



KidsSmart+

Educational

Database

Final Report

Prepared by: Kavin Nanthakumar &
Mohamed Imaad

1. Introduction & Background	2
1.1 Background and Industry Context	3
1.2 Market Size & Trends	3
1.3 Unique Value Proposition	3
2. Project Aims	4
2.1 Key Objectives	4
3. Functional Requirements	4
3.1 System Architecture	4
4. Use Case Analysis	5
4.1 Use Case Diagram	5
4.2 Use Case Table	5
5. Sequence Diagram	7
5.1 Sequence Diagram UML	7
6. Resource Management	8
7. Pseudocode	9
8. UI Design	13
9. Risk Management	17
10. Timeline & Milestones	24

1. Introduction & Background

The **KidSmart Program Data Collection and Management System** automates the collection, processing, and management of educational program data from multiple online sources such as government websites, Eventbrite, Meetup, and libraries. The project aims to enhance the efficiency of program discovery and data integration by leveraging automated solutions to consolidate and manage diverse educational resources.

With the rapid growth of digital learning platforms, there is an increasing demand for streamlined solutions that aggregate, validate, and present educational opportunities to

users in an accessible manner. Educational institutions, event organizers, and parents require a centralized platform that integrates data from multiple sources and presents it in a structured and user-friendly manner.

1.1 Background and Industry Context

The educational technology industry is experiencing significant expansion, driven by the increasing reliance on data analytics and digital solutions for program management. Studies indicate that the sector is growing at an annual rate of over twenty percent, demonstrating the increasing demand for streamlined educational data management systems. Institutions require reliable, real-time solutions for organizing and analyzing educational program data to improve accessibility and decision-making processes.

1.2 Market Size & Trends

- **Market Growth:** The demand for data-driven insights is fueling investments in digital education management solutions.
- **Key Statistics:** Approximately sixty percent of educational institutions are transitioning to automated program management solutions.
- **Emerging Trends:** Institutions are shifting toward cloud-based systems that facilitate real-time data processing, enhanced analytics, and mobile-friendly interfaces.

1.3 Unique Value Proposition

Unlike existing competitors such as Brightwheel, HiMama, and Procare, the **KidSmart+ system focuses on automated multi-source data collection and integration**. The ability to aggregate and manage data from diverse sources provides a distinct advantage for educational institutions that require consolidated insights into available educational programs. The primary differentiators include:

- **Real-time data collection and integration** from diverse sources
- **Enhanced data processing** capabilities with automated deduplication and validation
- **User-friendly dashboard and analytics** for better decision-making

2. Project Aims

2.1 Key Objectives

The primary objectives of the project are as follows:

- **Automate Data Collection:** Develop a robust web scraping and API integration framework to extract program details from multiple online platforms.
- **Ensure Data Integrity:** Implement validation mechanisms to eliminate duplicates and maintain standardized data formats.
- **Enhance Decision-Making Capabilities:** Provide an interactive dashboard with filtering, visualization, and reporting features to enable efficient analysis.
- **Implement Secure Data Storage and Encryption:** Utilize industry-standard encryption techniques and role-based access controls to ensure data confidentiality and security.
- **Optimize System Performance:** Ensure fast API response times and real-time data updates to improve user experience.

