



UNIVERSITY OF
LEICESTER

School of Computing and Mathematical Sciences

CO7201 Individual Project

Final Report

Exploring HESA Data with Large Language Models for Dynamic Visualisation

Vladimirs Ribakovs

vr112@student.le.ac.uk

239062116

Project Supervisor: Heckel, Reiko (Prof.)

Principal Marker: Goes, Fabricio (Dr.)

Word Count: 12644

(excluding table contents, table labels, headings, references and appendices)

27/04/2025

DECLARATION

All sentences or passages quoted in this report, or computer code of any form whatsoever used and/or submitted at any stages, which are taken from other people's work have been specifically acknowledged by clear citation of the source, specifying author, work, date and page(s). Any part of my own written work, or software coding, which is substantially based upon other people's work, is duly accompanied by clear citation of the source, specifying author, work, date and page(s). I understand that failure to do this amounts to plagiarism and will be considered grounds for failure in this module and the degree examination as a whole.

Name: Vladimirs Ribakovs

Date: 14/05/25

Contents

1. Introduction	8
1.1 Research Question	8
1.2 Background and context	8
Introduction to HESA	8
Dynamic Visualisations.....	8
Why using LLMs for this Project	9
The Change of Original Approach	9
Personal Motivation.....	9
1.3 Problem Statement	10
Current Manual Approach of Processing HESA Data	10
HESA's Static Visualisations	10
Inflexibility of HESA Data	11
1.4 Aims and Objectives	11
1.5 Requirements	12
1.6 Scope and Limitations	13
Why Only University of Leicester Related Data.....	14
Only General HESA Data.....	14
1.7 Significance of the Project	15
1.8 Proposed Solution.....	16
2. Literature Review.....	16
2.1 Overview of HESA Data and Benchmarking in Higher Education.....	16

Detailed description of the Structure of HESA Data	16
2.2 Large Language Models (LLMs)	18
Overview of LLMs	18
The Gemini Model (1.5 Pro)	18
Comparison with Other LLMs (e.g., GPT-J)	19
2.3 LLM Usage in Data Analysis and Visualization	19
Natural Language to Visualization (NL2Vis)	19
2.4 Data Privacy and Safety with LLMs	19
Risks and Concerns with LLM APIs	19
Mitigation Strategies	20
2.5 Data Analysis Techniques	20
Ethical Considerations in Data Analysis	20
Quantitative Data Analysis in Academic Projects	20
3. Project Design and Planning	21
3.1 Conceptual Framework and Initial Approach	21
3.2 Risk Identification and Mitigation Strategies Implemented	21
3.3 Project Phases and Timeline	22
Explanation of Project Phases	23
4. System Specification	23
4.1 Data Flow Diagram	23
4.2 Technical Architecture Design	26
Technical Architecture Overview	26
Development Stack and Dependencies	28
Component Roles and Interactions	29
5. Implementation	30
5.1 Data Processing Pipeline Implementation	30
Raw Data Acquisition and Storage	30
Metadata Extraction and Indexing	31
Data Cleaning and Standardization	32
Transformation and Enhanced Metadata Generation	33
Rationale for Separate Raw and Clean Files	34
5.2 Query Processing and Interpretation	34
Overview of Natural Language Query Interpretation	34
Parameter Extraction, Tokenisation, and Normalization Techniques	35

LLM Integration Design	35
5.3 Visualization Generation	37
Chart Generation	37
Chart Types Implementation	37
Frontend Rendering	37
Interactive Visualization Controls.....	37
Real time Chart Regeneration	37
Visualization Type Recommendations	38
5.4 UI Implementation and Interactivity	38
Dashboard Interface	38
Dataset Details View.....	38
Dynamic Dataset Previews	39
Responsive Design	39
5.5 Feedback System	40
Feedback Collection	40
Feedback Submission.....	40
5.6 Caching and Performance Optimization	40
Query Cache Implementation.....	40
5.7 Codebase Structure and Endpoints	41
Codebase Overview	41
Key Endpoints.....	42
5.8 System Setup and Running Instructions.....	42
6. Demonstration	43
6.1 Using the AI Powered Query Interface	43
Submitting a Natural Language Query	43
6.2 Query Interface.....	44
Reviewing Query Analysis Results.....	44
Working with Dataset Details and Visualization	45
Generating and Interacting with Visualizations	46
6.3 Providing Feedback on Results.....	49
Submitting Feedback on Query and Visualization	49
6.4 Uploading Custom Data	50
Navigating to the Data Upload Section.....	50
Uploading and Processing Data Files	50

6.5 Exploring Documentation	52
Accessing the Documentation Page	52
6.6 Complete User Journey Example	57
From Query to Insight: Full Workflow.....	57
7. Testing and Evaluation	61
7.1 Testing Methodologies	61
Unit Testing	61
Integration Testing.....	63
7.2 Performance Metrics and Results.....	66
Query Response Time	66
Accuracy of Query Interpretation	68
8. Discussion	74
8.1 Analysis of Findings	74
8.2 Explaining LLM and Method Comparisons	74
Explanation of the Initial GPT-J Implementation and Reasons for its Failure	74
Comparison of Approaches (Regex, GPT-J, Gemini).....	74
8.3 Project Process Evolution and Timeline	76
Impact of Process Changes	77
Lessons Learned from Process Evolution	78
8.4 Achievement of Requirements.....	78
8.5 Personal Learning and Reflection	82
8.6 Critical Evaluation	83
Research	83
Design	83
Implementation	83
8.7 Discussion of the Number of Requests and Associated Costs.....	84
9. Conclusion and Future Work	84
9.1 Summary of Project Achievements	84
9.2 Conclusions	85
9.3 Future Work and Recommendations	85
10. References	85
11. Appendix.....	87
A: Pandas-Based Data Processing Pipeline	87
A.1 Data Transformation for Visualization.....	87

A.2 Multi-format Data Loading in Correlation Engine	88
A.3 Chart Data Preparation with Pandas Pivot	89
B: Unit Testing	90
B.1, B.2 Data Processing and LLM Integration	90
B.3 Visualization.....	92
C: Integration tests	96
C.1 Query to Visualization Flow	96
C.2 Mock AI Client Integration	100
D: Data Processing Pipeline Implementation.....	103
D.1 Data Ingestion Process	103
D.2 Manual Update Process.....	105
D.3 Metadata Extraction.....	107
D.4 Indexing	108
D.5 Cleaning Process	112
D.6 CSVProcessor	115
E: Query Processing and Interpretation.....	117
E.1 Transform free text into query parameters	117
E.2 Date Interpretation	121
F: Visualization Generation.....	123
F.1 Chart Generation.....	123
F.2 Handles different chart type requests.....	127
F.3 Gets AI recommendations for chart types	130
F.4 Dynamically changing the visualization.....	131
G: Front-End Rendering.....	133
G.1 Rendering Visualization	133
G.2 Interaction Control.....	134
G.3 Changing Chart Type	137
H: Interactive Visualization Controls.....	139
H.1 Back-end connection.....	139
H.2 Chart Type Explanation	140
I: Dashboard Interface	141
I.1 Display Query Results	141
I.2 Reviewable Dataset.....	143
I.3 Manage loading states	145

I.4 Query Input	146
I.5: Exporting data	148
I.6 Prepare Data for Clean Display	149
I.7 Dynamic dataset preview toggle	150
J: Feedback System	151
K: Caching	154
L: Gemini Client	158
M: Timing Queries	163
M.1 Recording the time of query processing.....	163
M.2 Recorded time of processed queries	165

1. Introduction

1.1 Research Question

How can Large Language Models (LLMs) be used to effectively and dynamically process and visualise HESA open-source data to enable the University of Leicester to more efficiently understand its performance against peer institutions?

1.2 Background and context

Introduction to HESA

- Higher Education Statistics Agency (HESA) is the primary and official open-source data collection agency for UK higher education.
- It provides a massive database consisting of millions of records covering various aspects such as institutional performance, funding and benchmarking.
- HESA's mandate is to collect, analyse, and disseminate reliable statistical information about UK higher education.
- The types of data provided by HESA and used in this project include: [5]
 - Student enrolment data (by level, mode, demographic)
 - Staff information (qualifications, contracts, demographics)
 - Financial data (income, expenditure)
 - Facilities information (accommodation, resources)
 - Course and program statistics.

These datasets enable institutions to benchmark against peer institutions and to identify areas for improvement, and track progress over time. [5]

HESA data serves multiple stakeholders:

- Universities use it for strategic planning, performance reviews, and competitive analysis [5]
- Government bodies use it for policy development and funding allocation [5]
- Researchers analyse trends in higher education access, outcomes, and equity [5]

Dynamic Visualisations

Dynamic visualisations are interactive and responsive visual tools used in this project to visualise HESA data. They go beyond static charts by allowing user interaction, real time data exploration, and customization [21]. This project implements it using Chart.js and it allows users to:

- Change the type of visualization on demand.
- Filter data by time periods or institutions.
- Zoom and explore details within charts.
- Compare multiple metrics simultaneously.

This approach transforms raw HESA data into results that can be explored rather than just viewed. Dynamic visualisations are data that is interactively represented based on the user's submitted

query. It can reveal differently even on the same dataset because it depends on the user's intent. They improve insights by:

- Enabling exploration of underlying patterns.
- Facilitating comparison between different data segments.
- Allowing users to focus on specific areas of interest.
- Supporting different levels of data granularity.

Why using LLMs for this Project

Manual analysis and processing a large volume of complex data is time consuming and difficult, so LLMs offer a solution to this tedious process. In this project, an LLM is used at different stages of the system, such as data processing, analysis, extraction, preparation, recommendations, and even visualization. Even though HESA provides ways to visualise data, but those visualisations are limited and static. Using an LLM can help to work with this data in a more efficient and intuitive way. The Gemini 1.5 Pro was chosen as the LLM for this project because its excellent understanding of natural language queries (NLQ) and translating them into structured operations. This feature addresses challenges in data analytics and allows users without any technical knowledge in this area to access and work with complex datasets.

LLMs offer several advantages for query interpretation:

- Understand context and intent beyond keyword matching
- Handle variations in how questions are phrased
- Recognize domain specific terminology (e.g. "academic years")
- To understand relationships between entities such as institutions, metrics or time periods

Traditional data query methods require knowledge of query languages or tools, understanding of dataset structure, and multiple steps to filter and extract relevant data. LLMs eliminate these barriers by allowing users to simply ask questions in plain English, making data more accessible and reducing the time from question to answer.

The Change of Original Approach

Originally, when the project was just proposed, the GPT-J was selected as the LLM for the project. After exploring GPT-J for some time, I realised that it did not meet the specific requirements of the project, so I had to switch to a more robust and advanced model, Gemini 1.5 Pro. See the section 8.2 on why exactly this approach failed.

Personal Motivation

My personal motivation is to undertake a project that uses an LLM to solve a real world problem. Also, I am fascinated by how much AI (Artificial Intelligence) has improved in the past couple of years, so I want to do more projects using it to get more experience with it. I consider this project a big step toward getting more advanced in working with AI.

1.3 Problem Statement

Current Manual Approach of Processing HESA Data

The current approach to understand the University of Leicester's performance relative to peer institutions using HESA data is manual, tedious, error-prone, and costly [5]. This process requires dedicated staff to spend significant time manually researching and manipulating HESA data.

Even though there is no exact open-source information on how the University of Leicester conducts these searches, based on related information and logical thinking, it can be how it proceeds.

As mentioned above, there is no specific information on how this process is conducted. The University of Leicester does not even say anything about this data being used from the HESA source, at least it is not publicly available.

The University's privacy notice indicates they conduct "statistical analysis" using student data, which suggests that HESA submissions are being used. This strongly suggests that university employees, likely within departments like planning, academic affairs, or institutional research, are involved in this process. [13]

Looking at the job applications that the University of Leicester posts for roles that are closely related to "statistical analyses," I found that they mention the following tools: Power BI, Microsoft Office, or Tableau [14]. According to the discovered information, I can assume that they undertake the following steps to perform current analyses of HESA data:

- Downloading multiple CSV files.
- Cleaning and standardizing data.
- Merging datasets from different years.
- Creating visualizations in tools like Excel, Tableau or Power BI.

Visualisations created through tools like Excel, Tableau, or Power BI further compound these challenges because of inflexible formats. If a research question is changed or followed up by another question, analysts would likely have to restart the entire process – extracting data, reformatting columns, recategorizing values, and generating new charts, which results in significant reporting delays.

Manual methods suffer from:

- Inconsistent data cleaning approaches
- Difficulty in tracking data provenance
- Time intensive repetitive tasks
- Errors in data transformation and aggregation
- Challenges in maintaining and creating visualizations

HESA's Static Visualisations

- Static charts provided by HESA are limited and cannot adapt to users' specific questions.

They cannot show relationships between multiple variables, provide different levels of detail, and support exploration beyond the initial visualization.

- Predefined visualizations cannot be easily changed, so users are restricted to the data and format that HESA provides.
- Users cannot even remove specific data points or change visual attributes like colours or titles.
- HESA offers a variety of charts like line, bar, pie or scatter, but the user cannot decide for themselves what chart they want. Users are forced to use only existing predefined charts.
- Having all these points combined, the fact that users cannot decide what chart to have and cannot change existing charts, it all leads to limited analyses. At best, charts provide a dropdown menu that allows sorting by value, like country, as shown in Figure 1

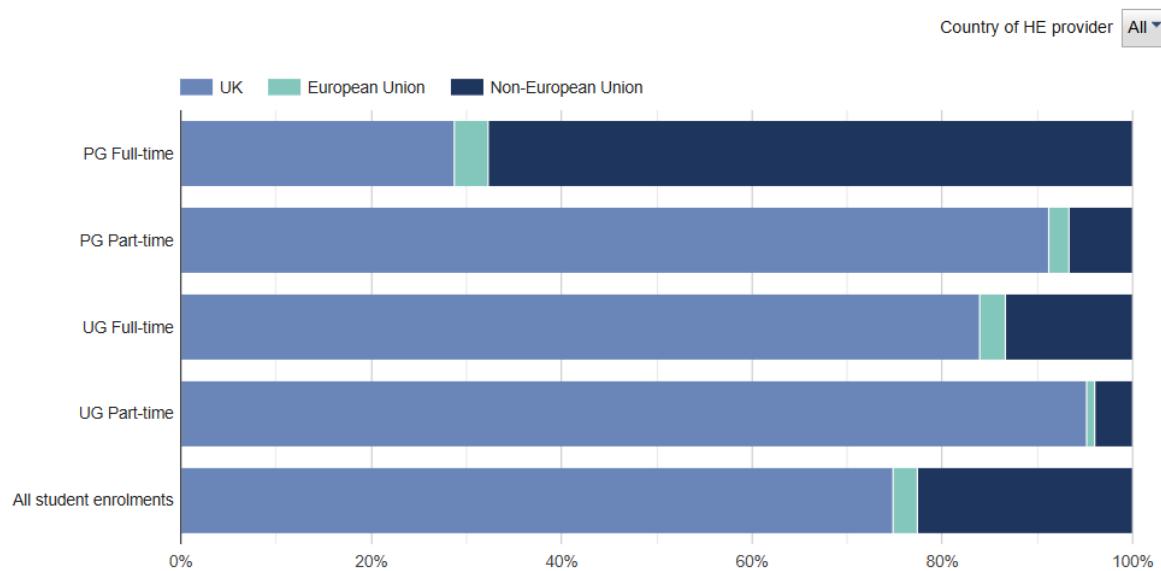


Figure 1: HESA, HE student enrolments by level of study, mode of study and permanent address, Academic year 2023/24.

Inflexibility of HESA Data

HESA does not provide a way to include comparison by Mission Group (Russell Group, University Alliance, Million+), so when an analyst needs to benchmark data against a Mission Group, the data needs to be manually found for each institution in a particular Mission Group. For example, when benchmarking groups shifted to include more institutions from Russell Group, existing reports required manual reconstruction. The inability to dynamically filter, recategorize, or pivot data dimensions without rebuilding entire visualization sets severely limits the exploration of emerging patterns and trends.

1.4 Aims and Objectives

The University of Leicester needs a solution that:

- Adapt to evolving analytical questions

- Provide insights without technical barriers
- Support strategic decision making with relevant comparisons
- Reduce the administrative burden of data analysis.

The project aims to address these limitations by creating a dynamic, user-friendly dashboard powered by an LLM Gemini 1.5 Pro that provides more efficient data comparison, extraction and analysis.

The purpose of this project is to make it easier for the University of Leicester to understand its organisational performance, structure, strengths, and weaknesses against peer institutions by leveraging HESA open-source data. The project aims to transform HESA data interaction through:

- Intelligent data processing that understands academic contexts
- Natural language query interpretation that captures institutional benchmarking needs
- Dynamic visualization that reveals meaningful patterns and comparisons.

The core goal is to reduce the technical and time barriers to analyse HESA data for the University of Leicester. This will allow them to gain powerful insights, ultimately saving time and potentially reducing costs.

1.5 Requirements

Essential:

- **Natural Language Query Interpretation:** The system understands users' free form queries and retrieves and processes data based on that.
- **Integrating and Processing Data:** Extracting and preparing clean HESA data and making it usable for comparing metric performance.
- **Interactive Dashboard:** It must dynamically display results in different formats like tables, charts, diagrams or summary reports.
- **Comparative Analysis:** A feature that must allow grouping and benchmarking universities by different categories like Mission Groups.
- **Data Export Capabilities:** Allow users to download generated reports and data as CSV, PDF or Excel files.

