

Jaypee Institute of Information Technology, Noida

Department Of Computer Science & Engineering and IT



Major Project Title: Marvin - Scalable AI Search Engine for Efficient Webpage Discovery and Retrieval

Enroll. No.	Name of Student
21103240	Arnav Teotia
21103261	Rohan Siwach
21103268	Milind Choudhary

Course Name: Major Project 1

Program: B. Tech. CS&E

4th Year 7th Sem

2024 - 2025

(I)
DECLARATION

We hereby declare that this submission is our own work and that, to the best of our knowledge and belief, it contains no material previously published or written by another person nor material which has been accepted for the award of any other degree or diploma of the university or other institute of higher learning, except where due acknowledgment has been made in the text.

Place:

Signature:

Date: 21/11/2024

Name: Rohan Siwach

Enrollment No: 21103261

Place:

Signature:

Date: 21/11/2024

Name: Arnav Teotia

Enrollment No: 21103240

Place:

Signature:

Date: 21/11/2024

Name: Milind Choudhary

Enrollment No: 21103268

(II)

CERTIFICATE

This is to certify that the work titled "**Scalable AI Search Engine for Efficient Webpage Discovery and Retrieval**" submitted by "**Arnav Teotia**", "**Rohan Siwach**" and "**Milind Choudhary**" in partial fulfilment for the award of degree of Bachelor of Technology of Jaypee Institute of Information Technology, Noida has been carried out under my supervision. This work has not been submitted partially or wholly to any other University or Institute for the award of this or any other degree or diploma.

Signature of Supervisor :

Name of Supervisor : Dr. Parmeet Kaur

Designation : Associate Professor

Date : 21/11/2024

(III)

ACKNOWLEDGEMENT

We would like to express our special thanks of gratitude to our project supervisor "**Dr. Parmeet Kaur**" as well as Jaypee Institute of Information Technology, Noida who gave us the golden opportunity to do this wonderful project on the topic "**Scalable AI Search Engine for Efficient Webpage Discovery and Retrieval**", which also helped us learning about so many new things which we are really thankful for. Secondly, we would also like to thank our friends who helped us a lot in the collection of the text descriptions required for this project within the limited time frame.

Signature of Student :

Name of Student : Arnav Teotia

Enrollment Number : 21103240

Date : 21/11/2024

Signature of Student :

Name of Student : Rohan Siwach

Enrollment Number : 21103261

Date : 21/11/2024

Signature of Student :

Name of Student : Milind Choudhary

Enrollment Number : 21103268

Date : 21/11/2024

(IV)

SUMMARY

Our project, Marvin: AI Search Engine, combines an intuitive Chrome extension interface with advanced data handling to create an efficient, AI-powered bookmark manager. Designed to capture and store both selected and scraped text from webpages. These bookmarks can later be queried using both semantic search for understanding intent and keyword-based search for precision, allowing users to retrieve information effortlessly by asking specific questions.

The project's frontend is built in TypeScript, while the backend utilizes FastAPI and Langchain to seamlessly handle data processing and retrieval. Saved bookmarks are stored in PostgreSQL and Vespa, the latter hosted in a Docker container, enabling rapid, vectorized search capabilities. Using a Retrieval-Augmented Generation (RAG) approach, AI Search Engine employs a language model via API integration to retrieve and generate responses to user queries, allowing for an intelligent and precise search experience.

By leveraging advanced technologies in a user-friendly Chrome extension, AI Search Engine serves as a valuable tool for users who need quick access to saved content, enhancing their research and productivity workflows through effective information retrieval.

Signature of Student

Name : Arnav Teotia

Date : 21/11/2024

Signature of Supervisor

Name : Dr. Parmeet Kaur

Date : 21/11/2024

Signature of Student

Name : Rohan Siwach

Date : 21/11/2024

Signature of Student

Name : Milind Choudhary

Date : 21/11/2024

List Of Figures

1. Use Case Diagram	
.....	29
2. Class Diagram	
.....	33
3. Control Flow Diagram	
.....	35
4. Sequence Diagram	
.....	37
5. Schema Design	
.....	39

(V)
TABLE OF CONTENTS

Chapter No.	Topics	Page No.
Chapter-1	Introduction	8-11
	1.1 General Introduction	
	1.2 Problem Statement	
	1.3 Significance/Novelty of the problem	
	1.4 Brief Description of Solution Approach	
	1.5 Comparison of existing approaches to the proposed system	
Chapter-2	Literature Survey	14-15
	2.1 Summary of papers studied	
	2.2 Integrated Summary of the Literature studied	
Chapter -3	Requirement Analysis and Solution Approach	16-18
	3.1 Requirement Analysis	
	3.2 Proposed System	
Chapter -4	Modelling and Implementation Details	18-22
	4.1 Design Diagrams	
	4.1.1 Control flow Diagrams	
	4.1.2 Data Flow Diagram	
	4.1.3 Use Case Diagram	
	4.1.4 Class Diagram	
	4.1.5 Sequence Diagram	
	5.2 Implementation Details and issues	
	5.2.1 SQL and Vector Database Connection	
	5.2.2 Google Gemini API Integration	
	5.2.3 Chrome Extension Interface	
	5.2.4 Error Handling and Robustness	

5.3 Risk Analysis and Mitigation

Chapter -5	Testing (Focuses on Quality of Robustness and Testing)	27-36
5.1	Testing Plan	
5.3	Error and Expectation Handling	
5.4	Limitations of the Solutions	
Chapter -6	Findings, Conclusion and Future Work	36-40
6.1	Findings	
6.2	Conclusion	
6.3	Future Work	
References		40-40

1. Introduction

1.1 General Introduction

The project, Marvin: AI Search Engine, is designed to revolutionize the way users manage and retrieve information from their saved web content, offering an AI-powered bookmark manager within a Chrome extension. Developed in collaboration with LangChain for advanced natural language processing, Marvin AI empowers users to capture selected and scraped text from web pages, save it, and later retrieve it using intelligent search queries.

With Langchain's cutting-edge language models, Marvin AI enables accurate, contextually relevant responses to user queries, creating an intuitive and precise retrieval experience. Marvin AI's user-friendly interface allows users to save and query bookmarks effortlessly, harnessing both PostgreSQL and Vespa (in a Docker container) for efficient data storage and fast, vectorized search capabilities. Through its Retrieval-Augmented Generation (RAG) setup, Marvin AI combines advanced AI tools with backend integration in FastAPI to ensure responsive, insightful answers.

Marvin - AI Search Engine aims to be a valuable tool for researchers, professionals, and information enthusiasts, simplifying content retrieval and maximizing productivity. This project exemplifies the potential of AI-driven solutions in making saved information readily accessible and meaningful, ultimately enhancing the way users engage with their own curated content.

1.2 Problem Statement

- a) Users often struggle to efficiently organize and retrieve relevant content from the vast amount of web data they encounter daily, making traditional bookmarking tools insufficient.
- b) Bookmarking tools provide a basic way to save links but lack sophisticated retrieval options, making it difficult to locate and utilize saved content effectively.
- c) Researchers, professionals, and students face significant challenges in quickly accessing specific information within their bookmarked resources.
- d) Traditional bookmark managers fail to capture meaningful context or extract key information from saved web pages, requiring users to manually sift through large lists.
- e) These tools lack AI-powered capabilities for querying bookmarks based on context or specific questions, resulting in inefficient and time-consuming retrieval processes.
- f) The absence of natural language interaction with saved content leads to productivity losses, as users struggle to transform bookmarked information into actionable insights.
- g) Bridging this gap requires a tool that combines bookmarking with advanced AI and NLP technologies to enable contextual search, retrieval, and interaction with saved web data.
- h) This project, Marvin-AI Search Engine, addresses these challenges by offering an intuitive Chrome extension with integrated NLP, leveraging large language models and a vector database setup for enhanced bookmark management.
- i) The tool enables users to capture, store, and retrieve specific content from bookmarks using natural language queries, streamlining workflows and empowering users to transform saved information into actionable knowledge.