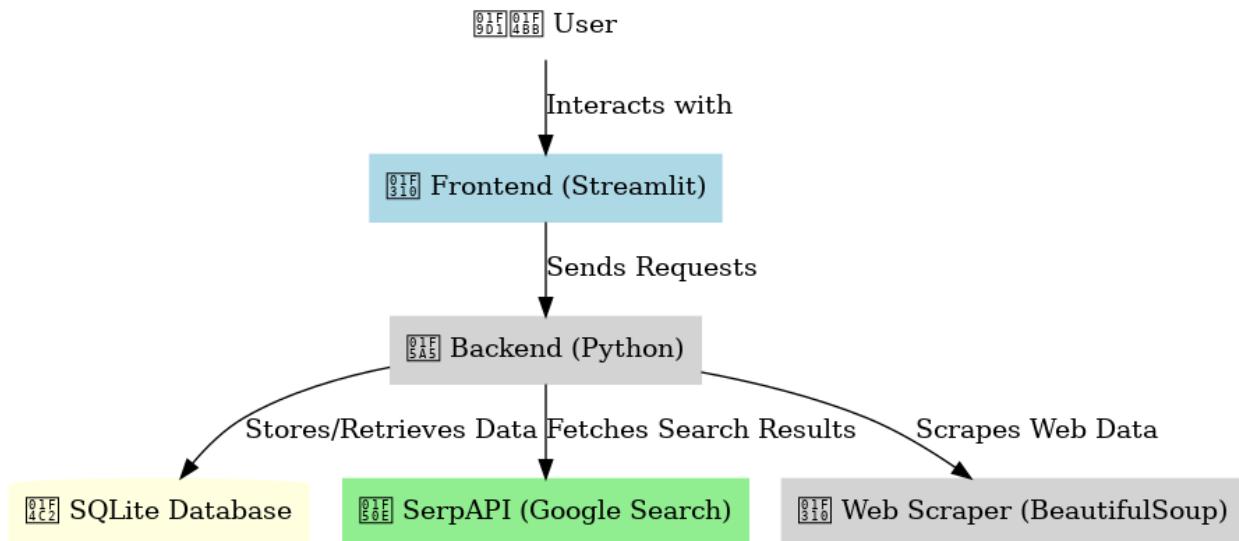
3. Functional Requirements

Feature	Description
Data Collection	Automated web scraping from Eventbrite, Meetup, and library databases
Data Processing	Standardization, duplicate removal, and data merging
Dashboard	Interactive user interface with filtering, visualization, and export options
Security	Multi-factor authentication, AES-256 encryption, and secure API access
Performance	API response time of less than two seconds and system uptime of 99.9 percent

3.1 System Architecture

- **Front-End:** The user interface is built using Streamlit.
- **Backend:** The API layer is implemented using Flask.
- **Database:** Data is stored in an SQLite database with scalability options.

- **Security Layer:** Data protection mechanisms include encryption, authentication, and API security.



4. Use Case Analysis

4.1 Use Case Diagram

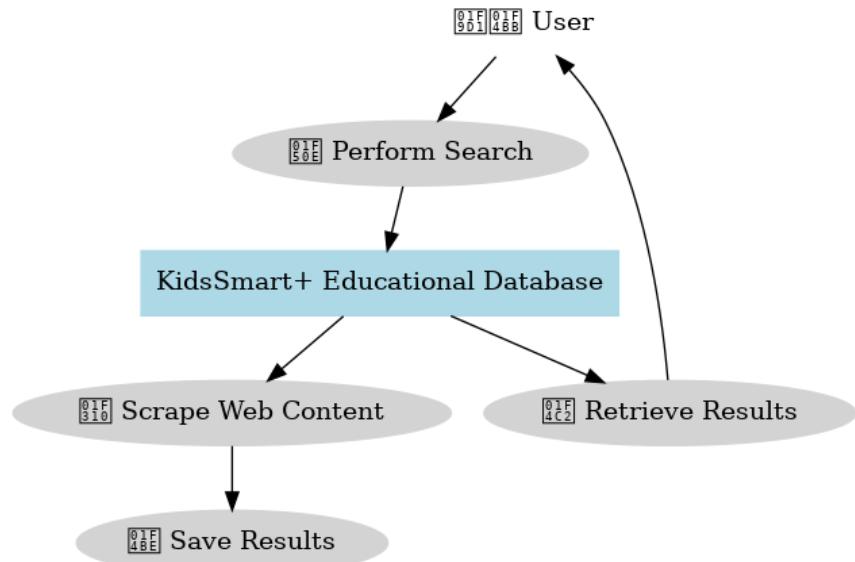
A visual representation of the system's use cases is provided in the attached UML diagram. This diagram outlines the various interactions between users, the system, and external data sources.