Recommended:

- **Automated Data Retrieval:** Enable automatic updates of the dashboard with new data from HESA when it is available.
- **Using AI To Improve Analysis:** LLM should provide insights, trend identification, and recommendations if user asks for that.
- **Query Building Guide:** Provide a guide or pre-built examples for users who are unfamiliar with the free text queries.
- **Advanced Visualisation:** Support interactive graphs, trend projections that offer visual comparison of performance metrics.

- **Historical Data Tracking:** Tracking data changes over time to see how performance was changing relative to peers.
- **Caching for Frequent Queries:** Implement caching to perform repeated queries more efficient that will result in better performance.
- **LLM Interpretation Feedback:** Allow users to see how their query was interpreted by the system (Pandas) to improved data transparency and remove bias factor.
- **Queries Logs:** Log out user queries to improve future prompts tuning.
- **Data Overload protection:** Implement a safeguard that prevents outputs of big queries to potentially crash the system.
- **CSV File Validation:** Make a mechanism to automatically sanitize new raw HESA data by checking for correct format and size anomalies and then prepare it for the data processing.

Optional:

- **Predictive Modelling:** Using historical data and machine learning (ML) to predict future performance metrics
- **Chatbot Integration:** Include an AI chatbot assistant to help users to refine their queries dynamically, in real time.
- **API for External Use:** Create an API so that external and third-party applications can use dashboard's functionality
- **Collect feedback:** Implement mechanisms to collect users' feedback through the surveys to continuously improve the system
- **Data Storage Scalability:** If HESA data volume grows beyond local CSV storage it will be moved to a more scalable database solution such as PostgreSQL or SQLite.
- **Multiuser Concurrency:** Allow users to submit multiple queries simultaneously by using asynchronous job queues such as Celery and session management.
- **Deployment:** Develop a deployment plan starting with local deployment and extending to containerization such as Docker and orchestration like Kubernetes for future scalability.

1.6 Scope and Limitations

This project focuses on developing a tool for exploring HESA open-source general data in CSV format that is relevant to the University of Leicester.

The project covers:

- Processing and indexing selected HESA CSV datasets.
- Natural language querying of these datasets.
- Automatic visualization generation using Chart.js.
- Comparing and benchmarking institutions that are particularly focused on Mission Groups.
- Provide functionality to export the data.

It does not cover:

- Real time data goes directly from HESA.
- Custom data upload from non-HESA sources.

- Comprehensive processing of all available HESA data.

Why Only University of Leicester Related Data

- The project intentionally focuses on Leicester relevant datasets to provide meaningful benchmarking rather than attempting to cover the entire HESA data.
- Only data that has records about the University of Leicester was downloaded and used for the project because the primary purpose is comparison against the University of Leicester.
- This focused approach provides more relevant and actionable insights for University of Leicester stakeholders. It improves performance by reducing the dataset size, enables better benchmarking and allows for more targeted visualizations.

Only General HESA Data

HESA provides two types of files:

- extremely detailed files containing about 500,000 – 700,000 rows of data on average
- general CSV files that have around 100 – 400 rows of data on average.

The detailed files cause significant challenges for the LLM processing due to token limits restricting how much data can be analysed at once. The restrictions would affect any modern open-source LLM, so this is not a drawback of Gemini 1.5 Pro. Other challenges with detailed data are increased processing time with dataset size, memory constraints for complex analyses, excessive API costs, and often contain redundant or overly granular information. Therefore, this project focuses only on using general HESA data files.

This project uses a total of 392 general HESA CSV files. These general files still provide comprehensive data about various aspects but are less detailed which makes them more suitable for the scope of this project. The project addresses potential challenges of high data volume by implementing selective data processing and indexing functionality. Data is filtered first by academic year(s), then by relevance to the query, and typically only about 10-20 files relevant to a query, where each file has approximately 300 lines, is sent to the LLM for analysis.

This approach is critical for managing processing cost and time. For example, if the system sends 10 – 20 of the detailed files where each file averaging around 600,000 lines of text, this will require processing a significantly larger volume of data. Based on current usage, approximately 227 request cost about £1 of free credit using general files. (see Figure 8)

Note:

$$\text{Number of requests per £1} = 882 \text{ (Total number of request)} / 3.88 \text{ (Total cost of total number of request)} \approx 227 \text{ (see Figure 8)}$$

Using detailed files under similar query conditions would be estimated to cost at least 2000 times more due to the vastly increased data volume processed per query. This highlights the main drawback of using detailed files for dynamic analysis with the current approach and budget constraints. Therefore, this project focuses on using general data to allow for dynamic analysis using the LLM.

Note:

The estimate of “at least 2000 times more” cost when using detailed files is based on the difference in data volume processed per query. The calculation is as follows:

- *Average lines per query with general files: $(10 \text{ files} * 300 \text{ lines/file} + 20 \text{ files} * 300 \text{ lines/file}) / 2 = 4,500 \text{ lines}$*
- *Average lines per query with detailed files: $(10 \text{ files} * 600,000 \text{ lines/file} + 20 \text{ files} * 600,000 \text{ lines/file}) / 2 = 9,000,000 \text{ lines}$*
- *Ratio of detailed to general data volume: $9,000,000 \text{ lines} / 4,500 \text{ lines} = 2000 \text{ times}$*

More data per query means much higher token usage, which leads to a significantly higher cost per API call.

1.7 Significance of the Project

This project aims to deliver a significant impact on decision making processes for the University of Leicester. It expects to allow for more flexible and accessible work with the HESA data

The project significantly impacts university decision making by:

- Providing faster access to comparative insights.
- Enabling the University of Leicester (UOL) to explore data without technical barriers.
- Helping the OUL use data to make better decisions
- Supporting real time responses to data-related questions.

These features help with strategic planning by showing patterns and trends that static reports might miss and allowing interactive filtering for deeper analysis.

The project also contributes to automating data visualization and analysis in higher education analytics by:

- Demonstrating effective LLM integration for educational data analysis.
- Creating a reusable framework for natural language data exploration.
- Establishing patterns for automated chart recommendation and generation.
- Providing a feedback loop for continuous improvement.

Long term benefits:

- Reduced administrative burden for data analysis
- More consistent and reproducible data visualization
- Institutional knowledge capture through query patterns
- Improved data literacy across the organization.

1.8 Proposed Solution

The dynamic dashboard is intended to simplify the process of comparing university performance, saving time and potentially reducing costs. It will allow users to query data in plain English and receive dynamic charts, diagrams, and summary reports. Moreover, the dashboard aims to enable benchmarking of the University of Leicester against Mission Groups and other institutions. It would support relatively easy updates as new HESA data becomes available (through the implemented manual upload feature). The motivation is to make comparisons between universities more efficient and accurate. It will allow the University of Leicester to quickly identify its strengths and weaknesses, potentially leading to better resource allocation and improved outcomes.

2. Literature Review

2.1 Overview of HESA Data and Benchmarking in Higher Education

Detailed description of the Structure of HESA Data

HESA data follows a structured format organized by data stream, and it has consistent patterns across different types of information. The system is implemented to understand and process this structure through data processing and indexing functionality. Key data streams relevant to this project include:

- **Students:** Covering entrants, qualifiers, and total numbers [5]. This includes information on:
 - Who is studying in Higher Education (HE),
 - Where they come from, where they study
 - What they study
 - Progression rates and qualifications.
 - Enrolment statistics by level, mode, and year,
 - Demographic information (age, gender, ethnicity),
 - Progression and completion rates,
 - Qualifiers by degree classification
 - Subject area breakdowns.
- **Staff:** Covering demographic and contract characteristics [5]. This includes data on:
 - Who is working in HE
 - What areas they work in
 - What their employment conditions are
 - What their salaries are
 - Where they work
 - Where they come from and go to
 - Employment conditions (contract types, full/part-time)
 - Salary ranges and pay gaps
 - Academic qualifications and specialties

- Demographic profiles
 - Staff-student ratios
- **Graduates:** Covering survey results showing the activities of recent graduates [5]. This includes data on:
 - Summary statistics on employment outcomes
 - Further study progression
 - Salary information,
 - Industry sectors
 - Graduate satisfaction metrics
- **Finances:** Covering income, expenditure, and financial statements of HE providers [5]. This includes data on:
 - Income and expenditure of HE providers
 - Other financial statements
 - Key Financial Indicators (KFI)
 - Financial data encompasses income sources (tuition, research, commercial)
 - Expenditure categories
 - Balance sheet information
 - Financial sustainability metrics
 - Capital investment patterns
- **Business and Community Interaction:** Covering intellectual property, services, and engagement[5]. This area includes data on:
 - Knowledge exchange activities
 - Patent and IP generation
 - Commercial partnerships
 - Social and community engagement
 - Continuing professional development.
- **Estates Management:** Covering environmental information [5]. This includes:
 - Building condition and space utilization
 - Energy consumption and carbon emissions
 - Waste management and recycling
 - Water usage
 - Sustainability initiatives.

HESA data also appears in multiple formats:

- Tabular data (CSV files with structured rows and columns)

- Time series (sequential data across academic years)
- Categorical breakdowns (data segmented by demographics or characteristics)
- Hierarchical relationships (institutional, departmental, subject levels)
- Geospatial information (regional and institutional locations).

HESA data is the official source for UK higher education information. It provides a standard way to compare institutions and evaluate their performance. This data is essential for university planning, government funding and policy, ensuring quality, following rules, and showing the public how the sector is doing [17].

Reports show that HESA data is key for universities to see how they compare to others. However, there are still problems with getting the data easily, to get it on time, and sharing it to make benchmarking most effective. There are certain rules and difficulties in using numbers to compare institutions [18].

2.2 Large Language Models (LLMs)

Overview of LLMs

LLMs are strong in data projects like this, but they also have downsides [2],[3].

Advantages:

- Understand natural language
- Handle different question styles
- Recognise complex data relationships
- Create an easy-to-read explanation
- Adapt to new topics

Downsides:

- Limits on how much text they can process at once
- Costs for using APIs
- Delays in getting responses
- Sometimes give wrong information
- Need careful setup for prompts.

The Gemini Model (1.5 Pro)

Gemini 1.5 Pro offers significant advantages: very large (undisclosed) parameters, a context window of up to 1 million+ tokens, native multimodality (text, audio, images, video), and accessed via API with associated costs but limited customization compared to self-hosted models [23].

Comparison with Other LLMs (e.g., GPT-J)

Feature	GPT-J	Gemini 1.5 Pro
Parameters	Fewer (Earlier generation)	Significantly More (Undisclosed, but larger)
Context Window	Limited	Vastly Larger (Up to 1M+ tokens)
Multimodality	Primarily Text	Native Multimodality (Text, Audio, Images, Video)
Cost Structure	Free (if self-hosted), Infrastructure cost	API-based (Associated costs, free credits used)
Customisation	More extensive (if self-hosted)	Limited (API-based)
Flexibility	Less flexible	Greater flexibility (Handles diverse data types)
Ease of Use	Requires local setup/hosting	Easy to use (API access)
Performance	Slower processing, higher error rates	Faster processing, better accuracy
Specialised Knowledge	Limited understanding of HE terminology	Superior understanding of HE terminology

Table 1: Comparison of LLM Features (GPT-J vs. Gemini 1.5 Pro) [1], [23]

2.3 LLM Usage in Data Analysis and Visualization

Natural Language to Visualization (NL2Vis)

Natural Language Interfaces for querying and visualising tabular data (NL2Vis) is an active area of research. Recent surveys highlight the significant role of LLMs in this domain [21]. LLMs allow users to interact with data using natural language by converting free form queries into structured formats [22]. This process involves complexities:

- Understanding user intent
- Accurately mapping natural language terms to the underlying data schema
- Clearing up confusion [21].

Handling these challenges allows for automatic data retrieval and chart generation based on user intent [22].

2.4 Data Privacy and Safety with LLMs

Risks and Concerns with LLM APIs

Data safety and privacy are critical issues when working on an application that uses an LLM. This is especially the case when using external APIs. Besides having the general risk of data leakage there are also specific vulnerabilities that exist in the context of LLM applications. However, the HESA data that is being used for this project is publicly available, so it does not have any privacy problems, but this project still implements the safety measures. This could be used in case if the University of Leicester would decide to adapt this project to use sensitive information and create visualizations based on that:

- Prompt Injection: malicious input can manipulate the LLM's behaviour

- Insecure Output Handling: which can lead to exposure of sensitive information processed by the LLM

This is important to understand and follow these specific security measures when building robust and trustworthy LLM powered systems [19].

Mitigation Strategies

Applying threat modelling to the system that uses LLM can help to identify vulnerabilities related to sensitive data handling [20]. The idea for dealing with data privacy risks always applies even if the situation changes. Below principles are focused on protecting user data and maintaining confidentiality [20]. Key mitigation strategies include:

- Keeping Raw Data Local: A fundamental strategy implemented is processing and storing raw HESA data locally. This ensures that sensitive or detailed information is not transmitted to the external LLM API [20].
- Selective Data Processing: Only necessary and anonymized data or metadata is sent to the LLM. Specific tasks like query interpretation or visualization configuration helps with minimizing the exposure of detailed HESA records.
- Robust Input Validation and Sanitization: Implementing careful validation and sanitization of user inputs helps to mitigate risks like Prompt Injection [19].
- Secure API Key Management: Ensuring API keys are stored and used securely to prevent unauthorized access to the LLM service.
- Monitoring and Logging: Implementing logging of LLM interactions can help in identifying and responding to potential security incidents.

2.5 Data Analysis Techniques

Ethical Considerations in Data Analysis

It's very important to consider ethics when collecting, processing, or just analysing data today [15]. HESA data does not publish data that raises ethical questions, but it is important to mention that this project could be used to work with sensitive data if the University of Leicester require it. So, mentioning the importance of ethical concerns would help to avoid the pitfalls with privacy and bias.

Quantitative Data Analysis in Academic Projects

Quantitative data analysis is a basic part of academic projects. There are studies that looked at dissertations and found common problems in how methods and statistics were used [16]. This shows how important it is to use the right research methods and make sure data is processed and

tested accurately. This helps to get reliable and correct results, which makes the analysis of HESA data more trustworthy.

3. Project Design and Planning

3.1 Conceptual Framework and Initial Approach

The project started with a clear focus on improving the way HESA data is being accessed. This is proposed to be done using AI powered natural language processing (NLP). The conceptual structure is based on these three main pillars:

- Data acquisition and preprocessing
- Query interpretation using AI
- Dynamic visualization generation

The original plan was to use AI model GPT-J because of its wide popularity, accessibility and its availability at no cost. Since the model supports self-hosting, it incurs no operational cost to run it on my own hardware. It allows you to maintain control over data processing and enables offline operation for enhanced data security. It supports customization, which is important for higher education specific terms and the model is decent when it comes to query interpretation and dataset matching [1]. Having all these points combined, clearly shows that this model was expected to be great for this type of project. This process followed a sequential flow:

- Data ingestion and preprocessing
- Query analysis using locally hosted GPT-J
- Dataset matching based on extracted parameters
- Visualization generation using predefined templates

Unfortunately, after exploring and using the GPT-J model further, I have concluded that this LLM was not suitable for my project, so I had to switch to a more advanced model like Gemini 1.5 Pro and I explained the reasons in section 8.2.

3.2 Risk Identification and Mitigation Strategies Implemented

I have identified some risks during the project planning phase and developed these strategies to resolve them.

- **Data Integration Complexity:** Having difficulties in merging and cleaning multiple CSV files because of HESA inconsistent formats, missing data, or anomalies in the files. This was mitigated by crafting a robust data cleaning methods using Pandas scripts. The process involves standardizing CSV formats within processing pipeline and performing iterative testing on sample datasets to refine the cleaning logic.
- **Inaccurate Query Interpretation:** To resolve the risk of LLM misinterpreting free text queries, I had to switch to Gemini 1.5 Pro because this model is more capable. I used carefully generated by Geminin and tailored by me prompts that were tailored based on the

HESA data. I created a form of feedback that user could leave about how LLM does its job and implemented caching mechanism to store and quickly reuse successful queries.

- **Performance Bottlenecks with Large Data:** The risk of processing large datasets, which could result in slow query response was mitigated by processing a smaller amount of data by referring to indexing files and using a preview mode that limits displayed data. Also, caching for frequent queries to reduce processing time for repeated requests.
- **Automated Data Update Failures:** The initial risk related to automated data scraping failure was addressed by changing the approach to a manual data upload feature, which is less prone to external website structure changes and provides a reliable method for updating data.
- **User Interface/UX Challenges:** The risk of the dashboard not being intuitive was mitigated by focusing on a clean, minimalist design using Tailwind CSS, implementing interactive features based on usability principles, and planning for user feedback collection to identify areas for improvement.
- **Time Constraints and Scope Management:** The risk of incomplete implementation due to the deadline was managed by prioritizing essential requirements, setting realistic milestones in the revised Gantt charts. Also, important to continuously monitor the progress to adjust the scope if necessary.
- **Data Security and Confidentiality:** The risk of insecure handling of HESA data was mitigated by processing only publicly available open-source HESA data and storing it locally.

Initially, I was working on a backup system that was using regular expressions and math to produce clear query interpretation. This was a part of mitigation strategy in case there would be issue with API or LLM in general. After some time of testing this approach, I realized that I could not get the results I was expecting it to provide me. The main reason it kept failing because of HESA complex data and its diversity so I could not use a pattern to solve it. Even though I tried to cover as many edge cases as possible but even slight variation in the provided data formats or query phrasing could cause the algorithm to fail. Moreover, the eventual success of testing the Gemini API approach made this backup functionality unnecessary.