1.3 Significance/Novelty of the problem

- a) The significance of Marvin, the AI Search Engine, lies in its ability to revolutionize information retrieval, enhancing productivity and accessibility through advanced data handling. It offers a fresh solution to traditional limitations in bookmark management. Below are the key aspects of its significance and innovation:
- b) Enhanced Knowledge Access: Marvin democratizes access to users' stored web content by facilitating intelligent and contextual queries. This enables users to quickly locate relevant information, transforming saved bookmarks into actionable insights and boosting productivity and informed decision-making.
- c) Technological Innovation: Marvin integrates a Retrieval-Augmented Generation (RAG) approach using Langchain's language models, creating a novel synergy between AI, vector databases, and Chrome extensions. This integration of natural language processing directly into the bookmarking process offers a more intuitive, powerful way for users to interact with their saved content.
- d) Efficiency in Content Retrieval: Traditional bookmarking lacks advanced search capabilities, making it difficult to access specific information swiftly. Marvin overcomes this by combining AI and vector databases, enabling precise, context-driven search results that save time and enhance retrieval accuracy.
- e) Real-Time Information Retrieval: Marvin provides real-time responses to user queries based on their stored data, making it highly adaptable for dynamic and evolving research or project needs. This facilitates agile access to content, keeping users' workflows efficient.
- f) User-Centric Design: Through an intuitive Chrome extension interface, Marvin improves the bookmarking experience, making it easier for researchers, students, and professionals to interact with their saved content. Its seamless integration of advanced search within a simple UI ensures that users can quickly find relevant information without a steep learning curve.

- g) **Interdisciplinary Utility:** Marvin's applications extend across fields like academia, journalism, and professional research. By offering structured and contextual access to stored information, it serves as a valuable tool for a wide range of users seeking efficient information retrieval.

1.4 Brief Description of Solution Approach

- a) Development of an AI-powered Chrome extension designed to offer efficient bookmarking and intelligent search capabilities through an intuitive interface.
- b) Implementation of advanced search features using both semantic and keyword-based queries, allowing users to retrieve stored web content based on contextual relevance and precision.
- c) Utilization of Retrieval-Augmented Generation (RAG) with Langchain's language models to enhance the AI's ability to generate contextually accurate responses to user queries.
- d) Integration of PostgreSQL and Vespa for scalable, efficient storage and rapid vectorized search, with Vespa hosted in a Docker container for optimized data processing and retrieval.
- e) Seamless frontend and backend integration using TypeScript for the frontend, and FastAPI for the backend, ensuring smooth communication between the user interface and AI-driven backend systems.

Tech Stack used:

React, TypeScript, Tailwind CSS, FastAPI, Langchain, PostgreSQL, Vespa, Docker

Platforms:

VS Code, PostgreSQL, Docker

API Testing:

ThunderClient, Postman

1.5 Comparison of Existing solution to the problem framed

The proposed system, Marvin -AI Search Engine, introduces a transformative solution in bookmark management by integrating AI for advanced data retrieval. Below is a comparison of MarvinAI's innovative features against traditional bookmarking systems.

Existing Approaches:

Basic Bookmark Managers: Standard bookmarking tools such as Chrome's built-in manager allow users to save URLs but lack contextual metadata or advanced retrieval features. Users must manually scroll through lists, making it challenging to locate relevant content quickly.

Tagging and Folder Systems: Some tools offer tagging and folder categorization for bookmarks, providing limited organization but still requiring manual searching. These systems offer minimal assistance in retrieving bookmarks based on specific topics or detailed questions.

Third-Party Bookmark Managers: Advanced bookmarking platforms like Pocket or Raindrop.io allow users to save and organize content, often with visual previews. However, they do not support natural language querying or intelligent search capabilities, restricting efficient information retrieval and exploration.

Keyword-Based Search: Existing platforms may include simple keyword search options, but these are typically limited to title and URL searches, lacking the depth to interpret and retrieve based on contextual information or complex queries.

Proposed System (Marvin: AI Search Engine):

AI-Driven, Contextual Bookmark Retrieval: Marvin-AI Search Engine employs a Retrieval-Augmented Generation (RAG) setup that allows users to query saved content in natural language. This eliminates manual searching, enabling quick, context-aware responses based on the content and metadata of saved bookmarks.

Natural Language Interface: Unlike standard systems, Marvin-AI Search Engine uses Langchain's language models to support natural language queries. Users can retrieve specific content without needing exact keywords, making information access more accessible and intuitive.

Dynamic and Scalable Search Capabilities: By storing bookmarks in a vector database like Vespa, marvin-AI Search Engine supports fast, flexible, and scalable search, offering deeper insights into saved content compared to static folder or tag systems.

User-Friendly, Interactive Chrome Extension: marvin-AI Search Engine's Chrome extension provides an intuitive UI, accessible to both technical and non-technical users, and integrates directly into the browsing experience. Its seamless interface contrasts with more complex third-party systems requiring additional software or manual organization.

Real-Time, Customizable Insights: marvin-AI Search Engine's system provides real-time responses to user queries, enabling a more agile way to interact with saved information for dynamic research needs or project-based workflows.

2. Literature Summary

2.1 Summary of papers studied

- a) 1 "**Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks**" (Lewis et al., 2020) introduced RAG, a model combining retrieval and generation to improve knowledge-intensive tasks. By fetching relevant documents and conditioning the generative model on them, RAG enhances accuracy and creativity in responses. The approach showed significant improvements over traditional models in tasks like open-domain question answering but highlighted challenges in retrieving irrelevant data, which could impact the output's quality. This work laid the foundation for RAG's evolution, emphasizing the synergy between retrieval and generation (Lewis et al., 2020).
- b) "**LangChain: A Framework for Building Language Models with External Tools**" (Zhao et al., 2023) introduced LangChain, a framework for integrating language models with external tools like APIs and databases. LangChain enhances RAG capabilities by modularizing the workflow of combining retrieval and generation tasks. It allows developers to customize retrieval strategies, improve model efficiency, and create responsive AI systems for dynamic tasks. This work highlights LangChain's potential for building more powerful knowledge-based applications by bridging the gap between static language models and external data sources (Zhao et al., 2023).

- c) **"Enhancing Retrieval-Augmented Generation Models for Information Retrieval"** (Khandelwal et al., 2021) focused on improving retrieval strategies in RAG models to boost document retrieval and generation efficiency. By refining retrieval methods to better rank documents and reduce irrelevant information, they demonstrated improved accuracy in fact-based and long-tail queries. Their work also emphasized the importance of document selection quality for generating more contextually relevant answers. While their retrieval enhancements showed promising results, challenges like balancing retrieval precision and generation flexibility remain (Khandelwal et al., 2021).
- d) **"Improving Document Retrieval and Generation in RAG Systems"** (Yang et al., 2022) explored methods to refine both document retrieval and generation in RAG systems, proposing multi-stage retrieval and semantic matching techniques. They demonstrated that by improving the accuracy of document retrieval, the quality of generated responses in knowledge-intensive tasks could be enhanced. This work contributed to making RAG systems more efficient by minimizing irrelevant document retrieval while ensuring coherent and relevant output generation, although it pointed out the computational cost of more complex retrieval strategies (Yang et al., 2022).

2.2 Integrated study of literature studied

Retrieval-Augmented Generation (RAG) has emerged as a powerful technique for enhancing knowledge retrieval in AI systems by combining retrieval and generation in a unified model. The seminal work by Lewis et al. (2020) introduced the RAG framework, which leverages retrieved documents to condition the generation process, significantly improving the relevance of the generated outputs. This approach outperformed traditional methods by allowing the model to generate responses that are more grounded in external knowledge. However, the effectiveness of RAG can be compromised by irrelevant retrievals, as noted by Khandelwal et al. (2021), who explored ways to refine retrieval strategies to address this issue. In addition, Yang et al. (2022) further advanced RAG by incorporating semantic matching techniques to ensure that the retrieved documents are more relevant and coherent, enhancing the overall quality of generated responses.

Building on these foundational works, LangChain (Zhao et al., 2023) introduced a framework that extends RAG by facilitating the seamless integration of external tools, such as databases and APIs, into the language model's workflow. This integration allows for more dynamic and contextually aware data retrieval, improving the model's ability to generate contextually relevant responses. The development of RAG and frameworks like LangChain showcases a significant leap forward in knowledge-intensive NLP tasks, enabling AI systems to interact with external resources more effectively and generate more accurate, informed outputs.

3. Requirement Analysis and Solution Approach

3.1 Requirement Analysis

Our requirement for the implementation and execution of the project involved using the following:

1. **User Requirements:** Understanding the needs and preferences of AI Search Engine's target users is critical. This involves identifying user demographics, their experience with web browsing and information retrieval, and their expectations for a seamless, AI-powered bookmark management system..
2. **Functional Requirements:** These specific features and functionalities are essential for AI Search Engine to meet its objectives:
 - Natural Language Processing (NLP): Integrate language model capabilities (via Langchain) to interpret user queries in natural language..
 - Bookmark Storage and Retrieval: Store and retrieve bookmarks with associated metadata (selected and scraped text) using Vespa, a vector database, for precise and contextual search.
 - Real-Time Query Responses: Enable real-time responses to user queries, allowing users to find relevant information from saved content instantly.
 - User-Friendly Chrome Extension Interface: Develop a user-friendly extension UI for easy access, organization, and search of bookmarks.