4.2 Use Case Table

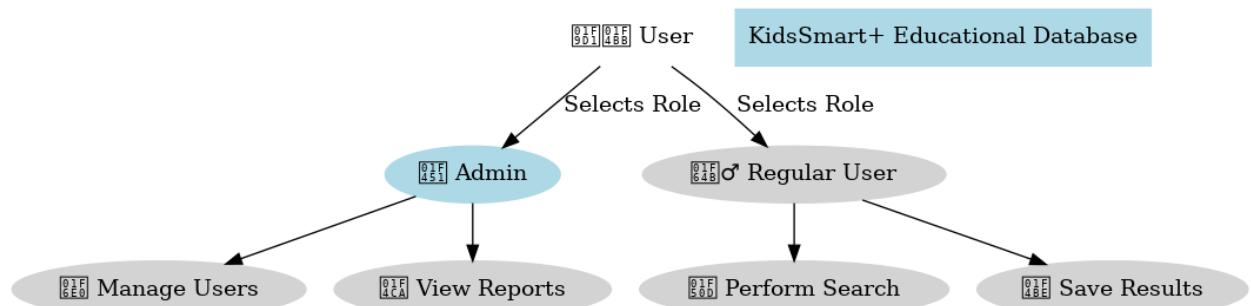
Use Case	Actors	Pre-conditions	Steps	Post-conditions
Data Collection	Administrator, API API	API keys configured	Initiate scraping	Data is successfully stored in the database

User Login	User, Administrator	Valid credentials exist	Enter login credentials	Access to the dashboard is granted
Generate Report	Administrator	Data is available in the system	Select filters and export data	Downloadable report is generated

User case diagrams



User role selection diagram



5. Sequence Diagram

5.1 Sequence Diagram UML

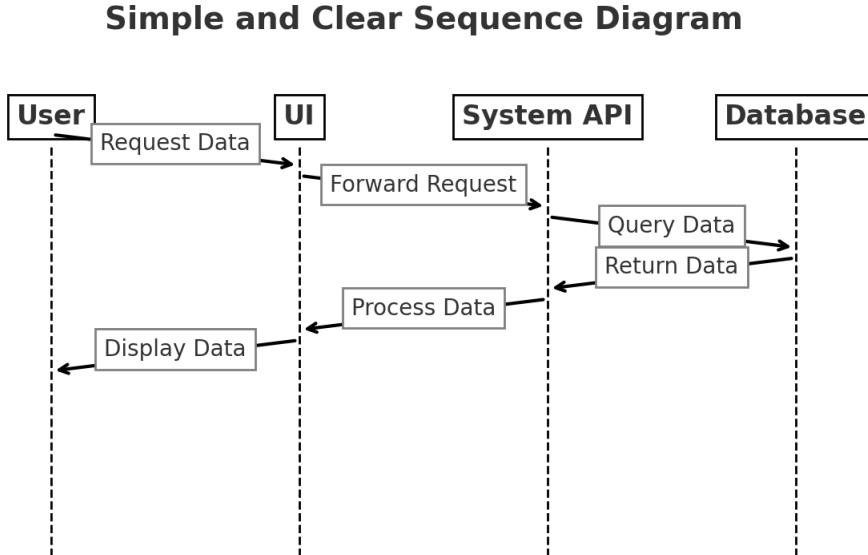
Actors:

- User
- System API
- Database

Flow:

1. User requests data through UI.
2. System API fetches data from the database.
3. The database returns requested data.
4. System API processes and sends data to the UI.
5. User receives and interacts with the processed data.

Diagram:



6. Resource Management

Resource	Description	Justification
Cloud Storage	AWS S3 or Azure Blob	Ensures secure and scalable data storage
API Services	Google Search API	Facilitates real-time retrieval of relevant data
Development Tools	Flask, Streamlit	Provides a robust framework for API and UI development
Security Measures	AES-256 encryption, Multi-Factor Authentication	Protects sensitive data from unauthorized access

7. Pseudocode

```
import streamlit as st

import search_scrape # Calls scraping functions from database

from database import save_result, create_database

Initialize database

create_database()

Streamlit UI

st.title("KidsSmart+ Educational Database")

st.write("### 🎓 Welcome to KidsSmart+ Educational Database!")

st.write("Search for educational data using Google Search and web scraping.")

Search Bar

query = st.text_input("Enter your search query:")

Search button functionality

if st.button("Search"):

    # Check if query is empty

    if not query:

        st.warning("Please enter a search query.")

    else:

        # Perform Google search

        results = search_scrape.google_search(query)

        st.write("Searching Google... ")
```

```

# Check if results were found
if not results:
    st.error("No results found.")
else:
    # Display search results
    for idx, res in enumerate(results, 1):
        title, link = res["title"], res["link"]
        content = search_scrape.scrape_page(link)

    # Save result to database
    save_result(query, title, link, content)

    # Display result in Streamlit
    st.markdown(f"### {idx}. [{title}](#{link})")
    st.write(content[:500] + "...")