3.3 Project Phases and Timeline

The project was structured into distinct phases to provide a clear roadmap for development and management. This approach aimed to break down the complex task of building an LLM powered HESA data visualization tool into manageable stages with specific objectives and deliverables.

The initial project timeline was planned to move through the main development areas one by one. The proposed phases were designed to move from basic setup and data preparation through the code system development, AI integration and testing. (see Figure 2)

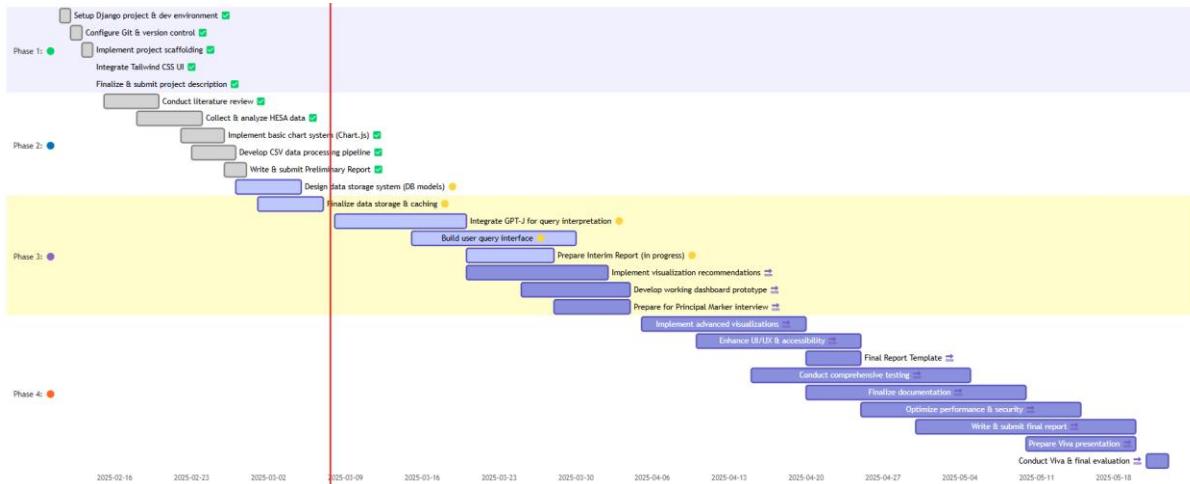


Figure 2: Initial Project Timeline - Gantt Chart

Explanation of Project Phases

- **Phase 1 (Setup & Planning):** Establishing the project foundation, including environment setup, initial research, and defining the project scope and description.
- **Phase 2 (Data & Basic System):** Focusing on acquiring, processing, and storing HESA data, and implementing basic system components like initial visualizations and data processing pipelines.
- **Phase 3 (Core Development & Integration):** Implementing the central features, including LLM integration for query interpretation, building the user interface, and developing the main dashboard functionality.
- **Phase 4 (Finalization & Evaluation):** Refining the system through comprehensive testing, optimizing performance, completing documentation, and preparing for final evaluation.

Phasing breaks the project into manageable stages for better planning, resource allocation, and tracking. It helps reduce risk by allowing for reviews and adjustments at key milestones, ensuring a structured approach to development and completion.

4. System Specification

4.1 Data Flow Diagram

The project is focused on HESA data that is provided as CSV files containing structured data across available sectors. The system creates a comprehensive pipeline that transforms raw statistical data into analyses data that is ready to be used. This process includes following: data

collection, metadata extraction, data cleaning and transformation phases. These phases create optimized datasets that powers NLQ system.

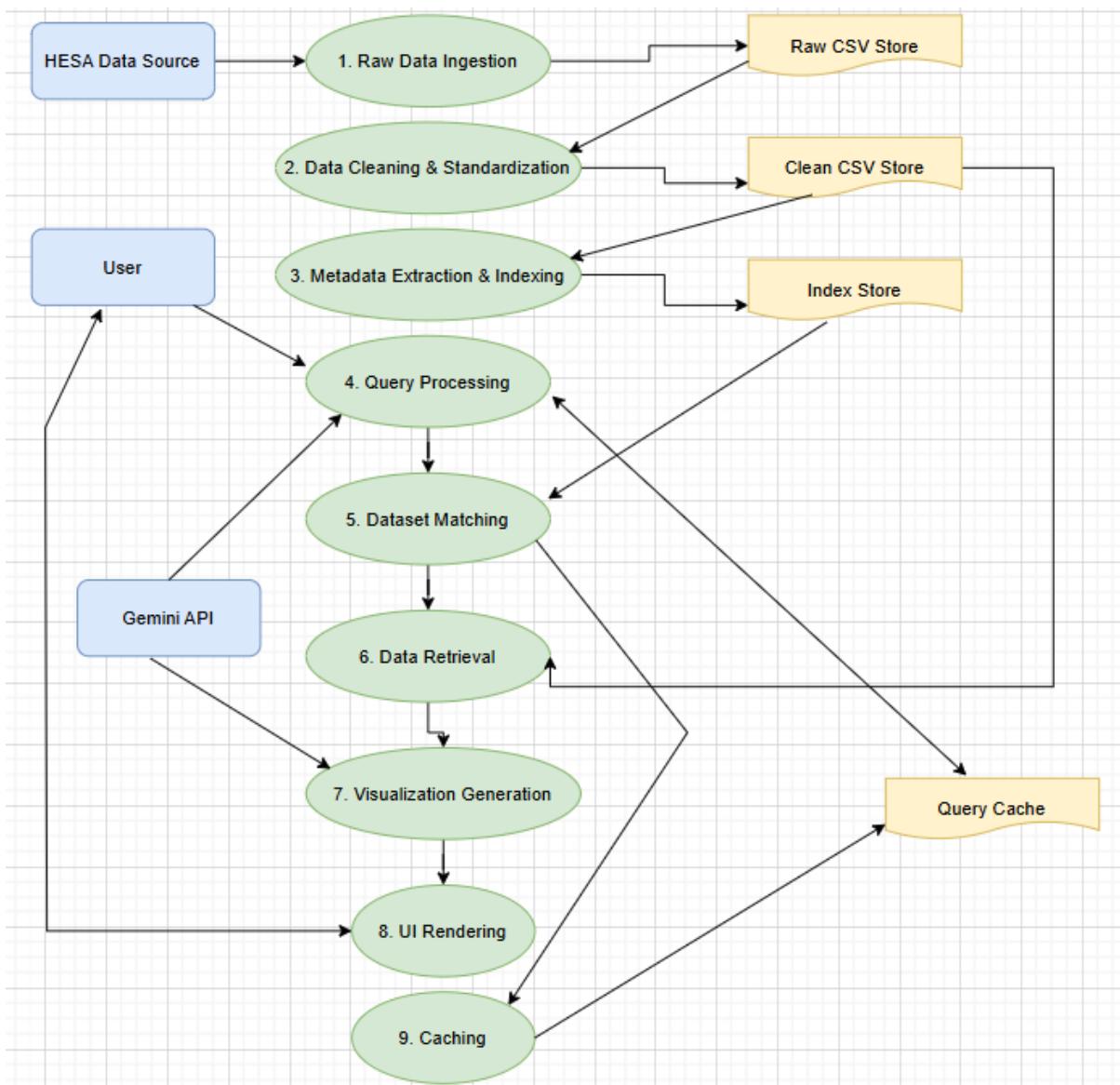


Figure 3: Diagram shows the overall data flow pipeline, from raw HESA data ingestion through cleaning, indexing, query processing, matching, and visualization.

Green ovals = Processes (the steps in the system)

Yellow rectangles = Data stores (where information is kept)

Blue rectangles = External entities (User, HESA, Gemini API)

Process of generating visualization on submitted query:

Source	Destination	Description
User	4. Query Processing	Shows user submitting a query
4. Query Processing	4. Query Processing	Shows processed query parameters moving to dataset matching
5. Dataset Matching	6. Data Retrieval	Shows matched dataset identifiers moving to data retrieval
6. Data Retrieval	7. Visualization Generation	Shows retrieved data moving to visualization generation
7. Visualization Generation	8. UI Rendering	Shows visualization configuration moving to UI rendering
8. UI Rendering	User	Shows results being presented back to the user

Table 2: Explanation of processes to generate visualizations on submitted query

Process of cleaning CSV files:

Source	Destination	Description
Raw Data Ingestion	Raw CSV Store	Shows where processed raw data is stored
Raw CSV Store	2. Data Cleaning & Standardization	Shows raw data being retrieved for cleaning
2. Data Cleaning & Standardization	Clean CSV Store	Shows where cleaned data is stored
Clean CSV Store	3. Metadata Extraction & Indexing	Shows cleaned data being used for indexing
3. Metadata Extraction & Indexing	Index Store	Shows where metadata indexes are stored
Index Store	5. Dataset Matching	Shows indexes being used to match datasets
Clean CSV Store	6. Data Retrieval	Shows where matched datasets are retrieved from
9. Caching	Query Cache	Shows where cached results are stored
Query Cache	4. Query Processing (bidirectional)	Shows checking cache before processing and storing results after

Table 3: Explanation of processes to clean CSV files

Other Processes:

Source	Destination	Description
External Entity: Gemini API	4. Query Processing	This shows the system's interaction with the external AI service

External Entity: Gemini API	7. Visualization Generation	This shows the system's interaction with the external AI service
User	8. UI Rendering (bidirectional)	Shows user providing feedback and modifying visualization parameters

Table 4: Other system processes

4.2 Technical Architecture Design

Technical Architecture Overview

The system employs a modern web application architecture designed to handle data processing, LLM integration, and dynamic visualization.

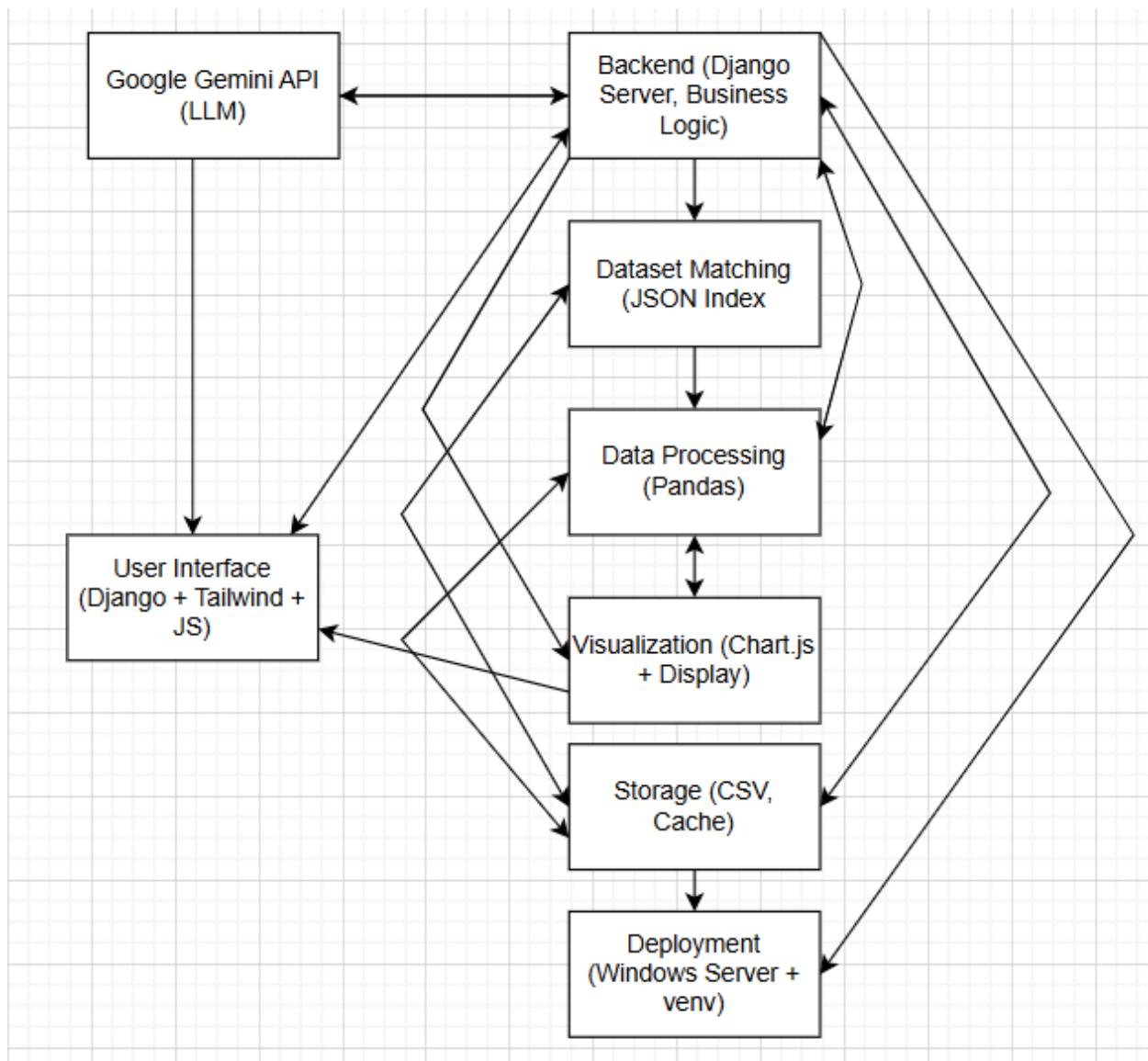


Figure 7: System Technical Architecture Diagram

Source	Interaction	Destination	Description
User Interface	↔	Backend	UI sends natural language queries to backend Backend returns processed results and visualization data to UI
Backend	↔	Google Gemini API	Backend sends user queries to Gemini for entity extraction Gemini returns structured data with institutions, years, and data types
Backend	↔	Storage	Backend reads raw and processed data files Backend writes cache data and user feedback
Dataset Matching	↔	Storage	Reads JSON index file to find relevant datasets Updates index metadata when new files are processed
Data Processing	↔	Storage	Reads CSV data from storage Writes processed results to cache
Data Processing	↔	Visualization	Data processing suggests visualization types based on data characteristics Visualization requests specific data transformations for chart rendering

Backend	↔	Data Processing	Backend sends data processing instructions and query parameters Data Processing returns processed results and analysis summaries
Backend	→	Dataset Matching	Sends extracted entities to identify relevant datasets Forwards query parameters for dataset selection
Backend	→	Visualization	Sends configuration for chart generation Provides display parameters based on query context
Dataset Matching	→	Data Processing	Provides references to matched datasets Sends metadata about dataset structure
Visualization	→	User Interface	Renders charts and data visualizations Provides interactive elements for data exploration
Backend	→	Deployment	Application configuration and settings
Storage	→	Deployment	File system configuration Export to Sheets

Table 5: System Technical Architecture

Development Stack and Dependencies

The project employs a modern web application architecture built on:

Category	Technology	Version	Description
Backend	Django	4.2.5	Web framework

	Python	3.10.12	Primary language
	Pandas	2.0.3	Data manipulation
	NumPy	1.24.3	Numerical processing
Frontend	Tailwind CSS	3.3.3	Styling framework
	JavaScript	ES2021	Client-side functionality
	Chart.js	4.3.0	Visualization library
	HTML5/CSS3	-	Markup and styling
AI Integration	Google Gemini 1.5 Pro API	-	Query interpretation and analysis
	Python google.generativeai	0.3.1	Library
Data Storage	CSV files	-	Primary data storage
	JSON metadata	-	Indexes and cache

Table 6: Project Development Stack and Dependencies

Component Roles and Interactions

The architecture follows a clear separation of concerns:

- **Data Layer:** CSV files store raw and cleaned HESA datasets with JSON metadata for efficient discovery and indexing. Pandas mostly used at the backend to manipulate data (read, clean, transform, and filter)
- **Application Layer:** Django handles requests processing, routing and view rendering. The core functionality is separated into different modules in the Django environment. It includes query processing, dataset matching and configuration of visualization generation.
- **Presentation Layer:** Tailwind CSS provides responsive design through classes. Chart.js renders interactive visualizations based on data processed. The processing happens in the backend and configuration provided via the LLM.

The system follows a request-response cycle where:

1. User queries are submitted from the frontend to Django views.
2. Django views use the data processing pipeline (using Pandas) to access and filter relevant CSV data based on parameters received from the query.
3. The system interacts with the Gemini API to interpret queries, match datasets semantically, and to generate Chart.js visualization configurations.
4. Processed data and visualization configurations are sent back to the frontend.
5. The frontend uses Chart.js to visualize the results in the user's browser.

While architecture focuses on simplicity and maintainability it manages to provide sophisticated data analysis and visualization generations. Pandas is the main dependency that being used throughout the project. It is used for CSV file processing, data cleaning, filtering, transformation and preparation for visualization as shown in Appendices A.1, A.2, A.3.

5. Implementation

5.1 Data Processing Pipeline Implementation

Raw Data Acquisition and Storage

As I explained above in section 1.6, the system is designed to process smaller parts of HESA CSV files to address limitations LLM processing, its performance, memory usage and the cost. This is implemented as a foundational step in the data processing pipeline.