- Bookmark Metadata and Categorization: Allow users to add tags, metadata, and descriptions to bookmarks for improved organization and retrieval.
3. **Non-Functional Requirements:** These quality attributes define AI Search Engine's performance, usability, and reliability::
- Performance: The extension should process and respond to queries promptly, minimizing user wait time.
 - Usability: The UI should be intuitive and straightforward, allowing users to navigate and manage bookmarks with minimal technical expertise.
 - Reliability: AI Search Engine should consistently provide accurate and relevant information when retrieving bookmark content
 - Security: Measures should be implemented to secure user data and prevent unauthorized access to stored bookmarks.
 - Scalability: AI Search Engine should handle a growing volume of bookmarks and user queries without compromising system performance.
4. **Technical Requirements:** These include the technologies and tools required to develop and deploy the chatbot. Technical requirements may include:
- Programming languages: TypeScript for the frontend (Chrome extension) and Python for the backend development.
 - Database management systems : PostgreSQL for data storage and Vespa (in a Docker container) for vector-based bookmark retrieval..
 - Web development frameworks : FastAPI for backend services, managing bookmark data, and handling user queries.
 - LLM Integration: Langchain's language models for processing user queries and generating contextual responses.

3.2 Proposed System

The proposed system, AI Search Engine, is an innovative Chrome extension and bookmark manager designed to simplify access to saved web content through a seamless, AI-powered querying experience. Built with advanced natural language processing (NLP) technologies, AI Search Engine enables users to organize, retrieve, and explore their bookmarks more efficiently by storing and processing selected text in a vector database.

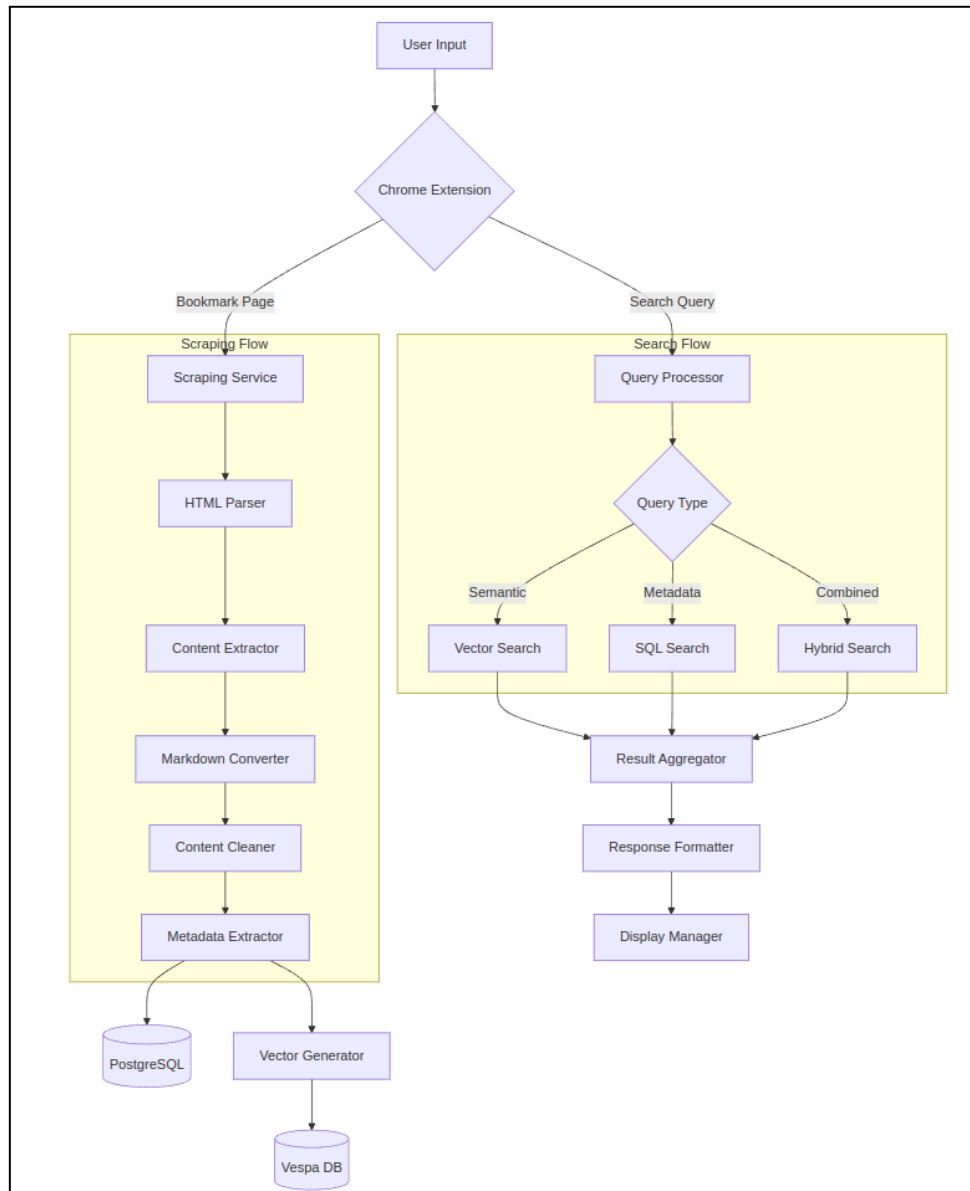
Key Components and Functioning:

1. **Natural Language Processing (NLP):** AI Search Engine uses advanced language models through Langchain to understand and process user queries in natural language. This capability enables users to search and retrieve bookmarks contextually, providing highly relevant results based on the query's intent and nuances.
2. **Integration with Vespa and PostgreSQL Databases:** AI Search Engine integrates with Vespa, a vector database, to store and retrieve bookmark data, and PostgreSQL for structured metadata storage. These databases hold selected and scraped text from bookmarks, making it possible for AI Search Engine to return accurate and meaningful responses to user questions based on their saved content.
3. **User-Friendly Chrome Extension Interface:** AI Search Engine's interface is designed as a Chrome extension using TypeScript, providing users with an intuitive platform to manage, organize, and search through bookmarks. This user-friendly design allows users to quickly add bookmarks, view metadata, and query stored content without navigating outside their browser.
4. **Bookmark Metadata and Categorization:** AI Search Engine provides features for tagging, categorizing, and adding metadata to bookmarks, enhancing the organization and retrieval of saved content. Users can add descriptions, tags, and other details to streamline the search process.