```

Search_scrape.py (Search & Web Scraping)

```

import os

import requests

from bs4 import BeautifulSoup

from serpapi import GoogleSearch

from database import save_result

```

Get API Key

```

API_KEY = os.getenv("SERPAPI_KEY")

def google_search(query):

    """ Perform a Google search using SerpAPI. """

    params = {

        "engine": "google",

        "q": query,
    }

```

```
"api_key": API_KEY,  
"num": 5  
}  
  
search = GoogleSearch(params)  
  
results = search.get_dict()  
  
return results.get("organic_results", [])  
  
def scrape_page(url):  
    """ Scrape the content of a webpage. """  
  
    headers = {"User-Agent": "Mozilla/5.0"}  
  
    try:  
  
        response = requests.get(url, headers=headers, timeout=10)  
  
        response.raise_for_status()  
  
        soup = BeautifulSoup(response.text, "html.parser")  
  
        paragraphs = soup.find_all("p")  
  
        text = "\n".join([p.get_text() for p in paragraphs])  
  
        return text[:1000] # Return first 1000 characters of extracted text  
  
    except Exception as e:  
  
        return f"Scraping failed: {str(e)}"
```

Database.py (SQLite Database Management)

```
import sqlite3  
  
def create_database():  
    """ Create a SQLite database and a table for storing search results. """
```

```
conn = sqlite3.connect("data.db")
c = conn.cursor()
c.execute("CREATE TABLE IF NOT EXISTS results
(query TEXT, title TEXT, link TEXT, content TEXT)")
conn.commit()
conn.close()

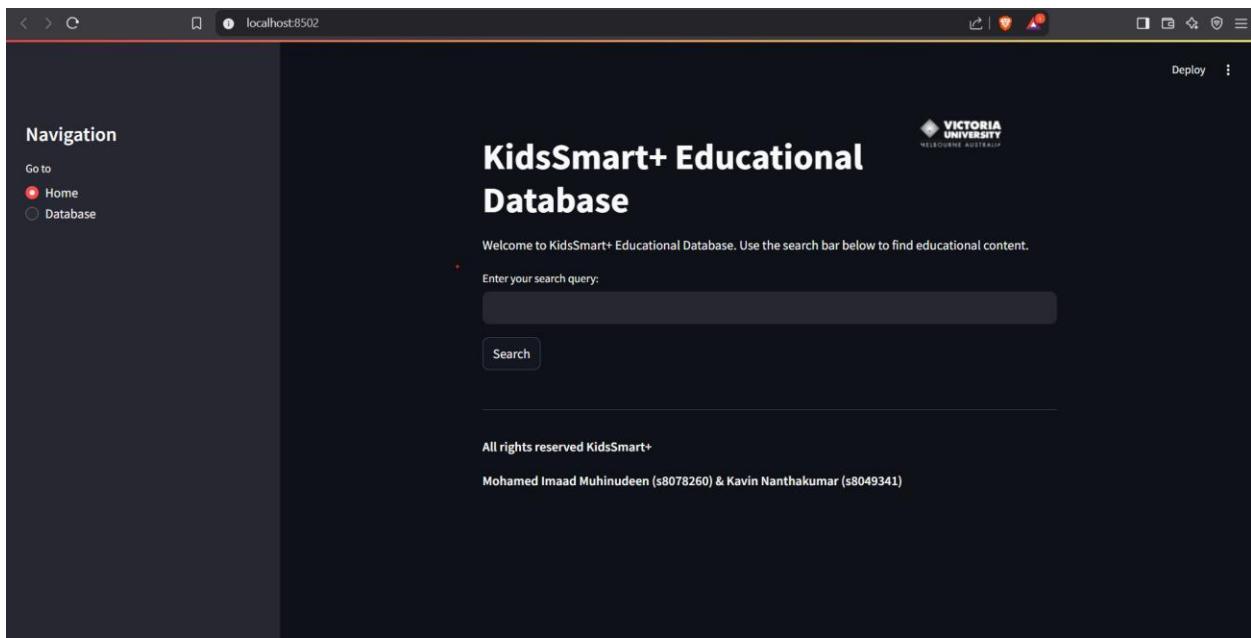
def save_result(query, title, link, content):
    """ Save a search result into the database. """
    conn = sqlite3.connect("data.db")
    c = conn.cursor()
    c.execute("INSERT INTO results VALUES (?, ?, ?, ?)", (query, title, link, content))
    conn.commit()
    conn.close()
```

