregator.

8.1 Containerization with Docker

8.1.1 Core Services Setup

```
1 global:
2   scrapeinterval: 15s
3
4 scrapeconfigs:
5   - jobname: 'smart-news-aggregator-api'
6     staticconfigs:
7       - targets: ['smart-news-aggregator-api:5050']
8       # - targets: ['localhost:5050']
9
10  - jobname: 'redis'
11    staticconfigs:
12      - targets: ['redis-exporter:9121']
13
14  - jobname: 'mongodb'
15    staticconfigs:
16      - targets: ['mongodb-exporter:9216']
```

Listing 8.1: Prometheus Configuration (prometheus.yml)

```
1 services:
2   redis:
```

```

3   image: redis:alpine
4   ports:
5     - "6379:6379"
6   volumes:
7     - redisdata:/data
8     - redisbackup:/backup
9   command: redis-server --save 60 1 --loglevel warning
10  restart: unless-stopped
11
12  redis-exporter:
13    image: oliver006/redisexporter
14    ports:
15      - "9121:9121"
16    environment:
17      REDISADDR: "redis://host.docker.internal:6379"
18
19  mongodb:
20    image: mongo:latest
21    containername: mongodb
22    ports:
23      - "27017:27017"
24    volumes:
25      - mongodbddata:/data/db
26      - mongoddbbackup:/backup
27    platform: linux/arm64
28    healthcheck:
29      test: echo 'db.runCommand("ping").ok' | mongosh localhost:27017/test
30        --quiet
31      interval: 10s
32      timeout: 10s
33      retries: 5
34    restart: unless-stopped
35
36  mongodb-exporter:
37    image: bitnami/mongodb-exporter:0.40.0
38    ports:
39      - "9216:9216"
40    environment:
41      MONGODBURI: "mongodb://host.docker.internal:27017"

```

```

41
42 prometheus:
43     image: prom/prometheus
44     ports:
45         - "9090:9090"
46     volumes:
47         - ./prometheus.yml:/etc/prometheus/prometheus.yml
48
49 grafana:
50     image: grafana/grafana
51     ports:
52         - "3000:3000"
53     volumes:
54         - grafanadata:/var/lib/grafana
55     dependson:
56         - prometheus
57
58 volumes:
59     redisdata:
60     redisbackup:
61     mongodbddata:
62     mongoddbbackup:
63     grafanadata:

```

Listing 8.2: Docker compose for core services (docker-compose.yml)

8.1.2 Backend Setup

```

1 # Dockerfile.prod
2 FROM python:3.11-slim
3
4 WORKDIR /app
5
6 COPY requirements.txt .
7 RUN pip install --no-cache-dir -r requirements.txt
8
9 COPY . .
10
11 ENV FLASKAPP=src/app.py

```

```

12 ENV FLASKENV=production
13
14 EXPOSE 5000
15
16 # Use Gunicorn for production
17 CMD ["gunicorn", "-b", "0.0.0.0:5000", "src.app:application"]

```

Listing 8.3: Dockerfile Backend (Production configuration)

```

1 services:
2   smart-news-aggregator-api:
3     build:
4       context: .
5       dockerfile: Dockerfile.prod
6     ports:
7       - "5050:5000"
8     environment:
9       - FLASKENV=production
10      - FLASKENVFILE=.env.prod
11     envfile:
12       - .env.prod

```

Listing 8.4: Docker Compose Backend (docker-compose.yml)

8.1.3 Frontend Setup

```

1 # Stage 1: Build the Vite app
2 FROM node:20 as builder
3
4 WORKDIR /app
5
6 COPY package*.json ./
7 RUN npm install
8
9 COPY . .
10 RUN npm run build
11
12 # Stage 2: Serve with Nginx
13 FROM nginx:stable-alpine

```

```

14
15 # Copy the build output to Nginx web root
16 COPY --from=builder /app/dist /usr/share/nginx/html
17
18 # Optional: custom nginx config for single-page app (SPA)
19 COPY nginx.conf /etc/nginx/conf.d/default.conf
20
21 EXPOSE 80
22
23 CMD ["nginx", "-g", "daemon off;"]

```

Listing 8.5: Dockerfile Frontend (Production configuration)

```

1 services:
2   smart-news-aggregator-fe:
3     build:
4       context: .
5       dockerfile: Dockerfile.prod
6     ports:
7       - "3000:80" # host:container

```

Listing 8.6: Docker Compose Frontend (docker-compose.yml)

Listing 8.7: (.yml)

8.2 Testing and Monitoring

To ensure the reliability and correctness of the platform, the services were thoroughly tested and monitored after deployment. All backend and frontend services were containerized using Docker, enabling isolated and consistent environments across development and production.

Once deployed, the services were interacted with directly through the exposed endpoints to verify the correct execution of key functionalities such as authentication, article retrieval, comment posting, and admin operations (e.g., deleting users [8] or moderating content). API responses were validated, edge cases were tested, and rate limits were checked.

Monitoring was achieved through an integrated dashboard, where metrics such as service uptime, API response times, error logs, and user activity were visualized. This allowed real-time observation of system behavior and rapid identification of anomalies or performance bottlenecks. Container logs were also inspected using docker logs and integration with monitoring tools like **Grafana** and **Prometheus**.