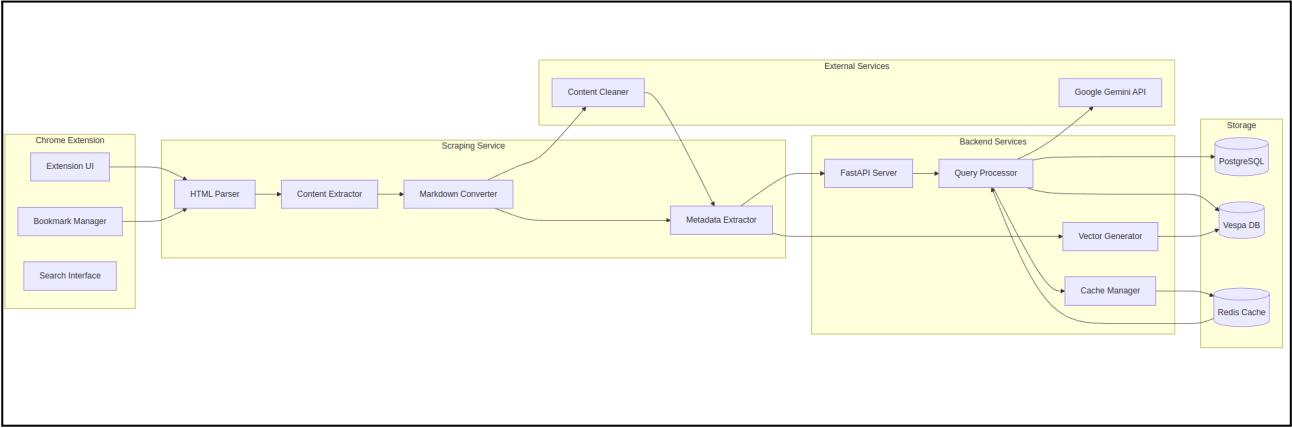
4. Modelling and Implementation Details

4.1 Design Diagrams

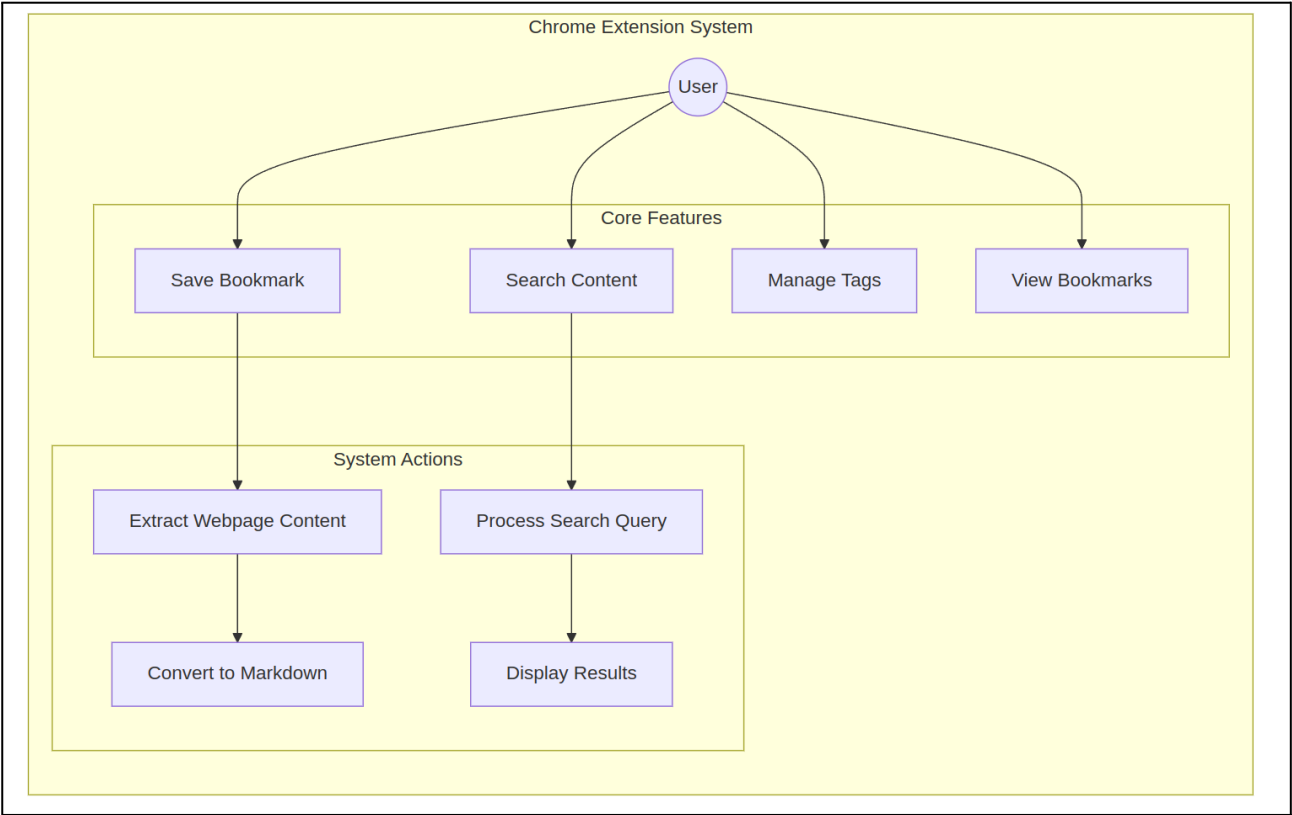
4.1.1 Control flow Diagrams



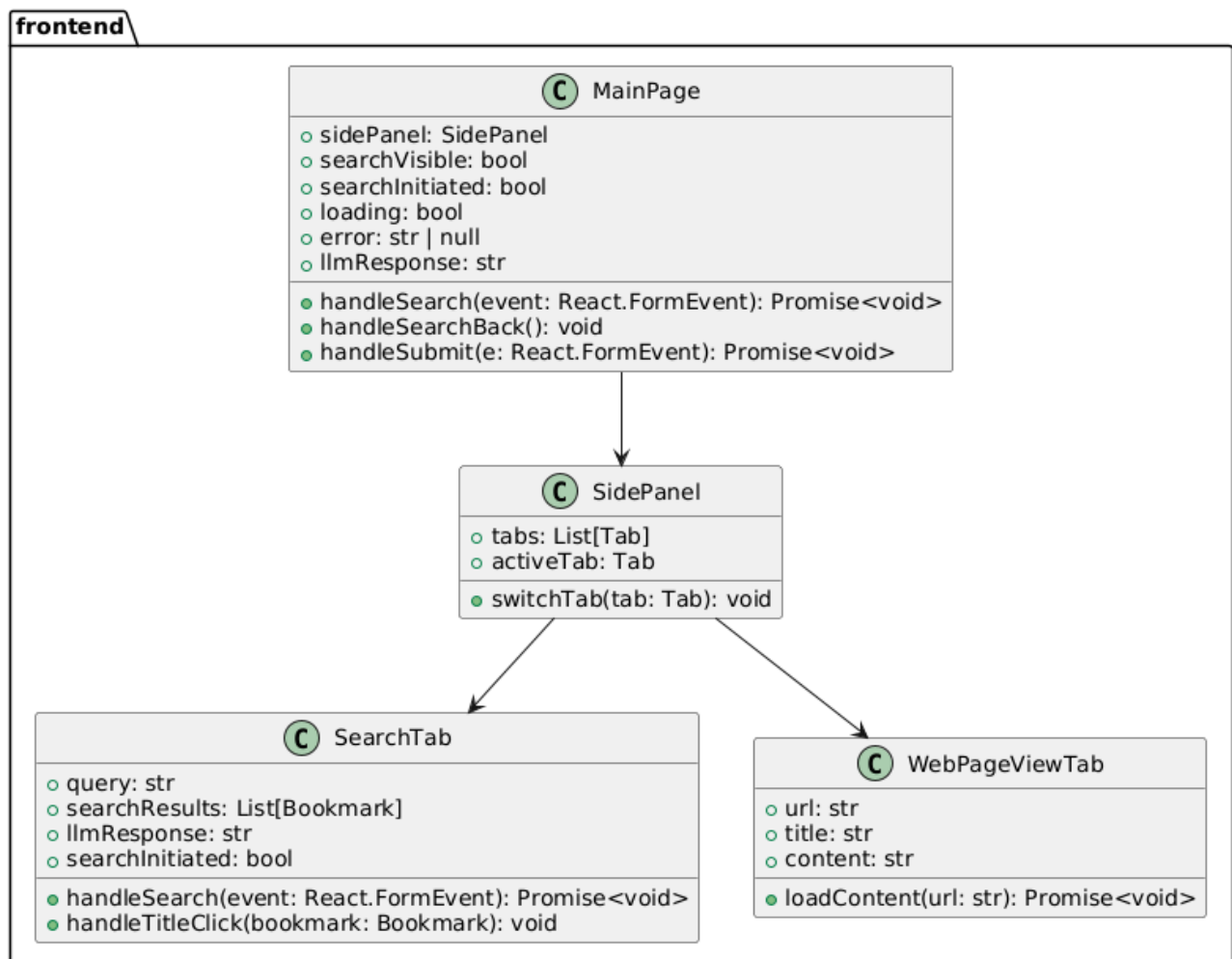
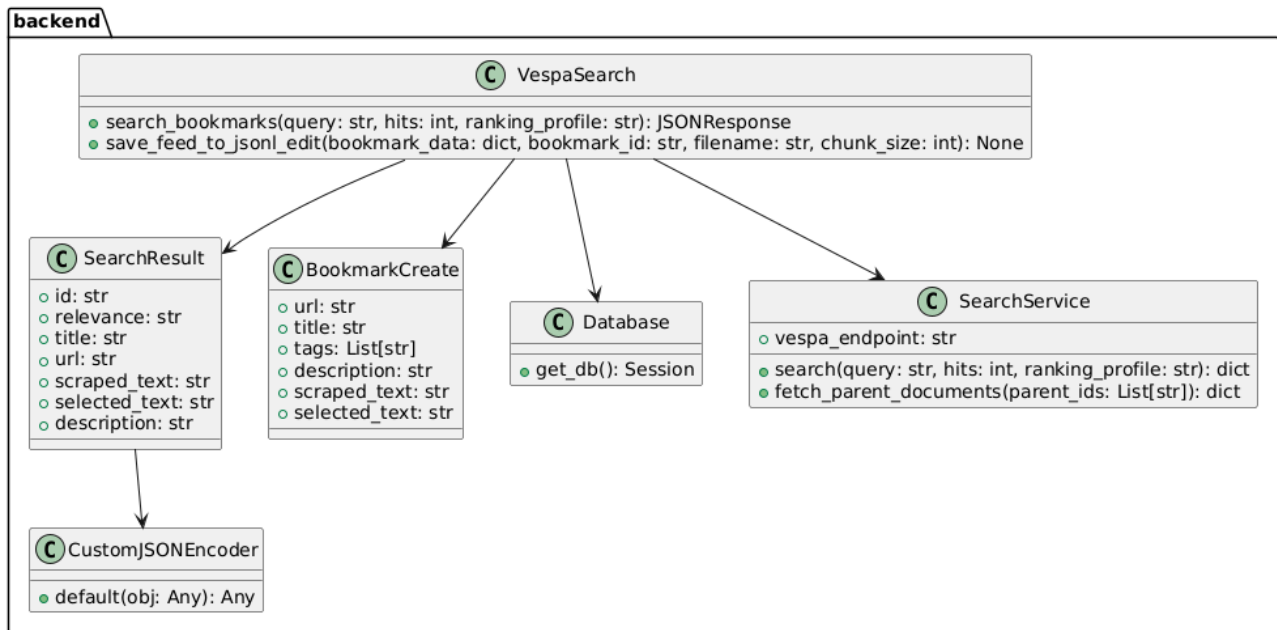
4.1.2 Data Flow Diagram