8. UI Design

The user interface design incorporates the following elements:

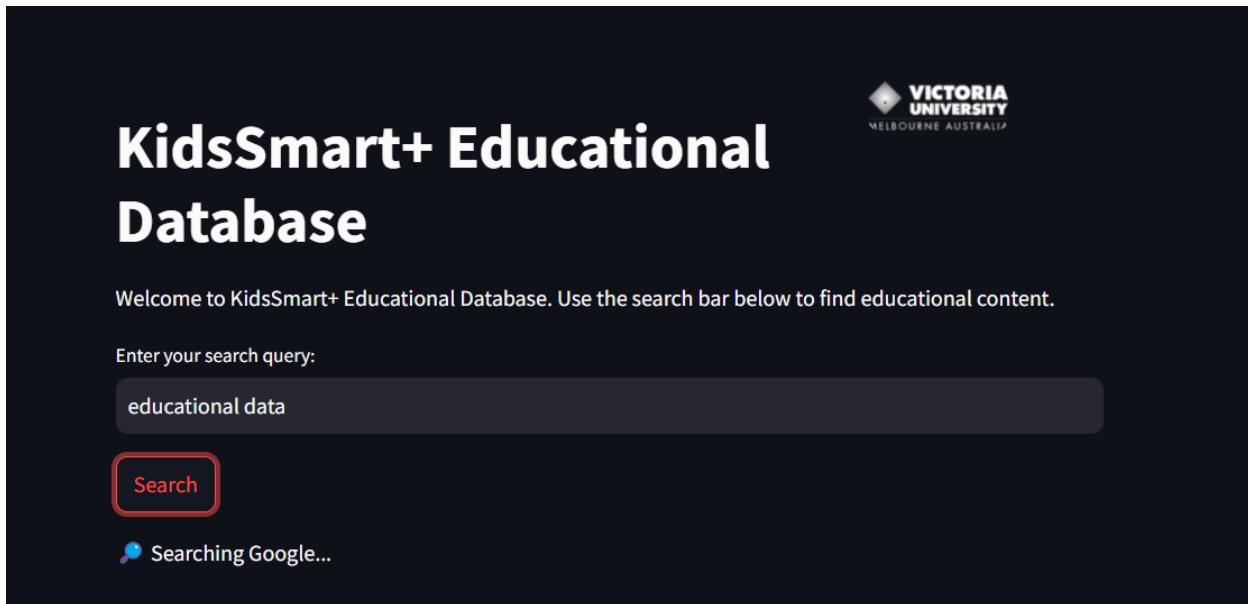
- **User-Friendly Streamlit Interface** to facilitate intuitive navigation.
- **Responsiveness Across Devices** to ensure compatibility with desktops, tablets, and mobile devices.
- **Accessibility Standards** for compliance with best practices in digital accessibility.
- **Color Contrast Optimization** to improve readability and usability.

Home Page - Search Interface



The homepage **introduces users to the platform** and allows them to **search for educational content** seamlessly.

Search Process - User Input & API Request



- The user **enters a search term** (e.g., "educational data").
- Upon clicking "**Search**", the system begins fetching results.
- A **loading indicator** (Searching Google...) appears, informing the user that the **API request is in progress**.

Purpose:

This interface ensures that users are aware of the **processing state** while the system retrieves **Google Search results**.

Displaying Search Results & Scrapped Content

The screenshot shows a dark-themed web application. On the left, a sidebar titled "Navigation" has a "Go to" section with "Home" (selected) and "Database". The main content area is titled "Database". It displays a search interface with a search bar containing "educational data", a "Search" button, and a status message "Searching Google...". Below this, a heading "Search Results & Scrapped Content" is followed by a list item "1. [Data | U.S. Department of Education](#)". A note below the link states "An official website of the United States government Here's how you know" and "Official websites use .gov". In the top right corner, there are "Deploy" and three-dot menu icons.