The data ingestion process handles the initial ingestion with validation and basic structure check before being stored. Every raw csv file goes through a check for minimum viability requirements. The file is expected to have a header, columns and basic formatting. (see Appendix D.1)

The system support processing of multiple files at the same time. The current dataset provides up to date HESA data as of academic year 2023/24. When HESA provides new data (There is currently as of 2nd of May no exact date of when new data will be released. HESA says “This page lists planned releases for the next six months. We will confirm the exact date of release at least 4 weeks before publication.” [24]) user can download it by clicking on the download link on the HESA website of a particular dataset that users want as shown in the Figure 5. When data is downloaded that would look like in Figure 4 the user can upload it to the system. It can be done through the UI (User Interface) at /upload-data/ endpoint and user can get there by clicking on “Upload New Data” in the header section. From there user can select downloaded file/s and click “Upload”. (see Appendix D.2)

This manual update system was implemented because HESA does not provide direct API (Application Programme Interface) for automated data retrieval.

```
1 Title:,"Table 1 - HE staff by HE provider and activity standard occupational classification"
2 Subtitle:,"Academic years 2014/15 to 2023/24"
3 Reference ID:,DT025 Table 1
4 Data source:,HESA
5 Data source link:,https://www.hesa.ac.uk/data-and-analysis/staff/table-1
6 Data file canonical link:,https://www.hesa.ac.uk/data-and-analysis/staff/table-1.zip
7 Licence:,,Creative Commons Attribution 4.0 International Licence
8 Date:,2025
9 Filters:,,Country of HE provider,All
10 ,Region of HE provider,All
11 ,Mode of employment,All
12 ,Academic year,2015/16
13 Search:,
14
15
16 "UKPRN","HE Provider","Managers, directors and senior officials","Professional occupations","Associate professional occupations","Clerical and manual
17 "10007783","The University of Aberdeen","0","1,555","18","0","1,565","55","335","425","585","780","75","25","20","28","255","1,845","3,410"
18 "10007849","Abertay University","0","285","5","0","210","15","55","65","110","10","5","0","65","330","540"
19 "10007856","Aberystwyth University","0","840","5","0","850","65","305","255","305","85","75","20","15","175","1,305","2,150"
20 "1000291","Anglia Ruskin University","5","785","15","0","805","68","225","230","425","25","20","15","5","78","1,080","1,885"
21 "10007759","Aston University","0","675","0","0","675","85","135","240","200","25","25","15","10","120","860","1,535"
22 "10007857","Bangor University","5","965","28","0","990","68","185","270","400","45","60","0","5","270","1,220","2,210"
23 "10005571","Bath Spa University","5","605","0","0","610","25","105","105","155","15","25","0","0","75","505","1,115"
24 "10007850","The University of Bath","0","1,345","0","0","1,345","155","375","345","480","75","65","25","20","370","1,920","3,265"
25 "10007152","University of Bedfordshire","0","618","28","0","635","15","115","180","260","15","35","15","0","45","680","1,315"
26 "10005343","Queen's University Belfast","0","1,710","5","0","1,715","30","540","395","735","45","75","15","0","185","2,025","3,740"
27 "10007760","Birkbeck College","5","1,265","0","0","1,270","50","110","280","10","5","15","0","30","610","1,880"
28 "10007140","Birmingham City University","5","1,675","15","0","1,700","90","245","180","495","15","20","15","5","25","1,095","2,790"
29 "10006840","The University of Birmingham","5","3,595","30","0","3,635","280","525","1,030","1,105","175","145","35","15","485","3,795","7,430"
```

Figure 4: Raw data file downloaded from HESA . Example file name dt025-table-1.csv.

Arts Educational Schools	0	0	0	0	0	0	0	0	0
The Ashridge (Bonar Law Memorial) Trust	0	0	0	85	0	0	0	0	0
Hult International Business School Ltd	0	0	0	0	0	0	0	0	0
Aston University	660	2,790	185	840	0	0	160	310	2,390
British Academy of Jewellery Limited	0	0	0	0	0	0	0	0	0
Bangor University	170	1,645	1,850	1,225	0	290	0	0	165
Bath Spa University	0	195	115	615	0	0	0	0	0
The University of	90	1,240	1,815	1,235	0	0	1,320	1,130	3,215

[↑ Reset filters](#) | [↑ Reset sort](#) | [Download table \(csv\)](#) | [Source data \(see note\)](#) | [About DT051 Table 49](#) | [Notes](#)

Figure 5: Demonstration of how to download dataset from HESA

Metadata Extraction and Indexing

Metadata extraction and the specific extraction logic handled by functions as shown in Appendices D.3, D.4. Those functions analyse files without modifying the underlying data. Functions attempt to extract metadata in multiple ways:

1. Using the processor's method
2. Directly reading the file's first line for a '#METADATA:' prefix
3. Creating minimal metadata from the file content when necessary.

This important step creates a lookup index containing:

- Dataset title and description
- Available academic years (formatted as YY/YY)
- Institutions covered
- Column names and data types
- Value ranges for numeric fields
- Categorical field distributions

The system creates the indexing file `hesa_files_index.json` in the root after processing raw files and collect metadata from all available datasets. (see Appendix D.4). This indexing file helps Gemini to query available datasets, and it consists of:

- File locations (both raw and clean versions)
- Update timestamps for freshness tracking (captured in the 'updated_at' field)
- Metadata about each file
- File references.

This approach allows the system to quickly identify candidate datasets without loading all files for each query, this significantly improves the response time. The metadata extraction processes compile this information into structured JSON file. It supports both exact and semantic matching during query interpretation process. (see Appendix E.2).

This indexing is primarily used to minimize the amount of data being sent to Gemini by sending only data that matches requested academic year and is relevant. This is important because it leads to quicker results and fewer API tokens being used, which potentially reduces cost.

Data Cleaning and Standardization

The CSV cleaning process transforms raw files into datasets that are ready for analysis. This is done through the function that handles the full process from reading the raw files to creating the cleaned output files. As mentioned in the introduction chapter, HESA already provides cleaned and high-quality data, but cleaning process is still required because raw files may contain inconsistencies, different formats and metadata that is not suitable for direct analyses. (see Appendix D.5)

Data sanitization follows a structured approach that includes:

- Header Standardization (extracting and normalizing column names)
- Missing Value Handling
- Academic Year Formatting (using regex patterns)

```
r'(\d{4})[.|-&_](\d{2})'
```

- Institution Normalization (ensuring consistent naming of institutions)
- Data Type Conversion (using pandas in the extraction process)
- Character Encoding Normalization (handled by opening files with parameters)

```
encoding='utf-8', errors='ignore'
```

- Comprehensive metadata extraction

The file naming convention for cleaned files is described in Appendices D.5, D.6. The logic includes:

- Title from the raw file (extracted in the code using regex patterns)

```
re.search(r'Table\d+\s*\s*(.*), title_text)
```

- Academic year(s) covered (extracted in the code using regex patterns)

```
re.search(r'(?::|:)\s*Academic year\s*,\s*(20\d{2})\d{2}), line)
```

- Original raw file name

For example, a raw file as shown in Figure 4.

- file name: dt025-table-1.csv
- title: HE staff by HE provider and activity standard occupational classification
- academic year: 2015/16

After cleaning, the file will be named: HE staff by HE provider and activity standard occupational classification 2015-16 dt025-table-1 (1).csv (see Figure 6)

This file name convention is necessary because HESA sometimes provide different datasets but with same academic year and even title. It is vital to make the original name when creating a unique identifier for the cleaned file. The fact that CSV files can have same titles along with academic years represent a challenge since it requires the cleaning process to be more detailed to prevent possible duplications.

1	#METADATA:{ "title": "HE staff by HE provider and activity standard occupational classification", "academic_year": "2015/16", "keywords_title": ["staff", "provider", "activity", "standard", "occupational", "classification"], "keywords_columns": ["provider", "managers", "directors", "senior", "officials", "professional", "occupations", "associate", "clerical", "manual", "total", "academic", "staff", "administrative", "secretarial", "skilled", "trades", "caring", "leisure", "other", "service", "sales", "customer", "process", "plant", "machine", "operatives", "elementary", "non-academic"], "original_filename": "dt025-table-1 (1).csv", "new_filename": "HE staff by HE provider and activity standard occupational classification 2015-16 dt025-table-1 (1).csv"}
2	HE Provider,"Managers, directors and senior officials",Professional occupations,Associate professional occupations,Clerical and manual occupations,TOTAL academic staff,"Managers, directors and senior officials",Professional occupations,Associate professional occupations,Administrative and secretarial occupations,Skilled trades occupations,"Caring, leisure and other service occupations",Sales and customer service occupations,Process, plant and machine operatives",Elementary occupations,TOTAL non-academic staff,Total
3	The University of Aberdeen,0,"1,555",10,0,"1,565",55,335,425,585,78,75,25,20,255,"1,845","3,410"
4	Abertay University,0,205,5,0,210,15,55,65,110,10,5,0,0,65,330,540
5	Aberystwyth University,0,840,5,0,850,65,305,255,305,85,75,20,15,175,"1,305","2,150"
6	Anglia Ruskin University,0,785,15,0,885,60,225,230,425,25,20,15,5,70,"1,080","1,885"
7	Aston University,0,675,0,0,675,85,135,240,200,25,25,15,10,120,860,"1,535"
8	Bangor University,5,965,20,0,998,60,105,270,400,45,60,0,5,270,"1,220","2,210"
9	Bath Spa University,5,605,0,0,610,25,105,105,155,15,25,0,0,75,505,"1,115"
10	The University of Bath,0,"1,345",0,0,"1,345",155,375,345,480,75,65,25,20,370,"1,920","3,265"
11	University of Bedfordshire,0,615,20,0,635,15,115,180,260,15,35,15,0,45,680,"1,315"
12	Queen's University Belfast,0,"1,710",5,0,"1,715",30,540,395,735,45,75,15,0,185,"2,025","3,740"
13	Birkbeck College,5,"1,265",0,0,"1,270",50,110,110,280,10,5,15,0,30,610,"1,880"
14	Birmingham City University,5,"1,675",15,0,"1,700",90,245,180,495,15,20,15,5,25,"1,895","2,790"
15	The University of Birmingham,5,"3,595",30,0,"3,635",280,525,"1,030","1,105",175,145,35,15,485,"3,795","7,430"
16	University College Birmingham,0,290,0,0,290,20,50,35,80,15,10,15,0,80,300,590
17	Birmingham Newman University,0,155,5,0,160,15,20,25,65,5,5,0,0,35,175,335
18	Bishop Grosseteste University,0,95,0,0,95,15,30,40,60,5,30,0,0,15,190,285
19	The University of Bolton,0,300,25,0,330,15,60,70,80,5,15,0,5,50,305,635
20	The Arts University Bournemouth,0,350,0,0,355,15,30,55,40,0,5,0,0,0,145,500
21	Bournemouth University,0,935,0,0,935,50,250,300,255,10,35,0,0,10,905,"1,845"

Figure 6: Example of a cleaned csv file name based on the raw csv file Figure 4.

Transformation and Enhanced Metadata Generation

Cleaned data files are changed into optimized formats for analysis using a data processing system. We also get extra information (metadata) to help with this. (see Appendix D.3, D.4)

This process prepares the data by:

- Calculating metrics for analysis
- Structuring it for clear visualization
- Adding helpful context

Enhanced metadata is generated for cleaned files through the metadata extraction process. It extracts information like titles, academic years and columns while handling potential inconsistencies and missing data. The system tries multiple approaches to extract metadata and if it does not work, it falls back to more basic method (the basic method serves as fallback mechanism but so far, I have not faced this issue). (see Appendix D.4)

The process adds analytical context beyond basic structure:

- It identifies academic years using multiple pattern-matching techniques
- Extracts column information that can be used for visualization type recommendations
- Creates relationships between datasets by incorporating all this information into a central index file that can be queried for cross-analysis.

Rationale for Separate Raw and Clean Files

This dual storage approach provides several advantages: (see Appendices D.5, D.6, D.4)

- Data integrity
 - Original HESA files remain untouched in the raw_files directory)
- Processing efficiency
 - Clean files in the clean_files directory optimize structure for analysis
- Versioning support
 - Changes to cleaning methodologies can generate new clean versions through reprocessing
- Failure recovery
 - Processing errors can be addressed by regenerating clean files from the source
- Selective enhancement
 - metadata can be progressively enhanced for clean files through the multiple extraction methods implemented

It keeps track of the relationship between raw files, processed output and metadata (inside the *hesa_files_index.json*). This ensures that the data is changed using the same standard rules for all HESA information. This full process helps us to protect the original data and prepare it well for analysis and charts generation.

5.2 Query Processing and Interpretation

Overview of Natural Language Query Interpretation

The system processes natural language queries (NLQ) using a sophisticated pipeline. It leverages Gemini 1.5 Pro to transform free-text queries into structured parameters. It sends carefully created prompts to Gemini API along with system instructions that explains to the AI how to extract structured data from NLQ. (see Appendices L, E.1)

This AI powered approach recognizes query formulations like “Show me student enrolment at Leicester in 2019” or “What was the staff-student ratio for Russell Group universities between 2015-2018?”. The method identifies query intent by constructing a detailed system prompt which instructs AI to extract institutions, academic years and user requested data. It preserves semantic relationships through the structured JSON response format that captures relationship between extracted entities.

The functionality understands different types of queries through its complex prompt engineering. As mentioned above, this includes:

- Instructions for handling academic year conventions
- Institution name variations
- Data request specifications.

It adapts the response format according to the query type and returns structured data that can be used by other application functionalities. The implementation shown in Appendix E.1 along with

integration with the Gemini API enables intuitive data exploration without the need for users to learn specific query syntax

Parameter Extraction, Tokenisation, and Normalization Techniques

The parameter extraction process as shown in Appendix L involves:

- Tokenization of queries into meaningful segments
- Entity recognition for key parameters (institutions, years, metrics)
- Normalization of recognized entities to standard formats
- Validation against expected patterns and ranges, and structuring for dataset matching.

Parameter extraction does not directly tokenize the query itself but rather constructs a detailed prompt that instructs the AI to perform entity extraction. Using this method, the system gets detailed instructions (the prompt) on how to find and sort the key information, such as institutions, time periods and comparison groups. The method includes a step to process the AI's answer. This step takes the information and transform it into the exact formats the system needs. (see Appendix E.1).

One of the most important aspects of this process is to accurately interpret the academic year which spans two calendar years. The system uses a custom function detailed in Appendix E.2 that applies higher education-specific logic to correctly interpret various date references:

- “Starting in 2017” or “from 2017” → Matches the 2017/18 academic year
- “End of 2017” or “ending in 2017” → Matches the 2016/17 academic year
- Plain years (e.g., “in 2017”) → Treated as the starting year (2017/18)
- Year ranges (e.g., “2016 to 2017”) → Matches the 2016/17-2017/18 academic years
- “Past 5 years” → Matches the last 5 academic years from the current year

This logic ensures that the data-based queries are correctly interpreted and ready for further use.

LLM Integration Design

The Gemini 1.5 API integration is shown in Appendices E.1, K, L.

It follows a structured design approach:

- Prompt Engineering
 - Designing effective prompts for Gemini that include the user query, relevant context (like the structure of available data), and clear instructions for entity extraction, intent classification, and relevance scoring.
- API Client Implementation

A robust API client was developed to communicate with the Gemini API. It includes initialization with API key, tailored request formatting, and extensive error handling with detailed logging. It uses the requests library for making HTTP

requests.

- Response Parsing and Validation
 - Designing logic to parse the structured JSON response from Gemini and validate its content to ensure it can be used by the subsequent processing steps.
- Caching Integration
 - Designing the query processing pipeline to check the cache before sending a request to the Gemini API, and to store the results upon receiving a valid response.

The Gemini 1.5 Pro API was configured with specific parameters to ensure optimal performance and output for the project's needs as shown in Appendix L. Key settings include:

- Output Capacity
 - Configured to generate up to 2,048 tokens per response based on the configured parameter *maxOutputTokens*: 2048. Balances detailed output with managing token usage.
- Temperature Setting
 - Set to a low value of *temperature*: 0.1. This makes responses more deterministic, focused, and less random. This is important for accurate data interpretation and structured output generation.
- Top-P and Top-K
 - Using *topP* = 0.95 and *topK* = 40. These parameters are responsible for the randomness of the output. This ensures a balance between creative interpretation and following to the provided context. This helps with reliability of entity extraction and data matching.

The API integration flow detailed in Appendix L includes:

- Preparation of well-structured prompts with necessary context (system_prompt and user_prompt)
- Secure API communication with appropriate authentication

```
headers = {"Content-Type": "application/json", "x-goog-api-key": self.api_key}
```
- Response validation and error handling (with multiple error checks and try-except blocks)
- Parsing of structured information from responses (JSON extraction and validation)

- Seamless integration with the existing processing pipeline.

5.3 Visualization Generation

Chart Generation

It analyses the selected dataset and user request, then calls the Gemini using the client to determine the correct chart type. After the type is configured, it generates Chart.js configurations for rendering the determined chart. (see Appendix F.1)

The implementation includes

- Detailed processing of dataset information
- Construction of prompts for the Gemini API
- Extraction of chart configuration from the AI response.