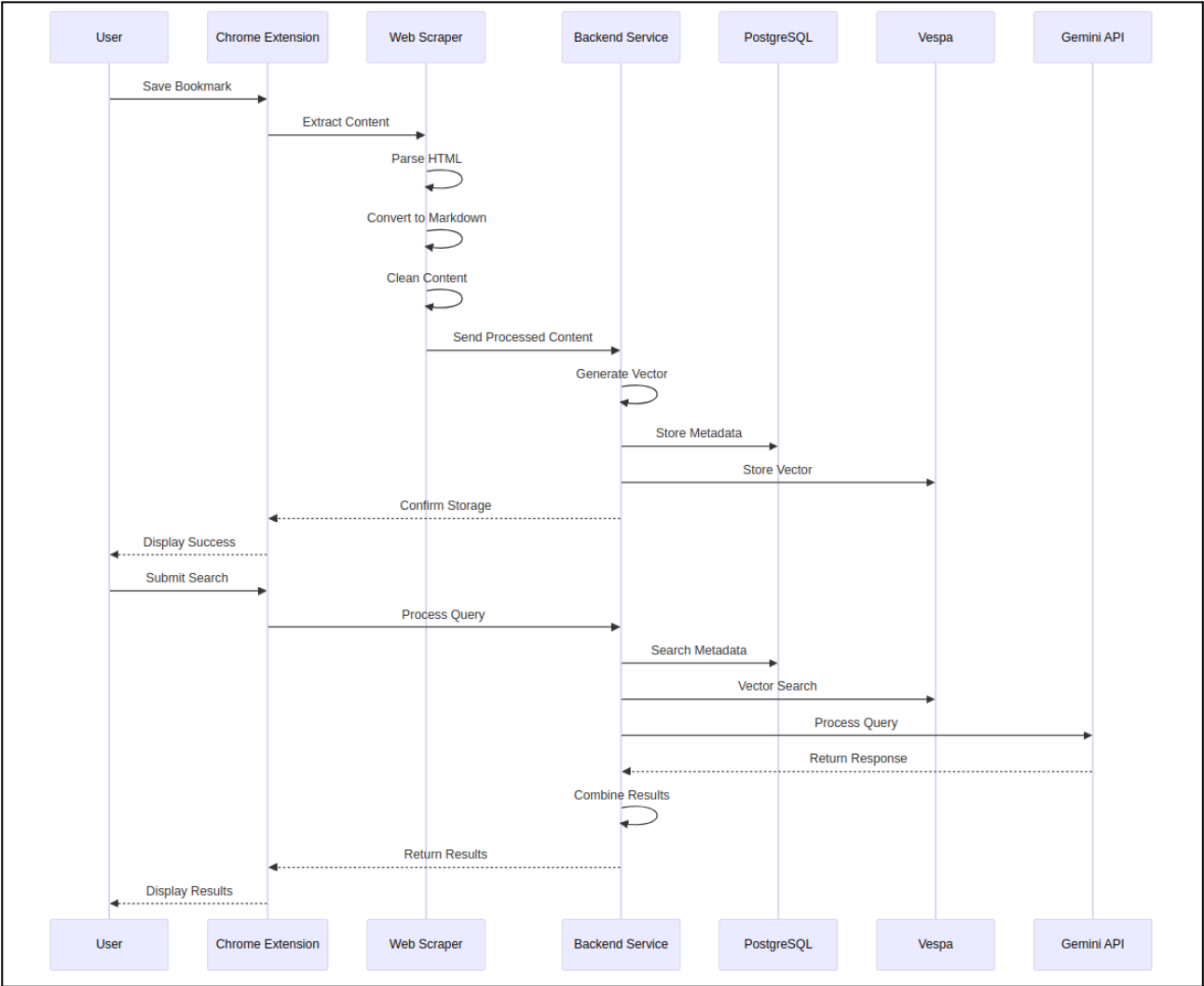
4.1.3 Use Case Diagram



4.1.4 Class Diagram



4.1.5 Sequence Diagram



4.2 Implementation Details and issues

4.2.1 SQL and Vector Database Connection

- **Implementation Details:** The project integrates both PostgreSQL and Vespa for managing and retrieving data. PostgreSQL is used for storing structured metadata about bookmarks, such as URLs, tags, and scraped content, while Vespa is used for vector-based semantic search. SQLAlchemy facilitates ORM-based interaction with PostgreSQL, while Vespa APIs handle operations related to storing and querying vector embeddings of selected and scraped text.
- **Potential Issues:**
 - **Connection Errors:** Errors may occur due to incorrect database connection strings, unreachable servers, or misconfigured Vespa endpoints.
 - **Authentication Errors:** Incorrect credentials or mismanaged access permissions can block connections to either PostgreSQL or Vespa.
 - **Database Schema Changes:** Updates to the database schema or vector indexing in Vespa without corresponding code changes may cause errors during queries or storage operations.

4.2.2 Google Gemini API Integration

- **Implementation Details:** The project integrates the Google Gemini API to handle natural language understanding and response generation. Renowned for its advanced

NLP capabilities, the Gemini API processes user queries and produces highly accurate, context-aware responses tailored to user inputs.

Potential Issues

While leveraging the Gemini API for language model inference offers numerous advantages, it also presents certain challenges that merit consideration:

- **API Rate Limits:** The functionality depends on the Gemini API, and exceeding usage limits could lead to throttling or temporary suspension of access.
- **Model Performance:** The relevance and effectiveness of responses rely on the API's performance. Inaccuracies or unexpected responses may require post-processing or error-handling mechanisms.
- **Response Latency:** API queries introduce latency, potentially impacting the responsiveness of the chatbot interface and user experience.

4.2.3 Chrome Extension Interface

- **Implementation Details:** The Chrome extension interface is developed using TypeScript and React to provide seamless bookmark management. Users can capture bookmarks, scrape content from webpages, and add metadata through a form interface. The extension communicates with the backend via FastAPI endpoints to store and retrieve data.
- **Potential Issues:**
 - **Browser Compatibility:** The extension may face compatibility issues with different versions of Chrome or Chromium-based browsers.
 - **Data Integrity:** Errors in transmitting selected or scraped text to the backend could lead to incomplete or incorrect storage.
 - **User Experience:** Poorly designed UI elements or unclear workflows might reduce user satisfaction and effectiveness.

4.2.4 Error Handling and Robustness

- **Implementation Details:** The project includes error handling mechanisms to gracefully handle exceptions and display error messages to users when database connections fail or API requests encounter errors.
- **Potential Issues:**
 - **Error Reporting:** Incomplete or unclear error messages may hinder users' ability to understand and troubleshoot issues. Comprehensive error reporting and logging can help diagnose problems more effectively.
 - **Edge Cases:** Unexpected inputs or corner cases might bypass the current error-handling logic, leading to unhandled exceptions..

4.3 Risk Analysis and Mitigation

Risk: Data Inconsistency

Description: The database may contain inconsistent or outdated information, leading to inaccuracies in the responses provided by the chatbot.

Mitigation Strategy:

1. **Data Validation:** Implement rigorous validation checks to ensure the integrity and consistency of data before processing queries.
2. **Regular Updates:** Establish procedures for regularly updating the database with the latest information from reliable sources to minimize the risk of outdated data.
3. **Error Handling:** Develop robust error-handling mechanisms to identify and flag inconsistencies in the data, providing users with transparent notifications when encountered.