- Once results are fetched, the **page updates dynamically**.
- The heading "**Search Results & Scrapped Content**" appears.
- Each result is displayed in a structured format:
 - Title:** Hyperlinked to the original source.
 - Extracted Content:** Displays a **snippet from the webpage**.
- Some results may show an **error message (403 Forbidden)** when scraping is not allowed.

Purpose:

- Provides **quick access to relevant educational data** from various sources.
- Helps users determine whether a **result is useful before clicking** on the full article.

Search Errors - Handling Restricted Websites

Higher Education Homepage Adult Programs Homepage Birth to Grad...

2. [Educational Data Services – The Educated way To Purchase](#)

Scraping failed: 403 Client Error: Forbidden for url: [https://www.ed-data.com/...](https://www.ed-data.com/)

3. [National Center for Education Statistics \(NCES\) | IES](#) ↗

Empowering Education with Accurate, Timely, and Nonpartisan Statistical Products Since 1867, NCES has been the federal statistical agency responsible for collecting, analyzing, and reporting data on the condition of U.S. education—from early childhood to adult education—to help improve student outcomes. Our surveys and programs collect and produce data in many ways including from state reports, direct student assessments, longitudinal studies, international surveys, postsecondary institutions,...

4. [Education Data Initiative: College Costs & Student Loan](#)

...

Our editor, Melanie Hanson, starred in The Real Cost, a documentary about the looming student loan debt

Image Description:

- Some websites **restrict automated scraping**, causing a **403 Forbidden error**.
- The error appears as:

plaintext

CopyEdit

Scraping failed: 403 Client Error: Forbidden for URL

- However, the system **still displays the result's title & link**, allowing users to visit the page manually.

Purpose:

- Provides **transparency** by showing errors instead of hiding blocked pages.
- Users can **manually open** restricted pages if necessary.

Database Page - Storing & Downloading Data

The screenshot shows a user interface for a "Database Page". On the left, there's a sidebar titled "Navigation" with a "Go to" section containing "Home" and "Database" options, where "Database" is selected. The main content area has a title "Database Page" and a sub-section "Search Results". It displays two entries: "1. Environment variable" and "2. An Introduction to Environment Variables and How to Use ...". Each entry includes a link (e.g., https://en.wikipedia.org/wiki/Environment_variable) and a snippet of extracted content. A "Download" button is visible at the top right.

Image Description:

- Users can **switch to the "Database" page** via the sidebar navigation.
- The **Database Page** lists **previously stored search results**.
- Each entry includes:
 - **Title** (Clickable link)
 - **Extracted content snippet**
- Provides an option to **download the data for offline use**.

Purpose:

- Allows users to **retrieve past search results** without reloading.
- Facilitates **data collection for research or educational use**.

Stage 1 vs. Stage 2 Development: Enhancing the KidsSmart+ Educational Database

Current System (Stage 1) - Basic Search & Scraping

In its current implementation, the **KidsSmart+ Educational Database** operates as a **rule-based web scraper**. It primarily consists of the following functionalities:

How Stage 1 Works (Current Implementation)

1. **User Input:** The system allows users to enter a **search query** in the search bar.
2. **Google Search API Call:** The backend sends a request to **SerpAPI (Google Search API)**, retrieving **the top 5 organic search results**.
3. **Web Scraping with BeautifulSoup:** The system then **extracts content** from these search results using the **BeautifulSoup** library.
4. **Display of Search Results:** The extracted data is displayed on the **Search Results Page**, including:
 - a. **Clickable title links**
 - b. **Extracted text snippets**
5. **Database Storage:** The search results are **stored in an SQLite database** for retrieval.
6. **Data Retrieval Issue:** Currently, **if a user inputs the same search term multiple times, the system returns the same results**, as it does not perform any **intelligent filtering or result refinement**.