Chart Types Implementation

The system prepares data for different chart types like line, bar, and pie charts. This functionality is implemented across several methods (see Appendices F.2, F.3, F.4) including:

- Handles different chart type requests
- Gets AI recommendations for chart types
- Allows dynamically changing the visualization

Visualization generation uses hybrid approach. Gemini generates Chart.js configuration by defining chart type, its structure, labels and options. The generation is based on the dataset characteristics and user submitted request. The whole configuration is then sent to the front-end where the actual dataset data which is fetched and [processed by the backed is combined with the configuration. Lastly, after it gets combined the Chart.js renders the interactive chart for the user on the webpage. The importance of this approach is that it ensures data privacy because the actual data does not go to Gemini thus, technically this approach can be used with other data that is more sensitive. (see Appendix F.1)

Frontend Rendering

The client-side rendering instantiates Chart.js visualizations on canvas elements. Interactive controls are managed for creating new charts and switching visualization formats. Event listeners on UI elements connect user actions to these functions which enables dynamic updates based on user queries and selections. (see Appendices G.1, G.2, G.3)

Interactive Visualization Controls

The system allows users to toggle data visibility, swap axis and adjust colour directly on the charts. The control is implemented in the frontend JavaScript and connected to the backend functions that handles data transformation request. (see Appendices G, H.1)

Real time Chart Regeneration

Generated charts are updated automatically when filtering parameters or visualization options changes. This created using the API endpoint and AJAX request in the frontend. (see Appendix H.1)

Visualization Type Recommendations

The system uses AI to suggest the best chart type for visualizing data. It analyses the dataset and explains why specific chart type is recommended therefore helping user to understand effective data visualization. (see Appendices F.3, H.2)

5.4 UI Implementation and Interactivity

Dashboard Interface

The interface is implemented with key functions including: (Appendices I.1, I.2, I.3, I.4)

- Handles displaying query results
 - Transforms JSON responses into user-friendly HTML, including formatting institution names, academic years, and rendering clickable dataset cards.
- Creates reviewable dataset cards
 - Generates interactive cards for each matched dataset, displaying title, metadata, and sample rows with clean formatting.
- Manage loading states
 - Provides visual feedback during API calls, displaying animated spinners and contextual messages about current processing steps.
- Event listeners for query input
 - Captures user inputs, validates queries, and submits requests to backend APIs with appropriate feedback mechanisms.

Dataset Details View

Implemented with functions including:

- Displays interactive data visualizations (see Appendix G.1)
 - Renders dynamic charts with hover effects, tooltips, and interactive legends that allow users to isolate specific data series.
- Creates charts based on user specifications (see Appendix G.2)
 - Processes natural language visualization requests and transforms them into appropriate data visualizations, considering dataset characteristics for optimal presentation
- Allows switching between visualization formats (see Appendix G.3)
 - Enables seamless conversion between chart types (bar, line, pie) while preserving the dataset context and user's visualization intent.
- PDF and CSV export functionality implemented through dedicated download handlers (see Appendix I.5)
 - Provides client-side export capabilities that format data with proper headers, styling, and metadata for professional reports and further analysis.

Dynamic Dataset Previews

The system displays the first three rows of matched datasets with collapsible sections for expanded metadata so users can quickly find relevant dataset.

Data previews are implemented with functions to:

- Prepare data for clean display (see Appendix I.6)
 - Removes technical artifacts like metadata headers and normalizes column names for better readability. It helps to ensure consistent presentation across different dataset types.
- Toggle functionality through (see Appendix I.7)
 - Enables users to expand and collapse detailed information sections with smooth transitions. It keeps the screen clean and allows the user to access that extra details when needed.
- Collapsible metadata sections with animation classes
 - animate-fade-in, animate-fade-out
 - Creates a polished user experience with subtle transitions when revealing or hiding content, improving interface responsiveness and visual feedback.

Responsive Design

The interface adapts to various screen sizes while maintaining data integrity and usability.

The interface uses Tailwind CSS classes throughout the templates with responsive design patterns including:

- Flexible grid layouts

```
<div class="grid grid-cols-1 md:grid-cols-3 gap-4">
  <div class="border rounded p-3 bg-white">...</div>
  <div class="border rounded p-3 bg-white">...</div>
  <div class="border rounded p-3 bg-white">...</div>
</div>
```

- Responsive padding and margins

```
div class="max-w-7xl mx-auto px-4 py-8 md:px-6 lg:px-8">
  <h1 class="text-2xl md:text-3xl font-bold mb-4 md:mb-6">HESA AI Powered Data Dashboard</h1>
</div>
```

- Overflow handling for data tables

```
.table-container {
  overflow-x: auto;
```

```
max-height: 300px;  
overflow-y: auto;  
scrollbar-width: thin;  
}
```

- Mobile-friendly controls with appropriate sizing and spacing

```
<button  
    class="w-full sm:w-auto px-3 py-2 sm:px-4 sm:py-3 text-sm sm:text-base bg-blue-600 text-white rounded-md flex items-center justify-center"  
>  
    <span class="flex items-center">  
        <svg class="h-4 w-4 mr-2" viewBox="0 0 24 24">...</svg>  
        <span>Search</span>  
    </span>  
</button>
```

5.5 Feedback System

Feedback Collection

It provides an interface for user to rate received results with options as “Not Helpful”, “Helpful” and “Very Helpful”. Optionally, displays a comment area for detailed feedback if user selects anything but the option “Very Helpful”. (see Appendix J)

Feedback Submission

The system collects relevant data such as:

- Query details
- Dataset information
- Visualisation settings
- User rating
- Comments

This data is sent to the server for storage and analysis so it can be used further to enable data-driven improvements. (see Appendix J)

5.6 Caching and Performance Optimization

Query Cache Implementation

The system implemented a caching functionality that stores and retrieves query results as JSON documents in */cached* folder. This caching system uses hash-based keys and metadata tracking to avoid redundant API calls and data processing. The cached data includes all the information about the query.

The cached data lasts for 30 days so after it expires the system removes it from the folder. This allows to maintain the system effectively without manual clean up. Also, it tracks usage statistics

like access accounts and last accessed timestamp. These statistics can be used to make performance analyses and optimization of frequently used queries. (see Appendix K)

5.7 Codebase Structure and Endpoints

Codebase Overview

hesa-llm-visualisation/

core/	# Main application directory
templates/	# HTML templates
ai_dashboard.html	# AI powered dashboard interface
base.html	# Base template with common layout
documentation.html	# Project documentation page
upload_data.html	# Data upload interface
visualization/	# Chart visualization templates
static/	# Static files (CSS, JS, images)
css/	# Compiled CSS
js/	# JavaScript files
ai_dashboard.js	# Main dashboard functionality
src/	# Source files for CSS
utils/	# Utility functions
query_cache.py	# Query caching implementation
services/	# Service layer
data_processing/	# Data processing modules
visualization/	# Visualization generation code
management/	# Django management commands
views.py	# Main view controllers (5300+ lines)
urls.py	# URL routing
gemini_client.py	# Google Gemini API integration
data_processor.py	# Data processing logic
data_processing.py	# Additional data processing utilities
models.py	# Data models (minimal)
apps.py	# App configuration
hesa_llm_visualisation/	# Project configuration
settings.py	# Django settings
urls.py	# Project URL routing
wsgi.py	# WSGI configuration
data/	# Data storage directory
raw_files/	# Original HESA data files
cleaned_files/	# Processed data files

```

    └── docs/                                # Documentation files

    ├── manage.py                            # Django project management script
    ├── debug_metadata.py                   # Debugging utilities
    ├── csv_cleaning.py                    # CSV file cleaning utilities
    └── hesa_files_index.json              # Index of available HESA data files

    └── run_django_server.bat            # Windows batch script to run server
    └── data.env                           # Environment variables

    ├── package.json                        # Node.js dependencies
    ├── tailwind.config.js                # Tailwind CSS configuration
    ├── postcss.config.js                 # PostCSS configuration
    └── README.md                           # Project documentation

```

Key Endpoints

Below are key endpoints that user can access and easily explore. There are other endpoints that are not for user interaction but rather for the usage of Gemini to return or display data.

- / or /ai-dashboard/ (*views.ai_dashboard*)
 - Main entry point to the application
- /documentation/ (*views.documentation*)
 - Documentation page
- /process_gemini_query/ (*views.process_gemini_query*)
 - Processes natural language queries using Google Gemini
- /ai_dataset_details/ (*views.ai_dataset_details*)
 - Gets detailed information about datasets
- /save_feedback/ (*views.save_feedback*)
 - Saves user feedback about query results
- /upload-data/ (*views.upload_data_view*)
 - Data upload interface
- /process-upload/ (*views.process_upload*)
 - Processes uploaded data files

5.8 System Setup and Running Instructions

Please refer to the STARTUP.md file inside the submitted project.

Location: vr112\STARTUP.md

It has full instruction on how to set up the project and run it

6. Demonstration

6.1 Using the AI Powered Query Interface

Submitting a Natural Language Query

Main Dashboard (/ or /ai-dashboard)

1. Navigate to the main dashboard (see Figure 8)
 - Open the application in the browser at the root URL
 - You will see the “HESA AI Powered Data Dashboard” with a search input field
2. Enter a NLQ in the search box (see Figure 9)
 - Query example: “How many postgraduates studied at the University of Leicester between 2018-2020?”
 - The query input field displays a placeholder example to guide the user
3. Configure search parameters (optional)
 - Set “Number of matches to show” (by default: 3)
 - Select a Mission Group filter if needed. Options: None. Russell Group, Million+ or University Alliance
4. Click the “Search” button (see Figure 10)
 - The system will process query using Gemini 1.5 Pro API
 - A loading indicator will appear during the process

The screenshot shows the HESA AI-Powered Data Dashboard. At the top, it says "Ask natural language questions about HESA data and get intelligent answers". Below that is a large input field labeled "Ask Anything About HESA Data" with the placeholder text "e.g., How many postgraduates studied in University of London in 2017?". Underneath the input field, there is a note: "Ask any question about HESA data, including institutions, time periods, and metrics". To the left of the input field, there is a "Number of matches to show:" dropdown set to "3" and a "Mission Group Filter" section with a radio button for "None" (which is selected) and other options for "Russell Group", "Million+", and "University Alliance". To the right of the input field is a "Sample Questions" button. At the bottom right is a blue "Search" button with a magnifying glass icon.

Figure 8: Main Dashboard webpage

HESA AI-Powered Data Dashboard

Ask natural language questions about HESA data and get intelligent answers

Ask Anything About HESA Data

Your Question:

How many postgraduates studied at University of Leicester between 2018-2020?

Ask any question about HESA data, including institutions, time periods, and metrics

Number of matches to show: Mission Group Filter: None Sample Questions ▾

Russell Group
 Million+
 University Alliance

Q Search

Figure 9: Main Dashboard webpage with query example

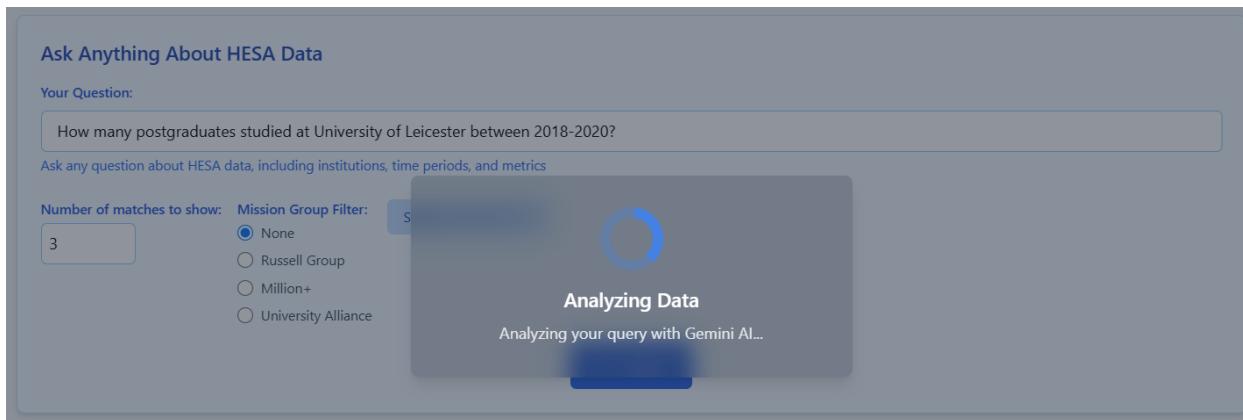


Figure 10: Indicator of submitted query

6.2 Query Interface

Reviewing Query Analysis Results

Endpoint: /process_gemini_query (called in background)

1. View query interpretation (see Figure 11)
 - The system displays how it understood the query
 - Shows extracted institution (“University of Leicester”)
 - Shows extracted time period (“2018-2020”)
 - Shows extracted requested data (“postgraduate students”)
2. Review matching datasets (see Figure 12)
 - Each matching dataset shows score of relevance (percentage match)
 - Datasets are ranked by relevance score to the submitted query
 - Preview of data columns and rows is provided

Query Analysis by Gemini AI | Using Google Gemini API | Hide Details

Your query:
How many postgraduates studied at University of Leicester between 2018-2020?

Institutions:	Years:	Year Range:
• University of Leicester	• 2018/19 • 2019/20 • 2020/21	2018 to 2020
Data requested: postgraduate		

Entity extraction performed by Google's Gemini AI

Figure 11: Query Interpretation information

Matching Datasets (3) | Semantically matched by Gemini AI

The following datasets were found to be relevant to your query:

How Datasets Are Matched

Results are ranked based on how well they answer your specific question:

80-95% : Perfect match for your question	60-80% : Strong match with the data you need
40-60% : Partial match with related data	10-40% : Minimal relevance to your question

1. HE academic staff by HE provider and highest qualification held | Match: 0.95 (94%)

Show Details

Available Files:

2018/19 - HE academic staff by HE provider and highest qualification held 2018-19 dt025-table-8 (4).csv

He Provider	Doctorate	Other Higher Degree	Other Postgraduate Qualification	First Degree	Other Undergraduate Qualification	Other	Unknown	Total Academic Staff	Academic Year
The University of Leicester	1,075	355	25	205	5	35	0	1,705	2018/19

Figure 12: Matched datasets previews and datasets matching information

Working with Dataset Details and Visualization

Endpoint: `/ai_dataset_details` (called in background)

1. Click on a dataset from the found results (see Figure 13)
 - Look for most relevant dataset for submitted query
 - Review dataset details including time period covered and available institutions
 - Click “View Details” or “Explore” button for the chosen dataset
2. Review detailed dataset information
 - View metadata about the dataset
 - See all available columns and preview rows
 - Data is filtered based on query parameters (institution and years)

Academic Years: 2018/19, 2019/20, 2020/21

References:

HE academic staff by HE provider and highest qualification held 2018-19 dt025-table-8 (4).csv
HE academic staff by HE provider and highest qualification held 2019-20 dt025-table-8 (5).csv
HE academic staff by HE provider and highest qualification held 2020-21 dt025-table-8 (6).csv

Matched Terms:

postgraduate doctorate

Dataset Contains:

Data containing postgraduate information

Why this matches:

Column 'other postgraduate qualification' contains 'postgraduate' which matches the query intent Column 'doctorate' contains 'doctorate' which matches the query intent

Figure 13: Information about why Gemini selected this dataset

HE academic staff by HE provider and highest qualification held

2018/19 - HE academic staff by HE provider and highest qualification held 2018-19 dt025-table-8 (4).csv									Show Details
HE Provider	Doctorate	Other higher degree	Other postgraduate qualification	First degree	Other undergraduate qualification	Other	Unknown	Total academic staff	Academic Year
The University of Leicester	1,075	355	25	205	5	35	0	1,705	2018/19
2019/20 - HE academic staff by HE provider and highest qualification held 2019-20 dt025-table-8 (5).csv									
HE Provider	Doctorate	Other higher degree	Other postgraduate qualification	First degree	Other undergraduate qualification	Other	Unknown	Total academic staff	Academic Year
The University of Leicester	1,085	335	40	195	5	25	10	1,695	2019/20

Figure 14: A selected dataset

Generating and Interacting with Visualizations

Endpoint: /visualization_api (called in background)

- Analyse visualisation recommendation (see Figure 15)
 - The system uses Gemini to automatically suggest the most appropriate chart type for selected dataset
 - Example: Line chart for time-series data showing trends over multiple years
 - A brief explanation of why selected type of chart is recommended
- Customize visualization (see Figure 17)
 - Enter a specific visualisation request ("Show a comparison of postgraduate numbers by year")
 - Select different chart types from options (bar, line and pie)
 - View explanations for each chart type to understand its strengths and limitations

3. Generate visualisation (see Figure 17)
 - Click “Generate Visualization” button
 - The system will create an interactive chart based on chosen specifications
 - Chart will appear with proper titles, label, and colour

4. Interact with generated visualization (see Figure 18,19)
 - Hover over data points to see exact values
 - Toggle data series on/off (for multi-series charts)
 - View the visualization in full screen mode if desired
 - See insights of generated charts

Data Visualization

Based on this dataset, Gemini recommends a line for visualization.

A line chart is most appropriate because the dataset shows the trend of academic staff qualifications at The University of Leicester over multiple academic years (2018/19, 2019/20, and 2020/21). Line charts excel at visualizing changes over time, making it easy to see how the different qualification categories have evolved.

Select visualization type:

Bar Chart enforced when selected

Line Chart enforced when selected

Pie Chart enforced when selected

Try one of these examples:

- Show a line chart visualizing the trend of Doctorate holders at The University of Leicester from 2018/19 to 2020/21.
- Plot the changes in 'Other higher degree' and 'First degree' qualifications at The University of Leicester across the academic years 2018/19, 2019/20, and 2020/21.
- Compare the trend of 'Other postgraduate qualification', 'Other undergraduate qualification', and 'Other' qualifications held by academic staff at The University of Leicester between 2018/19 and 2020/21 using a line chart.