Risk: Database Connectivity Issues

Description: Unforeseen network issues or server downtime may disrupt communication between the chatbot and the underlying database, causing service interruptions.

Mitigation Strategy:

1. **Connection Monitoring:** Implement monitoring systems to continuously monitor the status of database connections and promptly detect any connectivity issues.
2. **Automatic Reconnection:** Configure the system to automatically attempt reconnection to the database in case of connection failures, minimizing downtime.
3. **Fallback Mechanism:** Develop a fallback mechanism to provide users with alternative means of accessing information or services during database outages, such as cached responses or predefined fallback actions.

Risk: Security Vulnerabilities

Description: Unauthorized access, data breaches, or malicious attacks targeting the database infrastructure pose significant security risks, potentially compromising sensitive information.

Mitigation Strategy:

1. **Data Encryption:** Utilize encryption techniques to secure sensitive data stored in the database, preventing unauthorized access even in the event of a breach.
2. **Access Control:** Implement stringent access control measures, including role-based access control (RBAC) and strong authentication mechanisms, to restrict access to the database to authorized users only.
3. **Regular Security Audits:** Conduct periodic security audits and vulnerability assessments to identify and address potential weaknesses or vulnerabilities in the database infrastructure proactively.

By proactively identifying potential risks and implementing robust mitigation strategies, we aim to safeguard the integrity, availability, and confidentiality of data accessed and managed by our chatbot system.

The risks identified in the project are based on SEI Risk Taxonomy. These include personnel, budget, security, and external dependencies.

Risk ID	Classification	Description	Risk Area	Probability	Impact	Risk Exposure (P x I)
1	Personnel	Lack of expertise with Vespa or vector databases.	Development Team	High	Medium	15
2	Budget/Cost	Limited funding for Vespa hosting and scaling.	Financial Constraints	Medium	High	15
3	Security	Risk of unauthorized access to user bookmarks.	Application Security	High	High	25
4	External Dependency	Google Gemini API downtime affects queries..	Third-Party Dependency	Medium	High	15

Interrelationship Graph (IG)

The risk areas are connected in an interrelationship graph. For example:

- **Application Security** directly impacts **Third-Party Dependency** (weight = 9, significant).
- **Personnel Risks** influence **Financial Constraints** (weight = 3, medium).

Risk Area	# of Risk Statements	Weights (In + Out)	Total Weight	Priority
Application Security	4	25	50	1

Personnel Risks	3	9	27	2
Financial Constraints	2	15	15	3

5 Testing (Focuses on Quality of Robustness and Testing)

5.1 Testing Plan

Objective:

The testing plan is a comprehensive strategy devised to ensure the robustness, accuracy, and user satisfaction of AI Search Engine. It aims to identify and rectify potential issues across various stages of development and deployment.

Approach:

1. **Unit Testing:** Each component of the system, including database interactions, agent functionalities, and user interfaces, will undergo rigorous unit testing to validate its individual behavior and functionality.
2. **Integration Testing:** Following successful unit tests, the system components will be integrated and tested together to ensure seamless communication and interoperability.
3. **User Acceptance Testing (UAT):** Real users will evaluate the system to verify its usability, responsiveness, and adherence to requirements. Feedback will be collected to refine the system further.

4. **Performance Testing:** The system's performance will be assessed under different loads to determine its responsiveness, scalability, and resource utilization.
5. **Security Testing:** Comprehensive security testing will be conducted to identify and address potential vulnerabilities, ensuring the confidentiality, integrity, and availability of user data.

Scenarios:

1. **Basic and Complex Query Handling:** Evaluate the system's ability to accurately interpret and respond to both simple and complex user queries.
2. **Error Handling:** Verify the system's capability to gracefully handle errors and unexpected inputs, preventing system failures or data corruption.
3. **Concurrency Testing:** Assess the system's behavior under concurrent user interactions to ensure stability and data consistency.
4. **Edge Cases:** Test the system's response to unusual or extreme scenarios to ensure robustness and prevent unexpected behavior.

Tools:

1. **Automated Testing Frameworks:** Utilize tools such as pytest and Selenium for automated testing of system functionalities and user interfaces.
2. **Load Testing Tools:** Employ tools like Apache JMeter or Locust for simulating high loads and analyzing system performance under stress.
3. **Security Testing Tools:** Utilize security testing tools such as OWASP ZAP or Burp Suite for identifying and addressing security vulnerabilities.

Criteria:

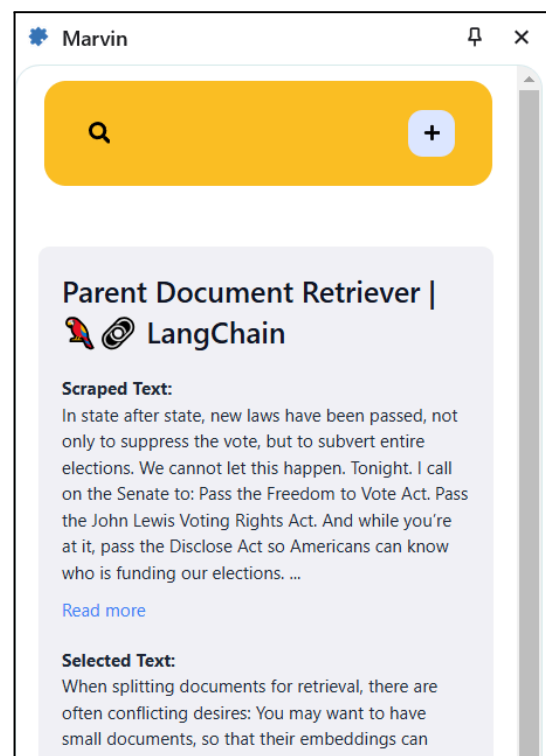
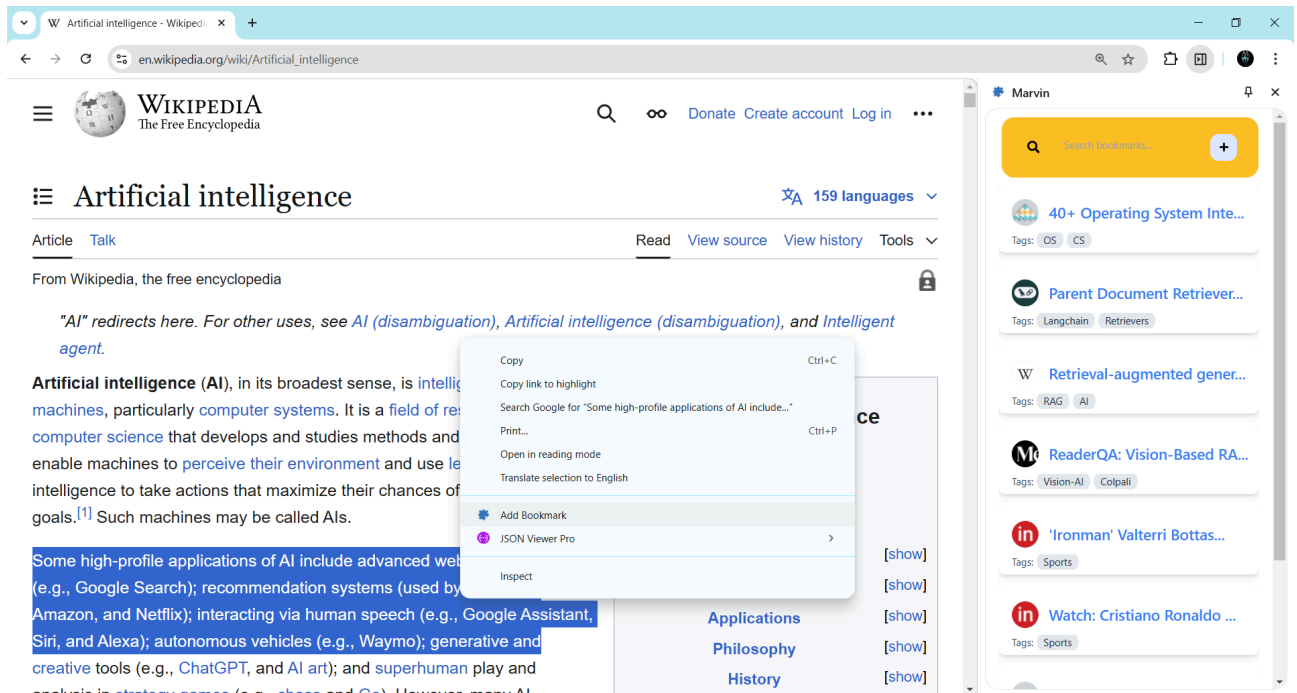
1. **Accuracy:** Ensure that the system provides accurate and relevant responses to user queries based on the underlying data and algorithms.