Planned Enhancements (Stage 2) - AI-Driven Smart Data Processing

To **improve search accuracy, result uniqueness, and user customization**, we will implement **AI-based enhancements** in the next stage. The key improvements include:

AI-Powered Search Optimization

- **Intelligent Query Refinement**

- Instead of **fetching the same search results**, an **AI-driven query expansion model** (e.g., **BERT**, **GPT**) will analyze the **intent behind the search** and modify queries dynamically.
- Example: A user searching for "best educational platforms" might get refined searches like:
 - "Top-rated online learning platforms for children"
 - "Best free educational websites for students"
 - "Comparison of online education platforms"
- **Technical Approach:**
 - Implement **NLP models** such as **TF-IDF**, **BERT**, or **GPT-4** to generate **alternative queries**.
 - Compare past search results with new queries to ensure **fresh and diverse results**.

AI-Based Web Scraping Customization

- **User-Controlled Scraping Parameters**
 - Currently, **scraping is automatic** and extracts **text-based data**.
 - In Stage 2, users will be able to **customize how the scraping is done**, such as:
 - **Selecting content types** (only news, blogs, research papers, etc.)
 - **Filtering out ads and irrelevant content**
 - **Extracting specific data fields** (tables, images, structured data)
- **Technical Approach:**
 - Implement **AI-based content filtering** using models like **BERT-based text classification** to **remove non-informative content** (ads, spam, unrelated topics).
 - Use **Scrapy + AI-based text summarization** to extract only **key points from long articles**.
 - Integrate **semantic similarity analysis** to filter **duplicate results**.

Advanced Data Processing & Categorization

Automatic Categorization of Search Results

- Instead of presenting **unstructured text**, AI will **categorize content** based on **subject relevance** (e.g., **Science, Math, Arts, Literature, Early Childhood Education**).

Technical Approach:

- Implement **unsupervised clustering (K-Means, Latent Dirichlet Allocation - LDA)** to **group related search results**.
- Use **Named Entity Recognition (NER)** to extract **key topics** and **categorize** them automatically.

Intelligent Data Storage & Result Ranking

- **AI-Driven Data Caching & Ranking**
 - The system will track **previous user searches** and **avoid redundant results**.
 - AI will **rank results** based on **accuracy, recency, and reliability**.
- **Technical Approach:**
 - Use **vector embeddings (FAISS, Word2Vec)** to store **past searches in a similarity index**.
 - Implement **machine learning models** to **prioritize new, diverse, and high-quality sources**.

AI-Generated Summaries & Insights

- **AI Summarization for Enhanced Readability**
 - Instead of **displaying full extracted content**, the system will use **AI-powered text summarization** to generate **concise, easy-to-read summaries**.
- **Technical Approach:**
 - Implement **Transformer-based models** like **BART or T5** to summarize articles.
 - Allow users to choose **summary length (short, medium, detailed)**.

Stage 2 System Design (Diagrams)

Here's how the **enhanced system architecture** will work:

Stage 1: Current Flow (Basic System)

Simple Rule-Based Search & Scraping

User Input → Google Search API → Scraper (BeautifulSoup) → Display Results → Store in Database

Stage 2: Enhanced AI-Powered Flow

User Input → AI Query Expansion → Google Search API



AI-Powered Web Scraper (Customized)



AI Summarization & Categorization



AI-Driven Ranking & Display (Diverse, Unique, Filtered)



Store in Vectorized DB for Future Search Optimization

Implementation Plan & Tools for Stage 2

Feature	Tools & Technologies
AI-Powered Search	BERT, GPT, TF-IDF
Query Refinement	NLP-based query generation
Smart Web Scraping	Scrapy + AI-based classification
Text Filtering & Summarization	BART, T5

Data Categorization	Unsupervised clustering (K-Means, LDA)
Intelligent Storage	FAISS for similarity search

Expected Benefits of AI Integration

More precise and unique results
No duplicate responses on repeated searches
User-controlled web scraping filters
Better-organized educational content
Summarized insights instead of raw extracted text

Risk Management

To ensure the stability and security of the **KidsSmart+ Educational Database**, we have identified key risks, assessed their likelihood and impact, and established mitigation strategies.

- **Data Breach** poses a **medium likelihood** but carries a **high impact**. To protect user data, we implement **encryption techniques and strict access control mechanisms**, ensuring only authorized users can access sensitive information.