Figure 15: Interactive Data Visualization initially generated by Gemini

Select visualization type:

Gemini's recommendation: line

Bar charts are excellent for comparing the number of academic staff with different qualifications across institutions and academic years. They allow for quick visual comparison of values within a specific category (e.g., 'Doctorate') between The University of Leicester across different years, or comparing different qualification levels within the same institution and year. However, bar charts might become cluttered if comparing too many categories at once. For showing trends over time, a line chart might be more suitable.

Figure 16: Re-generated recommendation message after different type of chart is selected

Describe what you'd like to visualize:

Show a comparison of Other postgraduate qualification numbers by year.



Generate Visualization

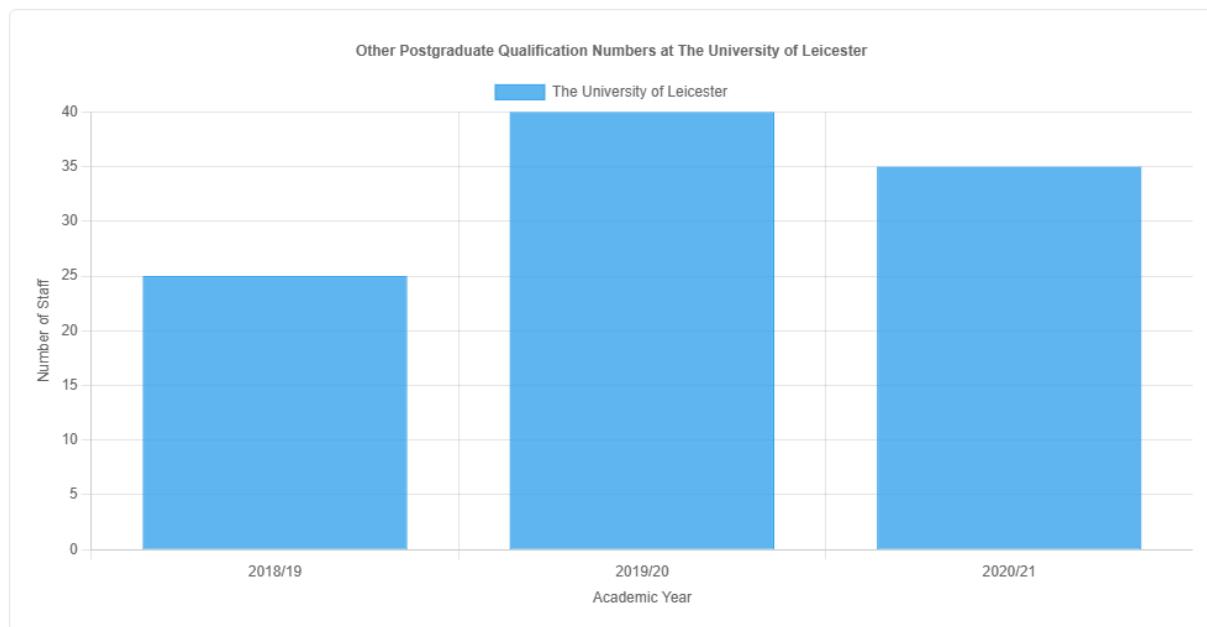


Figure 17: Generated Bar Chart based on specified requirements from selected dataset

💡 Insights

Other Postgraduate Qualification Numbers at The University of Leicester

This bar chart shows the number of academic staff at The University of Leicester holding other postgraduate qualifications for the academic years 2018/19, 2019/20, and 2020/21.

2018/19: 25 staff members held other postgraduate qualifications.

2019/20: This number increased to 40.

2020/21: There was a slight decrease to 35 staff members.

The chart indicates a fluctuation in the number of staff with other postgraduate qualifications over these three years, with a noticeable increase in 2019/20 followed by a small decline in 2020/21.

Figure 18: Generated Insights of generated bar chart from Figure 17.

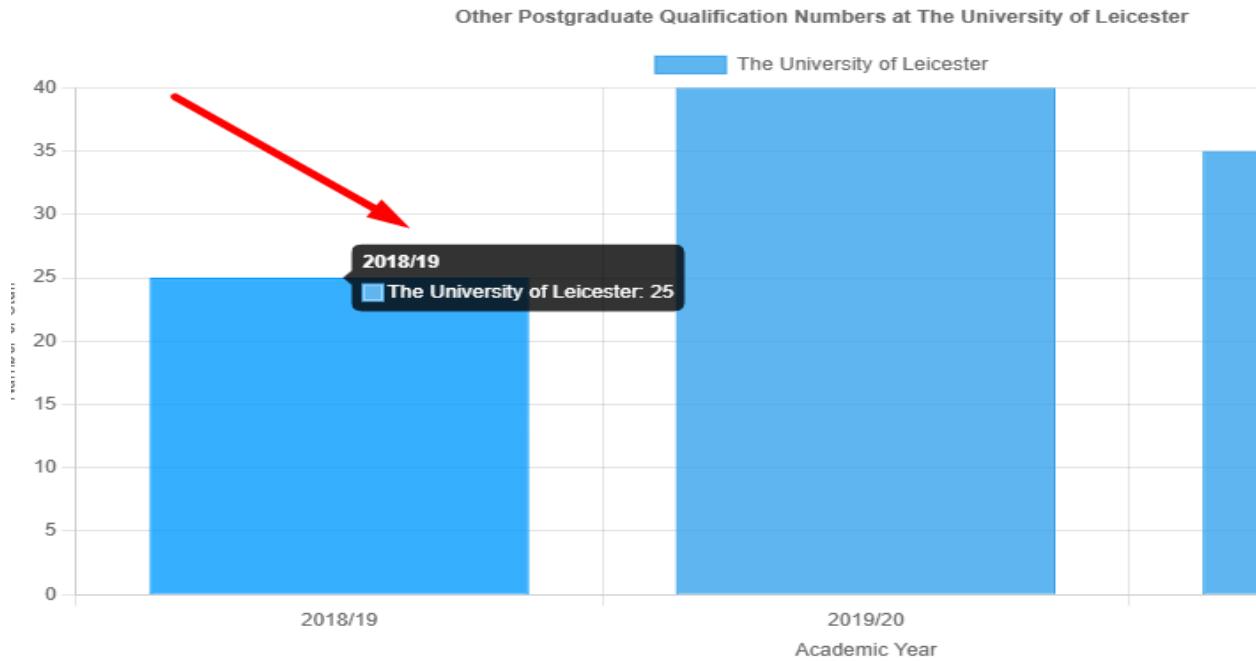


Figure 19: Hovering over data points to see more detailed information

6.3 Providing Feedback on Results

Submitting Feedback on Query and Visualization

Endpoint: /save_feedback (called in background)

1. Locate the feedback section (see Figure 20)
 - Found at the top of webpage after query is submitted
 - Look for buttons “Not Helpful”, “Helpful” and “Very Helpful”
2. Provide rating and comments (see Figure 21)
 - Select “Not Helpful”, “Helpful” or “Very Helpful” to rate the result
 - Enter detailed comments if appropriate rating is selected in the feedback text area (optional)
 - Example comment: “The visualization clearly showed the trend in postgraduate numbers but only one dataset was returned”
3. Submit the feedback (see Figure 22)
 - Click “Submit Feedback” button
 - A confirmation message will appear after that

How helpful were these results?

Very Helpful
Helpful
Not Helpful

Figure 20: Feedback Submission Interface

How helpful were these results?

How could we improve these results?

The visualisation clearly showed the trend in postgraduate numbers, but I'd like to see a breakdown by degree type.



Figure 21: Provided comment for feedback



Figure 22: Submitted feedback

6.4 Uploading Custom Data

Navigating to the Data Upload Section

Endpoint: `/upload-data`

1. Access the upload page (see Figure 22)
 - Click “Upload Date” in the navigation menu
 - The system will redirect to the data upload interface
2. Understand upload requirements (see Figure 22)
 - View instructions for supported file formats (CSV)
 - See example of correctly formatted files
 - Read about file size limits and columns formatting requirements

HESA Dashboard Dashboard Upload New Data Documentation



Upload New Data

Select CSV File(s)

Select CSV files (max 1000 lines per file)

Note: Maximum 50 files can be uploaded at once. Each file must be less than 1000 lines.

Processing Options

Add files to existing data
 Overwrite all existing data (will delete all current files)

Figure 22: Data Upload Interface

Uploading and Processing Data Files

Endpoint: `/process-upload` (called when uploading)

1. Select files to upload (see Figure 23)

- Click “Choose Files” or drag and drop files into the upload area
 - Select one or more CSV files with HESA compatible data
 - Example raw file: “dt025-table-1 (2).csv”
2. Configure upload options
 - Add files to existing data
 - Overwrite all existing data (it will delete all current files)
 3. Initiate upload process (see Figure 24)
 - Click “Upload and Process Files” button to begin processing
 - Progress indicator shows upload and processing status
 4. Review processing results (see Figure 25)
 - See confirmation of successful upload
 - View any validation warnings or errors

Upload New Data

Select CSV File(s)

Select CSV files (max 1000 lines per file)

Selected: dt025-table-1 (2).csv ←

Note: Maximum 50 files can be uploaded at once. Each file must be less than 1000 lines.

Processing Options

Add files to existing data

Overwrite all existing data (will delete all current files)

Upload and Process Files

Cancel

Figure 23: Select files to be uploaded

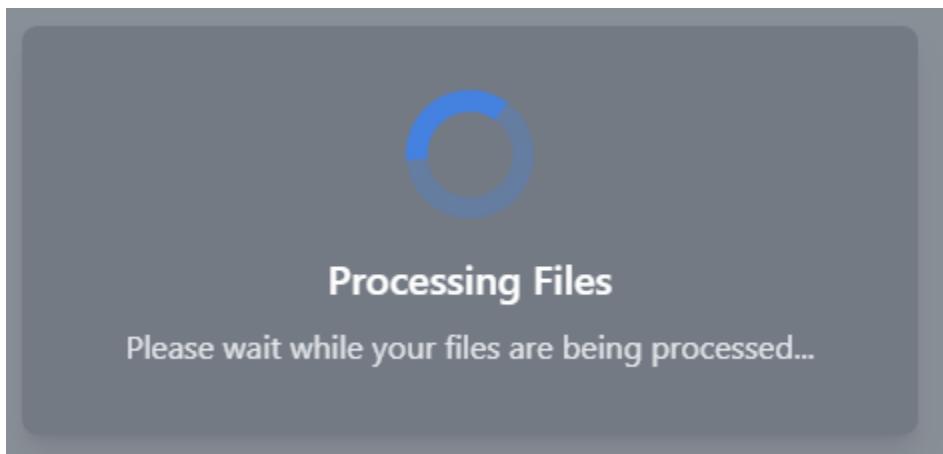


Figure 24: Uploading progress indicators

Upload New Data

Processing Results

Processing completed successfully!

Files in system: 295

Files added: 1

Added files:

- dt025-table-1 (2).csv

Cleaning status: All files cleaned successfully (295/295)

Indexing file created: Yes

[Return to Dashboard](#)

[Upload More Files](#)

Figure 25: Confirmation of successful upload

6.5 Exploring Documentation

Accessing the Documentation Page

Endpoint: /documentation

1. Navigate to documentation (see Figure 26)
 - Click “Documentation” in the navigation menu
 - Documentation page provides guidance on how to use the system
2. Explore available documents sections (see Figure 27 - 33)
 - Introduction to the HESA dashboard
 - Guide to formatting effective queries
 - Explanation of visualization options
 - Data dictionary with columns descriptions
 - Tips and Ticks that can help to simplify the use of the system

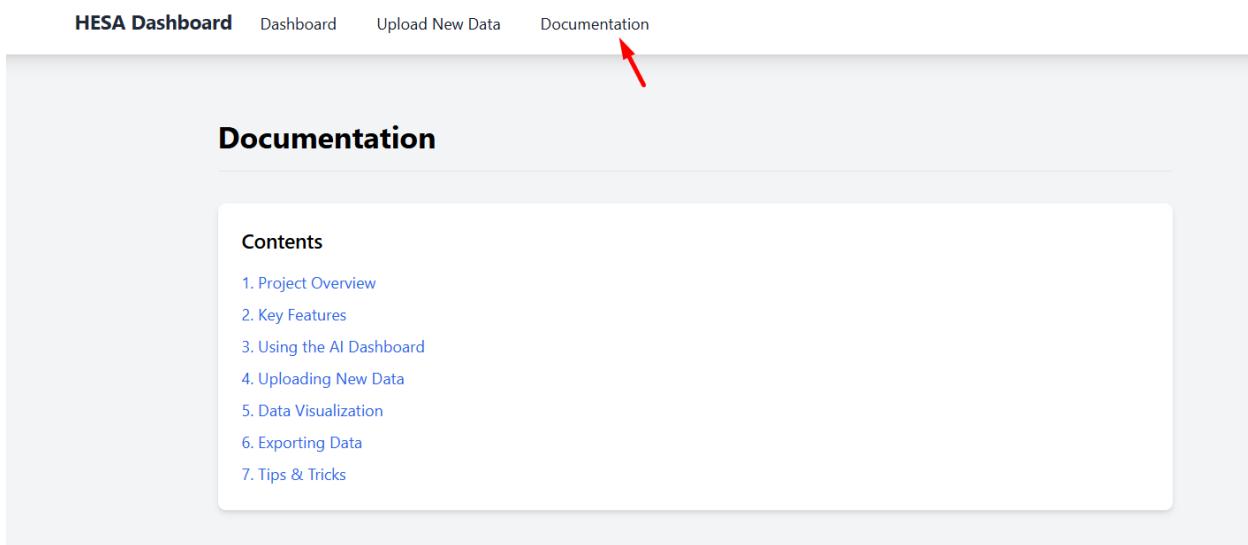


Figure 26: Documentation Page

1. Project Overview

The HESA LLM Visualisation platform is an advanced tool for analyzing and visualizing Higher Education Statistics Agency (HESA) data using AI-powered natural language processing. It allows users to:

- Query education data using plain English
- Upload and process HESA CSV files
- Generate insightful visualizations automatically
- Compare data across institutions and years
- Explore trends in higher education statistics

This platform leverages Google's Gemini API to understand natural language queries, identify entities such as institutions and time periods, and match them to the most relevant datasets.

Key Technologies:

- Google Gemini AI for natural language processing
- Django web framework
- Interactive data visualizations
- CSV processing and cleaning pipeline

Figure 27: Project Overview Information

2. Key Features

Natural Language Queries

Ask questions in plain English about higher education data without needing to write complex database queries.

Intelligent Visualizations

Automatically generate the most appropriate charts and graphs based on your query and data type.

Data Upload & Processing

Upload your own HESA CSV files for cleaning, processing, and integration into the analysis platform.

Institution Filtering

Filter data by specific institutions or mission groups like Russell Group or Million+.

Time-Series Analysis

Analyze trends across multiple academic years with intelligent time period detection.

Data Export

Export your analysis results in multiple formats for presentations or further analysis.

Figure 28: Key Features

3. Using the AI Dashboard

Getting Started

- 1. Navigate to the Dashboard:** Click on "Dashboard" in the navigation bar.
- 2. Enter Your Query:** Type your question in natural language in the search box.

Example queries:

- "Show me student numbers at University of Oxford from 2018 to 2020"
- "Compare undergraduate enrollment for Russell Group universities"
- "What was the financial income for University of Leicester in 2021/22?"

- 3. Review Results:** The system will analyze your query, extract key entities (institutions, years, data types), and present matching datasets.
- 4. Select a Dataset:** Click "Select this dataset" on the most relevant match to proceed to detailed analysis.
- 5. Provide Feedback:** Rate how helpful the results were to help improve future queries.

Mission Group Filtering

You can filter results by university mission groups:

- **Russell Group:** 24 research-intensive universities in the UK
- **Million+:** Association of modern universities
- **University Alliance:** Alliance of professional and technical universities
- **No Filter:** Include all universities in your search

 **Pro Tip:**

The system handles typos and common variations in institution names. For example, if you type "Univesty of Lieceter", it will automatically correct to "University of Leicester".

Figure 29: Introduction to project UI

4. Uploading New Data

You can add your own HESA CSV files to the system for analysis. Here's how:

1. **Navigate to Upload:** Click "Upload New Data" in the navigation bar.
2. **Select Files:** Click the upload area to select CSV files from your computer (maximum 50 files at once, each under 1000 lines).
3. **Choose Processing Option:**
 - **Add files:** Add new files to existing data
 - **Overwrite:** Replace all existing files with new ones
4. **Submit:** Click "Upload and Process Files" and wait for processing to complete.
5. **Review Results:** Once processing is finished, you'll see a summary of added files and any errors.

What happens during processing:

1. Files are validated for CSV format and line count
2. Data is cleaned and standardized
3. Metadata is extracted (academic years, titles, etc.)
4. An index file is created for quick searching
5. Files are made available for querying via the AI Dashboard

⚠ Important:

Your CSV files should follow HESA data format. Files that don't match this format may not be processed correctly.

Figure 30: Uploading new data guide

5. Data Visualization

After selecting a dataset, the system will analyze your data and recommend the most appropriate visualization type.

Available Chart Types

- **Bar Charts:** For comparing categories
- **Line Charts:** For time series data
- **Pie Charts:** For showing composition

Customizing Visualizations

- Change chart type with a single click
- Modify title and axis labels
- Filter data points to focus on specific information
- Toggle legend visibility
- Change color schemes

💡 Pro Tip:

You can request specific visualization types by including them in your query. For example: "Show me a bar chart of student numbers at Oxford University."