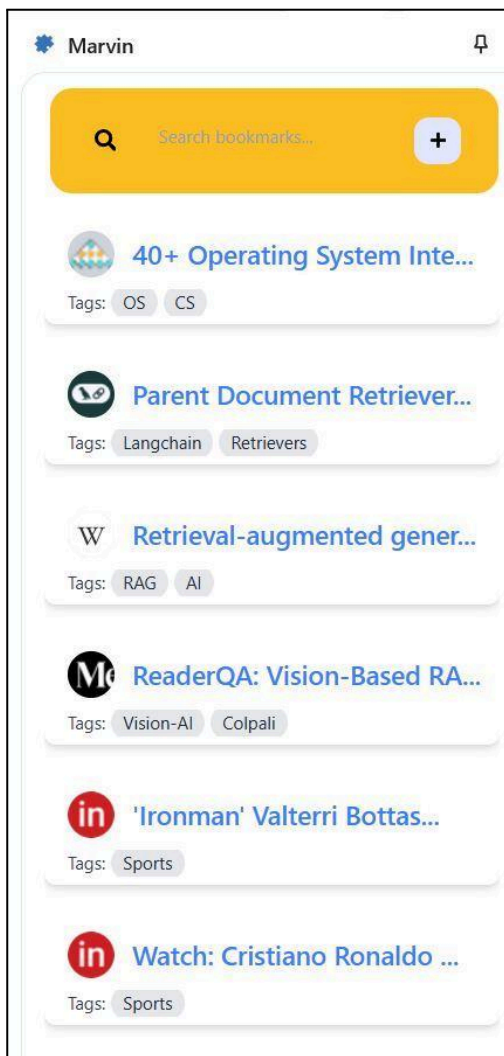
2. **Usability:** Evaluate the system's ease of use, intuitiveness, and user-friendliness to ensure a positive user experience.
3. **Performance:** Assess the system's responsiveness, latency, and resource utilization under normal and peak loads to ensure optimal performance.
4. **Security:** Identify and address potential security vulnerabilities to safeguard user data and prevent unauthorized access or data breaches.

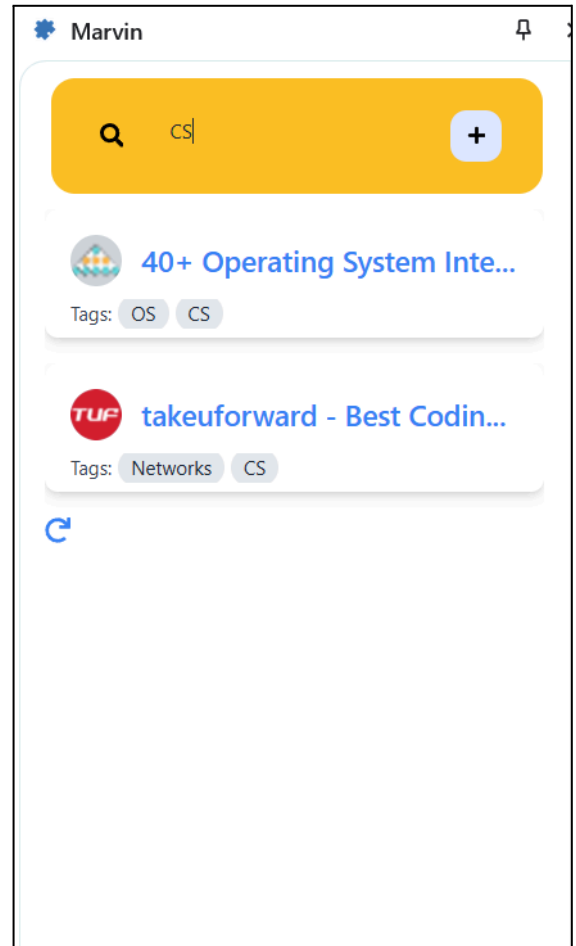
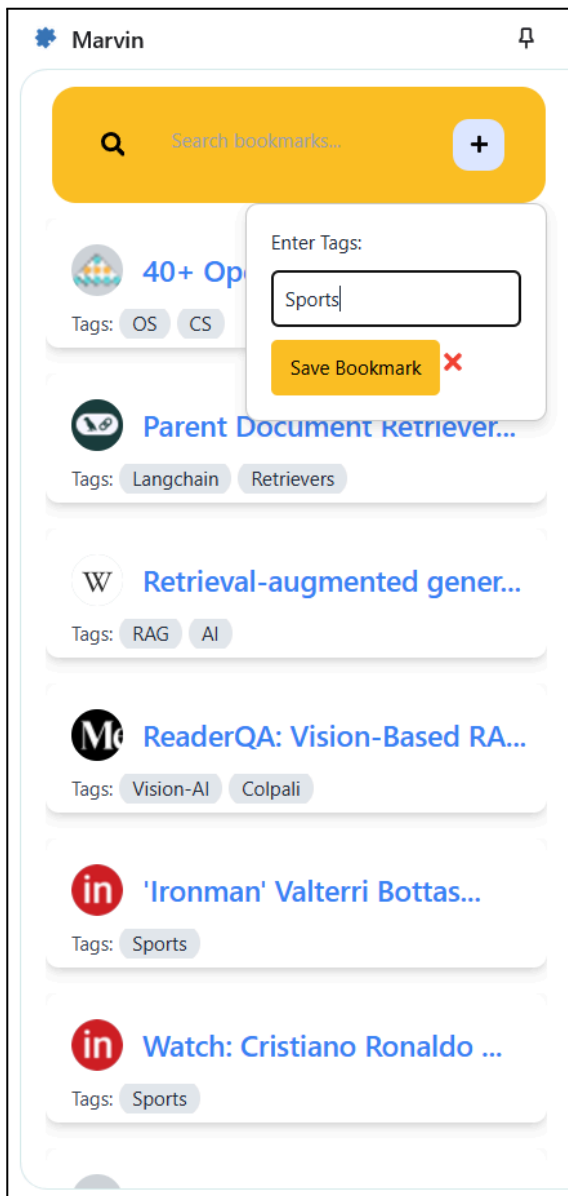
Schedule:

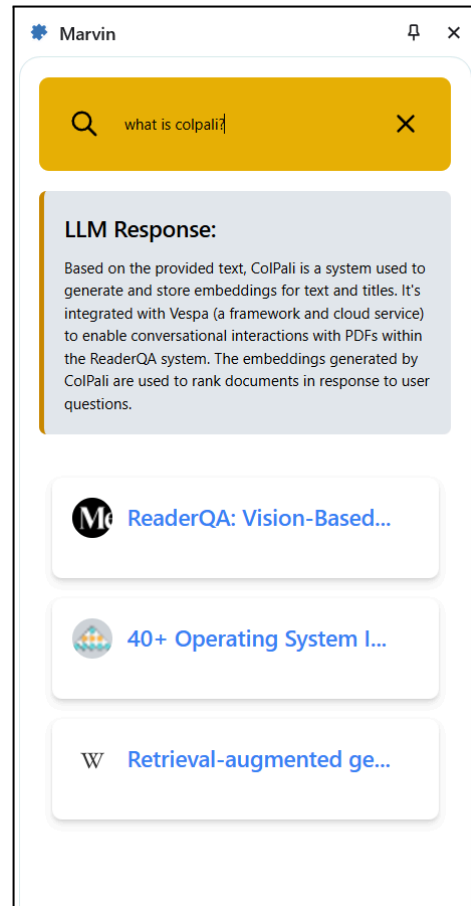
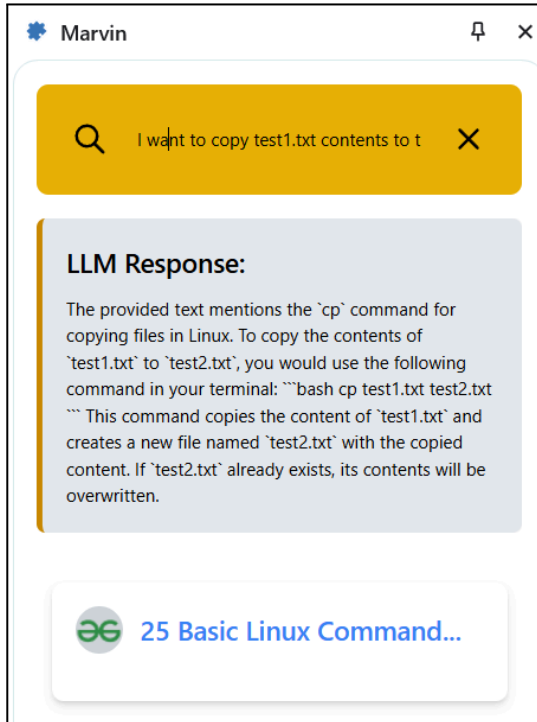
1. **Unit Testing:** Conducted iteratively during the development process to verify the correctness of individual components.
2. **Integration Testing:** Performed before UAT to ensure seamless integration and functionality across system components.
3. **UAT:** Scheduled after the development phase to gather feedback from real users and refine the system based on their inputs.
4. **Performance and Security Testing:** Conducted periodically throughout the development lifecycle to identify and address performance and security issues proactively.