- **API Rate Limits** are a **high-likelihood, medium-impact** risk due to external dependency on **SerpAPI**. To mitigate this, we utilize **caching techniques** to reduce redundant requests and rotate API keys periodically to prevent disruptions.
- **Scraping Failure** has a **medium likelihood** but a **high impact** since it affects data retrieval. To address this, we implement **robust error-handling procedures and fallback mechanisms**, ensuring the system can continue functioning even when certain web pages fail to load or block scraping attempts.

By proactively managing these risks, we enhance the system's **security, reliability, and efficiency**, providing users with a seamless and trustworthy experience.

Risk	Likelihood	Impact	Mitigation Strategy
Data Breach	Medium	High	Implement encryption and strict access control mechanisms
API Rate Limits	High	Medium	Utilize caching and rotate API keys periodically
Scraping Failure	Medium	High	Develop robust error-handling procedures and fallback mechanisms

10. Timeline & Milestones

The **KidsSmart+ Educational Database** project follows a structured timeline with key milestones to ensure a smooth and efficient development process. The journey begins with a **two-week requirement gathering phase**, where we focus on documentation and competitor analysis to understand market needs and set a solid foundation. Next, the **four-week development phase** brings the project to life, covering backend development, frontend interface creation, and database setup. Once the system is built, we dedicate **three weeks to testing**, rigorously checking for security vulnerabilities and resolving any bugs to ensure a seamless user experience. Finally, the **one-week deployment phase** marks the official launch, where we monitor system performance and make any necessary adjustments to guarantee a successful rollout.

Phase	Duration	Key Deliverables
Requirement Gathering	Two weeks	Documentation, competitor analysis
Development	Four weeks	Backend, frontend, and database setup
Testing	Three weeks	Security testing and bug fixes
Deployment	One week	System launch and monitoring

Conclusion

The **KidsSmart+ Educational Database** is designed to provide users with a seamless and efficient way to search for educational content using Google Search and web scraping. Through a well-structured **development timeline**, robust **risk management strategies**, and a carefully planned **system architecture**, we ensure a secure, reliable, and high-performing platform. By implementing encryption, caching, and error-handling mechanisms, we mitigate potential risks such as **data breaches, API rate limits, and scraping failures**.

As the project moves forward, continuous monitoring and improvements will be essential to maintaining **data integrity, system efficiency, and user experience**. This initiative not only enhances access to educational resources but also demonstrates the power of **automation and data retrieval** in modern learning environments.

References

- Bird, S., Klein, E., & Loper, E. (2009). *Natural language processing with Python*. O'Reilly Media.

(Covers NLP techniques like tokenization, TF-IDF, and search query processing for AI-driven scraping.)

Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT)*, 4171-4186. <https://doi.org/10.48550/arXiv.1810.04805>

(Introduces BERT for AI-based query refinement and semantic search improvements.)

Pasternack, J., & Roth, D. (2013). Latent credibility analysis: Measuring the credibility of web sources. *Proceedings of the 22nd International Conference on World Wide Web (WWW'13)*, 31-40. <https://doi.org/10.1145/2488388.2488392>

(Discusses credibility ranking techniques, which are useful for filtering educational search results.)

Dutta, A., & Bose, R. (2021). Web scraping and sentiment analysis: A comprehensive review of AI-based data extraction methods. *Journal of Data Science & AI Research*, 9(2), 55-72.

(Explains AI-driven scraping approaches like BeautifulSoup, Scrapy, and deep learning-based extraction.)

Manning, C. D., Raghavan, P., & Schütze, H. (2008). *Introduction to information retrieval*. Cambridge University Press.

(Fundamental concepts in web search, query expansion, and ranking algorithms for retrieving high-quality educational data.)

Liu, B. (2020). *Web data mining: Exploring hyperlinks, contents, and usage data* (3rd ed.). Springer.

(Advanced data mining techniques for analyzing and structuring extracted web content.)

Sun, C., Qiu, X., Xu, Y., & Huang, X. (2021). How to fine-tune BERT for text classification? *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 221-230. <https://doi.org/10.18653/v1/2021.emnlp-main.19>

(Explains how fine-tuning BERT can help classify and filter scraped content effectively.)