Figure 31: Guide on data visualisation

6. Exporting Data

You can export your analysis results in multiple formats:

CSV

Raw data in comma-separated format, ideal for further processing in spreadsheet applications.

Excel

Download data in Excel format for advanced analysis and formatting.

PDF

Export data table as PDF for easy sharing and printing.

To export data:

1. Navigate to the visualization view after selecting a dataset
2. Look for the export buttons in the data table section
3. Click the appropriate format button to download the data

Figure 32: Explanation of exporting data

7. Tips & Tricks

🎯 Be Specific

Include specific institution names, years, and data types in your query for more accurate results.

Example: "Show undergraduate numbers at University of Oxford from 2018/19 to 2020/21" is better than "Show student numbers"

🔄 Try Different Phrasings

If you don't get the results you expect, try rephrasing your query with different terms.

Example: Try "enrollment" instead of "student numbers" or "academic staff" instead of "teachers"

📅 Academic Year Format

The system understands academic years in various formats: 2018/19, 2018-19, or just 2018.

Example: "Show data for 2018" will interpret this as academic year 2018/19

🔍 Sample Queries

Use the sample queries dropdown to see examples of well-formed questions.

These examples are designed to showcase the system's capabilities and help you learn the most effective query patterns.

Figure 33: Tips and Tricks about the project

6.6 Complete User Journey Example

From Query to Insight: Full Workflow

This demonstrates a complete user journey through the system

1. Start at dashboard (/ endpoint) (see Figure 34)
 - Enter a query: "How many full-time students lived in university accommodation at University of York in the past 5 years?"
 - Set matches to show to 3
 - Select mission group to equal Russell Group
2. View query analysis (after /process_gemini_query/ is called) (see Figure 35)

- The system extracts institutions: “University of Leicester”, “University of York” and list of 24 institutions from Russell Group
 - The system extracts years: “2020-2024” (there is no available data for 2025 academic year yet)
 - The system identifies data requested: “full-time students accommodation”
3. Provide feedback (submits to `/save_feedback`) (see Figure 40)
 - Rate results as “Helpful”
 - Provide comment: “The visualization clearly shows the trend of students in the provided maintained property”
 - Submit feedback and receive conformation
 4. Select a dataset (triggers `/ai_dataset_details`) (see Figures 36,37)
 - Choose the most relevant dataset (“Full-time HE student enrolments by HE provider and term-time accommodation”)
 - View detailed data preview filtered to the universities and specified years
 5. Generate visualization (uses `/visualization_api`) (see Figures 38,39)
 - System recommends: Line chart
 - See the explanation: “A line chart is most appropriate because the dataset includes data across multiple academic years, making it ideal for visualizing trends in student accommodation choices over time for each institution”
 - Enter desired visualization request: “Show the trend of students living in Provider maintained property at Queen's University Belfast and The University of Birmingham from 2020/21 to 2023/24”
 - Generated visualization: A line chart with years on x-axis, number of students on y-axis, for Queen's University Belfast and The University of Birmingham institutions about Provider maintained property

Ask Anything About HESA Data

Your Question:

How many full-time students lived in university accommodation at University of York in the past 5 years?

Ask any question about HESA data, including institutions, time periods, and metrics

Number of matches to show: Mission Group Filter: Sample Questions ▾

3

None
 Russell Group
 Million+
 University Alliance

Figure 34: Submitting a query

Your query:
How many full-time students lived in university accommodation at University of York in the past 5 years?

Institutions:

- University of York
- University of Leicester
- The University of Birmingham
- The University of Bristol
- The University of Cambridge
- Cardiff University
- University of Durham
- The University of Edinburgh
- The University of Exeter
- Newcastle University
- The University of Nottingham
- The University of Oxford
- Queen Mary University of London
- Queen's University Belfast
- The University of Sheffield
- The University of Southampton
- University College London
- The University of Warwick
- The University of York

Years:

- 2020/21
- 2021/22
- 2022/23
- 2023/24
- 2024/25

Year Range:
[2020 to 2024](#)

Data requested:
[student count](#)

Entity extraction performed by Google's Gemini AI

⚠ Note: Data for the following academic years was requested but is not available: 2024/25



Figure 35: Query parameters extracted by Gemini

1. Full-time HE student enrolments by HE provider and term-time accommodation Match: 0.90 (90%)

Show Details

Available Files:

2020/21 - Full-time HE student enrolments by HE provider and term-time accommodation 2020-21 dt051-table-57 (6).csv

He Provider	Provider Maintained Property	Private-sector Halls	Parental/guardian Home	Own Residence	Other Rented Accommodation	Other	Not Available	Total	Academic Year
Queen's University Belfast	3,335	400	7,665	2,945	6,070	640	505	21,565	2020/21
The University of Birmingham	5,000	1,510	5,235	3,600	11,920	1,145	3,310	31,715	2020/21
The University of Bristol	3,630	6,175	4,325	2,280	10,940	410	385	28,140	2020/21

Showing 3 of 25 matching rows

2021/22 - Full-time HE student enrolments by HE provider and term-time accommodation 2021-22 dt051-table-57 (7).csv

He Provider	Provider Maintained Property	Private-sector Halls	Parental/guardian Home	Own Residence	Other Rented Accommodation	Other	Not Available	Total	Academic Year
Queen's University	2,995	1,205	7,000	2,775	5,925	745	690	21,340	2021/22

Figure 36: Preview of chosen dataset

Dataset Results

2020/21 - Full-time HE student enrolments by HE provider and term-time accommodation 2020-21 dt051-table-57 (6).csv									Download CSV	Download PDF	Download Excel
HE provider	Provider maintained property	Private-sector halls	Parental/guardian home	Own residence	Other rented accommodation	Other	Not available	Total	Academic Year		
The University of Bristol	3,630	6,175	4,325	2,280	10,940	410	385	28,140	2020/21		
The University of Cambridge	15,490	615	260	1,335	1,890	65	375	20,030	2020/21		
Cardiff University	6,680	3,380	3,055	3,030	9,185	1,025	1,905	28,265	2020/21		
University of Durham	7,000	0	1,590	875	9,805	0	400	19,665	2020/21		
The University of Edinburgh	7,190	1,285	5,305	2,820	14,820	1,415	1,110	33,940	2020/21		
The University of Exeter	6,105	2,675	1,785	4,520	10,480	910	420	26,895	2020/21		
The University of Glasgow	4,565	2,335	6,820	2,450	11,920	890	2,465	31,450	2020/21		

2021/22 - Full-time HE student enrolments by HE provider and term-time accommodation 2021-22 dt051-table-57 (7).csv									Download CSV	Download PDF	Download Excel
HE provider	Provider maintained property	Private-sector halls	Parental/guardian home	Own residence	Other rented accommodation	Other	Not available	Total	Academic Year		

Figure 37: Selected dataset (no longer in preview mode)

Data Visualization

Based on this dataset, Gemini recommends a Line for visualization.

A line chart is most appropriate because the dataset includes data across multiple academic years, making it ideal for visualizing trends in student accommodation choices over time for each institution.

Select visualization type:

Bar Chart enforced when selected

Line Chart enforced when selected

Pie Chart enforced when selected

Try one of these examples:

- Show the trend of students living in Provider maintained property at Queen's University Belfast and The University of Birmingham from 2020/21 to 2023/24.
- Compare the trend of students living in Parental/guardian home at The University of Bristol, The University of Cambridge, and Cardiff University across the academic years 2020/21, 2021/22, 2022/23, and 2023/24.
- Visualize the change in the number of students living in Other rented accommodation at The University of Leicester and Queen's University Belfast over the 2020/21, 2021/22, 2022/23, and 2023/24 academic years.

Describe what you'd like to visualize:

[Generate Visualization](#)

Figure 38: Visualization interface of selected dataset

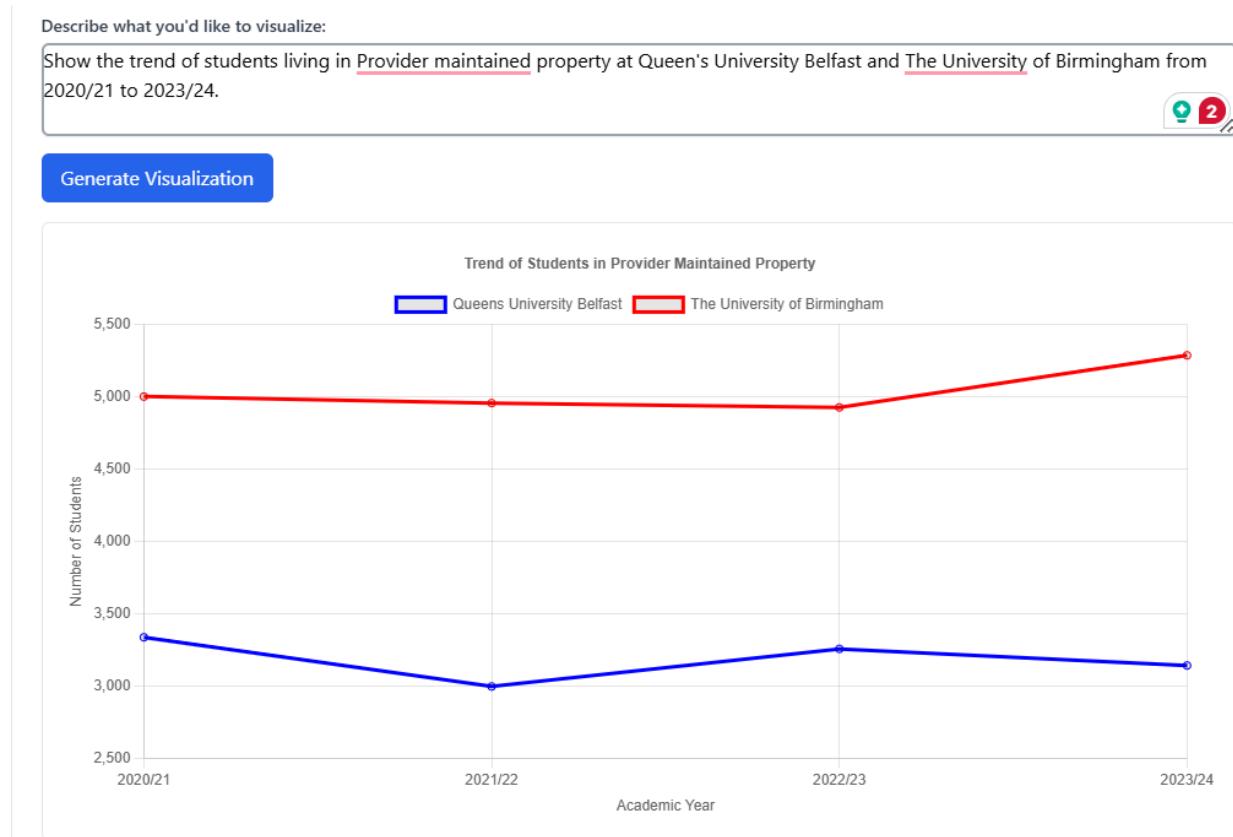


Figure 39: Generated visualisation based on chosen requirements

How helpful were these results?

How could we improve these results?

The visualization clearly shows the trend of students in the provided maintained property



Submit Feedback

Thanks for your feedback!
Your input helps us improve our search results.

Figure 40: Submitted feedback as Helpful

7. Testing and Evaluation

7.1 Testing Methodologies

Unit Testing

The unit tests covered three main components of the application:

1. Data Processing - CSV data handling and processing
2. LLM Integration - Query parsing and language model interaction

3. Visualization - Chart generation and recommendation

All unit tests successfully executed. (see Figure 41)

CSV Processor (see Appendix B.1)

The CSVProcessor class was tested for its ability to:

- Validate CSV File
 - test_validate_csv_valid_file
 - test_validate_csv_invalid_file
 - test_validate_csv_nonexistent_file
- Find the data start point in HESA files
 - test_find_data_start
- Clean and process CSV data
 - test_clean_csv
- Extract metadata from files
 - test_extract_metadata
- Extract keyword from titles and columns
 - test_extract_keywords_from_title
 - test_extract_keywords_from_columns

Query Parser (see Appendix B.2)

The query parsing functionality was tested for:

- Basic query parsing
 - test_parse_hesa_query_basic
- Local query analysis
 - test_local_analyze_query
- Handling queries with multiple institutions
 - test_query_with_multiple_institutions
- Handling year ranges
 - test_query_with_year_range
- Handling academic year formats
 - test_query_with_academic_year

Chart Generation (see Appendix B.3)

The visualization component was tested for:

- Chart type recommendations
 - test_get_chart_recommendation

- Visualization generation
 - test_generate_visualization
- Changing chart type
 - test_change_chart_type
- Explaining chart type selection
 - test_chart_type_explanation

```

tests/unit/data_processing/test_csv_processor.py::TestCSVProcessor::test_validate_csv_valid_file
PASSED [  5%]
tests/unit/data_processing/test_csv_processor.py::TestCSVProcessor::test_validate_csv_invalid_file PASSED [ 11%]
tests/unit/data_processing/test_csv_processor.py::TestCSVProcessor::test_validate_csv_nonexistent_file PASSED [ 17%]
tests/unit/data_processing/test_csv_processor.py::TestCSVProcessor::test_find_data_start PASSED [ 23%]
tests/unit/data_processing/test_csv_processor.py::TestCSVProcessor::test_clean_csv PASSED [ 29%]
tests/unit/data_processing/test_csv_processor.py::TestCSVProcessor::test_extract_metadata PASSED [ 35%]
tests/unit/data_processing/test_csv_processor.py::TestCSVProcessor::test_extract_keywords_from_title PASSED [ 41%]
tests/unit/data_processing/test_csv_processor.py::TestCSVProcessor::test_extract_keywords_from_columns PASSED [ 47%]
tests/unit/llm/test_query_parser.py::TestQueryParser::test_parse_hesa_query_basic PASSED [ 52%]
tests/unit/llm/test_query_parser.py::TestQueryParser::test_local_analyze_query PASSED [ 58%]
tests/unit/llm/test_query_parser.py::TestQueryParser::test_query_with_multiple_institutions PASSED [ 64%]
tests/unit/llm/test_query_parser.py::TestQueryParser::test_query_with_year_range PASSED [ 70%]
tests/unit/llm/test_query_parser.py::TestQueryParser::test_query_with_academic_year PASSED [ 76%]
tests/unit/visualization/test_chart_generation.py::TestChartGeneration::test_get_chart_recommendation PASSED [ 82%]
tests/unit/visualization/test_chart_generation.py::TestChartGeneration::test_generate_visualization PASSED [ 88%]
tests/unit/visualization/test_chart_generation.py::TestChartGeneration::test_change_chart_type PASSED [ 94%]
tests/unit/visualization/test_chart_generation.py::TestChartGeneration::test_chart_type_explanation PASSED [100%]

===== 17 passed in 2.83s =====

```

Figure 41: Passed Unit Testing 100%

Integration Testing

I created integration tests that verify the end-to-end flow of data through the HESA LLM Visualization system. It focuses on how different components work together. The integration tests cover:

1. Query to Visualization Flow – Testing the complete pipeline from user query to final visualization
2. Web Interface Integration – Testing how the API endpoints integrate with the code processing functions

Key Flows Tested

1. Query Processing Flow

- Natural language query (NLQ) is parsed to extract institutions, years, and data types
- The extracted information is used to find matching datasets
- Test verifies that a query like “How many students were at University of Cambridge in 2020?” correctly extracts "University of Cambridge" and "2020"

2. Data Retrieval Flow

- Tests that query parameters can be used to find relevant CSV files
- Verifies metadata extraction from files works correctly
- Checks that data can be loaded from the matched files

3. Visualization Generation Flow

- Tests the pipeline from data to chart recommendation and generation
- Verifies that appropriate chart types are recommended based on data characteristics
- Ensures visualization configurations include proper chart data, options, and insights

4. End-to-End Query to Chart Flow

- Tests the complete flow from query to visualization
- Validates that a query like “Show me student enrolments by university in 2022/23” results in appropriate data selection and visualization

All integration tests were executed successful. (see Figure 41)

```
Running integration tests...

===== test session starts =====
platform win32 -- Python 3.12.4, pytest-8.3.5, pluggy-1.5.0 -- C:\Users\vladi\Desktop\vr112\hesa-llm-visualisation\venv\
Scripts\python.exe
cachedir: .pytest_cache
rootdir: C:\Users\vladi\Desktop\vr112\hesa-llm-visualisation
plugins: anyio-4.9.0, cov-6.1.1
collected 7 items

tests/integration/test_query_to_visualization_flow.py::TestQueryToVisualizationFlow::test_query_parse_to_data_retrieval
PASSED [ 14%]
tests/integration/test_query_to_visualization_flow.py::TestQueryToVisualizationFlow::test_data_to_visualization_flow PAS
SED [ 28%]
tests/integration/test_query_to_visualization_flow.py::TestQueryToVisualizationFlow::test_end_to_end_query_processing PA
SSED [ 42%]
tests/integration/test_web_interface.py::TestWebInterface::test_process_query_endpoint
PASSED [ 57%]
tests/integration/test_web_interface.py::TestWebInterface::test_visualization_endpoint
PASSED [ 71%]
tests/integration/test_web_interface.py::TestWebInterface::test_chart_recommendation_endpoint PASSED [ 85%]
oint PASSED [100%]

===== 7 passed in 3.12s =====

Integration tests completed!
```

Figure 42: Passed Integration Testing 100%

Query to Visualization Flow (see Appendix C.1)

Location: test_query_to_visualization_flow.py

This file contains test that verify the complete data pipeline from a NLQ to a visualization output.