5.2 Working Screens









Objective:

The error and exception handling mechanism in AI Search Engine ensures robustness, reliability, and seamless user experience by detecting, handling, and recovering from errors and exceptional conditions.

Strategies:

Input Validation: User inputs are validated to ensure they meet specified criteria and formats, preventing invalid inputs from causing errors.

Exception Handling: Try-except blocks are implemented to catch and handle exceptions gracefully. Different types of exceptions, such as database errors or input/output errors, are handled appropriately to prevent system crashes.

Error Logging: Logging mechanisms record error messages, stack traces, and contextual information for debugging and troubleshooting purposes.

User-Friendly Error Messages: Clear and informative error messages guide users on resolving issues or seeking support assistance.

Fallback Mechanisms: Fallback strategies are implemented to handle unexpected scenarios or errors gracefully, such as using cached data or providing alternative responses.

Automated Monitoring: Automated monitoring and alerting systems detect and respond to errors in real-time, ensuring prompt identification and resolution of issues.

Implementation:

Input Validation: User queries are validated against predefined patterns or criteria to ensure they are well-formed and adhere to expected formats.

Exception Handling: Critical code blocks are wrapped in try-except blocks to catch and handle exceptions. Specific error messages or fallback behaviors are implemented based on the nature of the exception.

Error Logging: Logging libraries are integrated to log error messages, stack traces, and contextual information. Log files are stored centrally and periodically reviewed for identifying recurring issues or trends.

User-Friendly Error Messages: Standardized error messages are defined for common failure scenarios, providing guidance to users on potential resolutions or next steps.

Continuous Improvement: Error and exception handling mechanisms are continuously reviewed and improved based on ongoing monitoring, user feedback, and system performance analysis. Regular audits and post-mortem analyses identify areas for enhancement and optimization. Updates and patches are deployed promptly to address newly discovered vulnerabilities or error patterns.

5.3 Limitations of the Solution

While AI Search Engine offers valuable insights and assistance in querying the database for information highlighted by the user, it also has several limitations that users should be aware of:

Limited Knowledge Scope: The extension is restricted to the content saved by the user and the predefined query capabilities. It cannot answer questions that go beyond the stored data or require external domain knowledge.

Challenges with Ambiguity: To deliver accurate responses, the extension relies on well-defined and specific user queries. Vague or ambiguous questions may result in irrelevant or incorrect answers.

Lack of Natural Language Understanding: While the Chatbot can process user inputs in natural language, its understanding is limited to the predefined patterns and templates. It may struggle to grasp nuanced language, colloquialisms, or contextually complex queries.

Potential Data Inconsistencies: The accuracy of the Chatbot's responses is contingent upon the accuracy and completeness of the underlying data in the database. Inconsistencies, inaccuracies, or missing data in the database may lead to erroneous or incomplete responses.

Limited Error Handling: While the Chatbot incorporates error and exception handling mechanisms, it may not cover all possible failure scenarios. Users should be prepared to encounter occasional errors or unexpected behaviors and may need to refine their queries or seek assistance in such cases.

Sensitive Data Handling: The Chatbot does not incorporate mechanisms for handling sensitive or confidential data. Users should avoid querying or sharing sensitive information through the Chatbot interface.

6. FINDINGS, CONCLUSION, AND FUTURE WORK

6.1 Findings

The development and evaluation of Marvin: AI Search Engine revealed several key insights:

1. **Integration of Retrieval-Augmented Generation (RAG):** The incorporation of the RAG approach significantly improved search quality, allowing users to retrieve highly relevant information by combining the power of semantic and keyword-based search with the generative capabilities of language models. This hybrid method enabled more accurate and contextually relevant responses.
2. **Efficiency of Vectorized Search:** By leveraging Vespa's vector search capabilities, Marvin achieved faster and more efficient retrieval of saved bookmarks, ensuring that users could access stored content in real-time without delays.
3. **User Experience and Interaction:** Feedback from initial users highlighted the system's intuitive interface and the ease with which they could query saved bookmarks using natural language. This positive reception indicated that Marvin met its objective of improving user productivity and information retrieval efficiency.
4. **Challenges in Handling Complex Queries:** The system faced difficulties in interpreting highly ambiguous or multi-faceted queries. These challenges, along with occasional issues in the retrieval of optimal data, suggested areas for improvement in query parsing and response generation.

6.2 Conclusion

In conclusion, Marvin: AI Search Engine represents a groundbreaking approach to enhancing information retrieval through the integration of advanced AI technologies. By leveraging the LangChain framework alongside a powerful Retrieval-Augmented Generation (RAG) methodology and the robust capabilities of PostgreSQL and Vespa, Marvin offers an efficient, user-centric solution for managing and querying saved bookmarks. The system's hybrid search capabilities, combining semantic and keyword-based approaches, ensure users can access contextually relevant and precise information with minimal effort. Marvin's success in providing an intuitive user interface, coupled with its AI-driven language processing, positions it as an invaluable tool for users who require quick and accurate access to their stored content, ultimately transforming how bookmarked data is managed and retrieved.

6.3 Future Work

Enhancing Multimodal Capabilities: Incorporating multimodal data handling, such as integrating image, audio, and video data alongside textual content, would significantly broaden Marvin's application. This could allow users to query and retrieve information across various media types, providing a more comprehensive search experience.

Improved Natural Language Understanding: Continuous refinement of the system's language models to handle even more complex, context-rich queries would make Marvin more adaptable to diverse user needs.

Data Source Expansion: The integration of more real-time data sources could further enrich the system's responses, providing users with even more dynamic, up-to-date information.

Scalability and Performance Optimizations: Ensuring the system remains responsive and efficient as the volume of data and concurrent queries grows will be essential for maintaining performance.

Multi-Language and Localization Support: Extending the system's linguistic capabilities to support a broader range of languages will increase its accessibility to global users, ensuring inclusivity across diverse linguistic regions.

REFERENCES:

- [1] Langchain Documentation, *Langchain Documentation*, 2023. [Online]. Available: <https://langchain.com/docs/>. [Accessed: Nov. 19, 2024].
- [2] Vespa Search Engine, *Vespa Documentation*, 2023. [Online]. Available: <https://docs.vespa.ai/>. [Accessed: Nov. 19, 2024].
- [3] A. Bajaj and S. Sharma, "Leveraging Retrieval-Augmented Generation (RAG) for Information Retrieval Tasks," *arXiv preprint arXiv:2204.04578*, Apr. 2022. [Online]. Available: <https://arxiv.org/abs/2204.04578>. [Accessed: Nov. 19, 2024].
- [4] M. Chen and Y. Lin, "Towards Efficient Information Retrieval in Large-Scale Databases Using Postgres and Vespa," *arXiv preprint arXiv:2108.04667*, Aug. 2021. [Online]. Available: <https://arxiv.org/abs/2108.04667>. [Accessed: Nov. 19, 2024].
- [5] J. Lee and S. Park, "Applications of AI in Information Retrieval Systems," *arXiv preprint arXiv:2003.12607*, Mar. 2020. [Online]. Available: <https://arxiv.org/abs/2003.12607>. [Accessed: Nov. 19, 2024].
- [6] M. Lewis, et al., "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," *arXiv preprint arXiv:2005.11401*, May 2020. [Online]. Available: <https://arxiv.org/abs/2005.11401>. [Accessed: Nov. 19, 2024].
- [7] M. Sun, et al., "RAG: Retrieval-Augmented Generation for Long-Context Question Answering," *arXiv preprint arXiv:2104.07172*, Apr. 2021. [Online]. Available: <https://arxiv.org/abs/2104.07172>. [Accessed: Nov. 19, 2024].
- [8] V. Karpukhin, et al., "Dense Retriever for Question Answering with Retrieval-Augmented Generation," *arXiv preprint arXiv:2004.08908*, Apr. 2020. [Online]. Available: <https://arxiv.org/abs/2004.08908>. [Accessed: Nov. 19, 2024].