Chapter 9

Conclusion

9.1 Project Achievements

This project successfully designed and implemented a **Smart News Aggregator & User Personalization Platform**, leveraging modern full-stack technologies including **Flask, MongoDB, Redis, React, and Docker**. The platform allows users to browse, search, and interact with a large volume of news articles in real time while supporting advanced features such as:

- User account and preference management
- Article search and filtering based on user interests
- Commenting, liking, sharing, and saving interactions
- A full-featured admin dashboard for user, article, and error management
- Efficient use of **MongoDB** for document storage and **Redis** for caching frequently accessed data

Several **core achievements** of the project include:

- **Scalable API architecture** supporting multiple concurrent users with JWT-based authentication
- **Real-time aggregation** of user interactions (comments, likes, etc.) using MongoDB aggregation pipelines

- **Dynamic and responsive React UI** with pagination, filtering, and error handling
- **Deployment automation with Docker**, enabling isolated and reproducible environments
- **Monitoring and logging mechanisms** to observe platform behavior and debug failures efficiently

9.2 Challenges Overcome

During development, several technical and architectural challenges were encountered and addressed:

- **CORS preflight request issues** when interacting across frontend and backend containers were resolved through proper flask-cors configuration
- Handling **inconsistent data types** (e.g., ObjectId vs string) in MongoDB required careful pipeline design and casting
- Designing **efficient indexes** in MongoDB to optimize aggregation performance, especially for high-traffic collections like comments and interactions
- Implementing **data parsing and transformation** for various date/time formats using pydantic, datetime, and external APIs
- Managing **rate limits** and error logging during external API interactions with fallback and retry mechanisms

These challenges provided valuable insights into working with real-world data, full-stack integration, and production-level deployment.

9.3 Future Work

While the platform already demonstrates a solid foundation, several opportunities for enhancement exist:

- **GraphQL API Layer:** Flexible frontend data fetching
- **Advanced recommendation engine** using collaborative or content-based filtering to suggest articles based on user behavior and preferences
- Integration with **machine learning models** for topic classification, sentiment analysis, or fake news detection
- **Mobile app version** of the platform using React Native or Flutter to increase accessibility
- **Newsletter Automation:** Implementation of **email notification systems** to alert users of trending articles or replies to their comments
- **Role-based access control** with more granular permissions for moderators and administrators
- **PDF Article Export:** Enhanced accessibility
- Full integration with **monitoring tools** like Prometheus, Grafana, and ELK Stack for real-time performance and alerting
- **CI/CD pipeline setup** for automated deployment, testing, and delivery across environments
- Extension of data sources by connecting to more external news APIs and supporting multilingual content

Bibliography

- [1] Auth0: JWT.IO - JSON Web Tokens (2024), <https://jwt.io/introduction>
- [2] Colvin, S.: Pydantic documentation (2024), <https://docs.pydantic.dev>
- [3] Community, P.R.: Flask-restx documentation. <https://flask-restx.readthedocs.io/> (2024)
- [4] (formerly Facebook), M.: React – A JavaScript Library for Building User Interfaces (2024), <https://react.dev>
- [5] Inc., D.: Docker Documentation (2024), <https://docs.docker.com>
- [6] Inc., M.: MongoDB Manual (2024), <https://www.mongodb.com/docs/manual/>
- [7] Labs, R.: Redis Documentation (2024), <https://redis.io/docs>
- [8] RandomUser.me: Random User Generator API (2024), <https://randomuser.me/documentation>
- [9] Ronacher, A.: Flask Documentation (2024), <https://flask.palletsprojects.com>
- [10] Software, S.: Swagger: Openapi specification. <https://swagger.io/specification/> (2024)
- [11] Sullivan, D.: NoSQL for Mere Mortals. Addison-Wesley (2015)
- [12] You, E.: Vite – Next Generation Frontend Tooling (2024), <https://vitejs.dev>

Acknowledgments

First and foremost, I would like to express my deepest gratitude to **God**, whose grace, strength, and guidance sustained me throughout the duration of this project.

I sincerely thank all the **professors and academic staff** who provided the foundation, mentorship, and encouragement that made this work possible. Your knowledge and dedication have been invaluable.

I would also like to acknowledge the **authors of NoSQL for Mere Mortals** for their clear explanations and practical approach, which greatly enhanced my understanding of NoSQL data modeling.

My heartfelt thanks to the developers and contributors of the various **tools and frameworks** used in this project — including **Flask, React, MongoDB, Redis**, and **Docker** — as well as the extensive and well-maintained **documentation** that enabled effective development.

Special thanks to **DeepL**, which helped with accurate and nuanced **translations**, and to **ChatGPT** and **DeepSeek**, whose assistance in technical guidance, code generation, and documentation structuring significantly accelerated the progress and quality of this work.

I extend my heartfelt thanks to my friends for their moral support, encouragement, and constant motivation. Their presence has been a great source of strength.

Finally, I thank my girlfriend for her patience, love, and unwavering support throughout this journey.

Finally, I am grateful to all those who supported me, directly or indirectly, in the successful completion of this project.