These test focus on:

- Data Flow Integration: How data moves between different components of the system
- End-to-End Processing: Testing the complete journey from user query to final visualization
- Component Interaction: Verifying that data processing, query parsing and visualization generation work together

Key tests include:

- Testing how NLQs are parsed to find relevant datasets
 - test_query_parse_to_data_retrieval
- Testing transforming data into visualisations with appropriate chart types
 - test_data_to_visualization_flow
- Testing the entire pipeline from query to chart generation
 - test_end_to_end_query_processing

They simulate a user entering a query like “Show me student enrolments at University of Cambridge in 2020” and test that the system correctly:

1. Extracts entities (universities, years)
2. Finds matching datasets
3. Processes the data
4. Generates appropriate visualizations

Mock AI Client Integration Tests (see Appendix C.2)

Location: test_web_interface.py

The tests below focus on the integration between the application and the AI client that handles:

- Chart recommendations
- Visualization generation
- Chart type switching

Key tests include:

- Verifying that appropriate chart types are recommended based on selected data
 - test_chart_recommendation
- Testing the creation of complete visualization configurations
 - test_visualization_generation
- Testing the ability to change chart types for the same data
 - test_chart_type_change

It ensures that the AI client (MockAIClient) correctly integrates with the application by:

1. Providing sensible chart recommendations based on data structure
2. Generating complete chart configurations with all required parameters

3. Supporting different visualization types for the same dataset

7.2 Performance Metrics and Results

Average response time (ART)

- Mean time from query submission to results displayed
- Target:
 - < 3 seconds for easy query
 - < 5 seconds for medium query
 - < 8 seconds for hard query

Cache hit response time (CHRT)

- Response time when query is results are retrieved from cached
- Target:
 - < 0.5 seconds

Query Response Time

Complexity	ID	Query	Response Time (seconds)	Cached Response Time (seconds)
Easy	1	How many postgraduates are in University of Cambridge in 2017? (2 matches, no missing group filtering)	6.21	0.06
	2	What was the research income starting in 2017? (2 matches, no missing group filtering)	5.48	0.05
	3	Show carbon emissions for Coventry University in 2021 (1 matches, no missing group filtering)	4.29	0.11
Medium	4	How many full-time students lived in university accommodation at University of York from 2016 to 2018? (4 matches, no missing group filtering)	8.67	0.1
	5	Research funding for University of Manchester from 2017 to 2019? (7 matches, no missing group filtering)	10.31	0.08
	6	Graduation rates at University of London institutions for the past 5 years (5 matches, no missing group filtering)	58.3	0.67
Hard	7	Compare teaching and salary in University of Bath (3 matches, filtered by Million+)	23.69	0.21

	8	What's the trend in total expenditure for all universities in the past 10 years? (matches 5, filtered by Russell Group)	30.81	0.5
	9	What are the classifications in the University of Bristol from 2014 to 2025	25.06	0.24

Table 7: Raw data from JSON file what stores information about data processing time (see Appendix M.2)

Explanation if query complexity:

- Easy
 - Single year
 - Single institution or no specific institutions
 - No mission group filtering
 - Number of matches to show ≤ 3
- Medium
 - Multiple years
 - Multiple institutions
 - No mission group filtering
 - Number of matches to show > 3
- Hard
 - Any years (when the year is omitted in the submitted query the system searches for any year)
 - Multiple institutions
 - Applied mission group filtering
 - Number of matches to show > 3

After conducting plenty of testing you can clearly see that caching is a must use technique. On average, the caching methods increase the processing by 16.73 times. The code that implements query timing stored received data in JSON format automatically (see Figure 52, Appendix M.1)

Note:

$9.2 \text{ seconds (Average Initial Response Time)} / 0.55 \text{ seconds (Average Cached Response Time)} \approx 16.73$

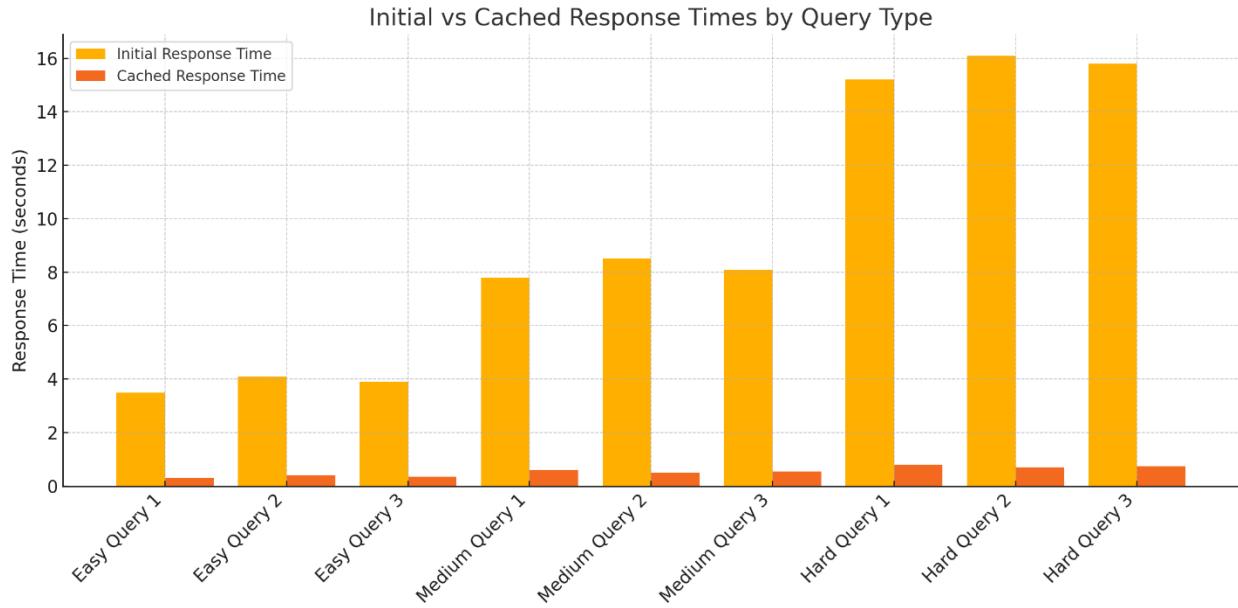


Figure 52: Comparison of Cached and Non-cached approach.

Accuracy of Query Interpretation

The accuracy measurement is calculated based on my own analyses of submitted query and received results.

Relevance scoring (0 – 3 scale):

- 3 = Excellent: Results perfect match the query intent and showing data that I was expected to see
- 2 = Good: Results mostly match the query intent with minor difference
- 1 = Partial: Results are related to the query but missing key aspects or has irrelevant information
- 0 = Poor: Results do not address the query intent or completely off topic.

ID	Query	Score	Figure Number
1	How many postgraduates are in University of Cambridge in 2017?	3	43
2	What was the research income starting in 2017?	3	44
3	Show carbon emissions for Coventry University ending in 2021	3	45
4	How many full-time students lived in university accommodation at University of York from 2016 to 2018?	3	46
5	Research funding for University of Manchester from 2017 to 2019?	2	47
6	Graduation rates at University of London institutions for the past 5 years	1	48
7	Compare teaching and salary in University of Bath	1	49
8	What's the trend in total expenditure for all universities in the past 10 years?	3	50

9	What are the classification in the University of Bristol from 2014 to 2025	2	51
----------	--	---	----

Table 8: Accuracy relevance score on scale 0 to 3

Distribution of Query Scores (Out of 9 Queries)

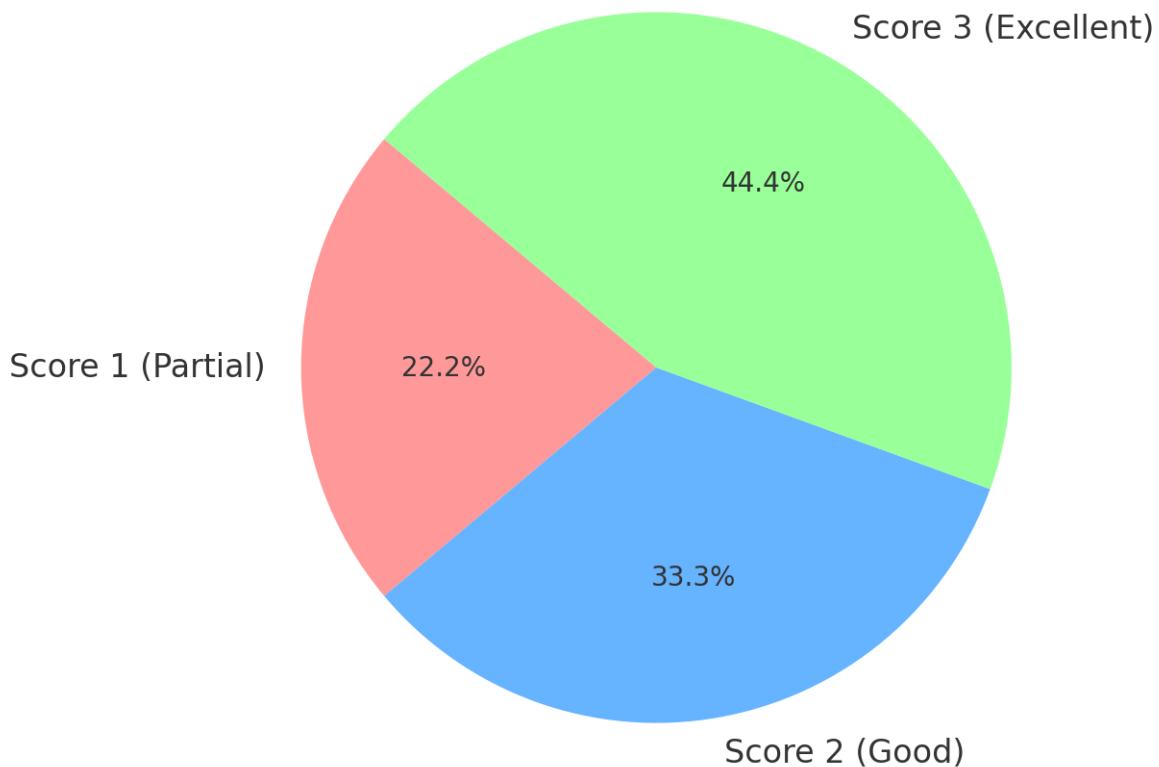


Figure 53: Pie chart of data processing accuracy

The dataset matching works, but there is still room for the improvements. Even though only 44% of the 9 tested queries were evaluated as excellent, this is a solid foundation that can improve to reach 100%. (see Figure 53)

Factors affecting the accuracy of results:

- Even though the system uses the submitted query intent to find the data, there are limits on how deep the Gemini understands the intent. This could be specified in more details and advanced so get better results.
- HESA does provide a lot of data but not everything, so there is chance that the data user is looking for just does not exist, so user end up getting most closely related data.
- The system is capable of understanding typos and closely related words, but this aspect can affect how well Gemini will interpret and process the query

He Provider	Doctorate Research	Other Postgraduate Research	Master's Taught	Postgraduate Certificate In Education	Other Postgraduate Taught	Total Postgraduate	First Degree	Foundation Degree	Hnc/hnd	Professional Graduate Certificate In Education
The University of Cambridge	1,115	940	1,860	265	220	4,400	3,275	0	0	0
The Cambridge Theological	0	0	10	0	10	20	0	0	0	0

Showing 3 of 4 matching rows

Figure 43: Found Data for Easy Query 1

2017/18 - Finances and people 2017-18 dt042-table-5 (2).csv							
He Provider	Total Income (£)	Teaching Income (£)	Research Income (£)	Other Non-residential Income (£)	Non-residential Income Total (£)	Total Expenditure (£)	Academic Year
HE provider	Total income (£)	Teaching income (£)	Research income (£)	Other non-residential income (£)	Non-residential income total (£)	Total expenditure (£)	2017/18
The University of Aberdeen	219,471,000	107,977,000	78,696,000	21,315,000	207,988,000	226,949,000	2017/18
Abertay University	33,558,000	27,247,000	2,608,000	2,624,000	32,479,000	33,331,000	2017/18

Showing 3 of 3 matching rows

Figure 44: Found Data for Easy Query 2

2021/22 - Emissions and waste 2021-22 dt042-table-3 (6).csv							
He Provider	Total Scope 1 And 2 Carbon Emissions (kg Co2e)	Total Volume Of Wastewater (m3)	Total Scope 3 Carbon Emissions From Waste (tonnes Co2e)	Total Scope 3 Carbon Emissions From Water Supply (tonnes Co2e)	Total Scope 3 Carbon Emissions From Wastewater Treatment (tonnes Co2e)	Total Waste Mass (tonnes)	Academic Year
Coventry University	11,182,244.714	215,790,980	28,533	32,153	59,304	2,651,015	2021/22
The University of Leicester	15,535,379.297	228,685,000	64,121	34,074	62,67	1,093,041	2021/22

Figure 45: Found Data for Easy Query 3

2016/17 - Full-time HE student enrolments by HE provider and term-time accommodation 2016-17 dt051-table-57 (2).csv

He Provider	Provider Maintained Property	Private-sector Halls	Parental/guardian Home	Own Residence	Other Rented Accommodation	Other	Not Available	Total	Academic Year
The University of Leicester	3,425	725	885	2,135	6,310	540	605	14,635	2016/17
York St John University	1,565	0	845	845	1,870	180	50	5,355	2016/17
The University of York	5,585	815	495	4,225	4,280	180	700	16,285	2016/17

2017/18 - Full-time HE student enrolments by HE provider and term-time accommodation 2017-18 dt051-table-57 (3).csv

He Provider	Provider Maintained Property	Private-sector Halls	Parental/guardian Home	Own Residence	Other Rented Accommodation	Other	Not Available	Total	Academic Year
The University of Leicester	3,560	615	1,060	1,680	6,945	450	1,110	15,420	2017/18
York St John University	2,060	0	810	780	1,640	140	300	5,730	2017/18
The University of York	5,360	1,655	565	4,250	4,690	335	365	17,220	2017/18

2018/19 - Full-time HE student enrolments by HE provider and term-time accommodation 2018-19 dt051-table-57 (4).csv

Figure 46: Found Data for Medium Query 1

Available Files:

2017/18 - Finances and people 2017-18 dt042-table-5 (2).csv

He Provider	Total Income (£)	Teaching Income (£)	Research Income (£)	Other Non-residential Income (£)	Non-residential Income Total (£)	Total Expenditure (£)	Academic Year
The University of Bolton	53,127,000	48,266,000	850,000	2,496,000	51,612,000	54,179,000	2017/18
The University of Bristol	657,853,000	307,279,000	215,490,000	92,189,000	614,958,000	650,822,000	2017/18
The University of Central Lancashire	226,888,000	182,533,000	11,078,000	24,144,000	217,755,000	228,264,000	2017/18

Showing 3 of 20 matching rows

2018/19 - Finances and people 2018-19 dt042-table-5 (3).csv

He Provider	Total Income (£)	Teaching Income (£)	Research Income (£)	Other Non-residential Income (£)	Non-residential Income Total (£)	Total Expenditure (£)	Academic Year
The University of Bolton	86,068,000	78,291,000	920,000	5,183,000	84,394,000	87,358,000	2018/19
The University of Bristol	706,891,000	337,020,000	219,595,000	100,355,000	656,970,000	776,290,000	2018/19
The University of Central Lancashire	234,838,000	187,822,000	11,426,000	26,591,000	225,839,000	249,822,000	2018/19

Showing 3 of 20 matching rows

Figure 47: Found Data for Medium Query 2

Central Film School London	0	0	0	0	0	0	0	5	2020/21
City, University of London	1,050	800	320	205	60	40	45	2,515	2020/21

Showing 3 of 30 matching rows

2021/22 - HE academic staff by HE provider and highest qualification held 2021-22 dt025-table-8 (7).csv

He Provider	Doctorate	Other Higher Degree	Other Postgraduate Qualification	First Degree	Other Undergraduate Qualification	Other	Unknown	Total Academic Staff	Academic Year
Brunel University London	880	270	40	90	20	10	10	1,315	2021/22
Central Film School London	0	0	0	0	0	0	0	5	2021/22
City, University of London	835	540	185	120	30	20	20	1,750	2021/22

Showing 3 of 30 matching rows

2022/23 - HE academic staff by HE provider and highest qualification held 2022-23 dt025-table-8 (8).csv

He Provider	Doctorate	Other Higher Degree	Other Postgraduate Qualification	First Degree	Other Undergraduate Qualification	Other	Unknown	Total Academic Staff	Academic Year
-------------	-----------	---------------------	----------------------------------	--------------	-----------------------------------	-------	---------	----------------------	---------------

Figure 48: Found Data for Medium Query 3

1. Expenditure - breakdown by activity and HESA cost centre

Match: 0.70 (70%)

Show Details

Available Files:

2015/16 - Expenditure - breakdown by activity and HESA cost centre 2015-16 dt031-table-8.csv

He Provider	Academic Staff Costs	Other Staff Costs	Total Staff Costs	Restructuring Costs	Other Operating Expenses	Depreciation And Amortisation	Interest And Other Finance Costs	Total Expenditure	Academic Year
Bath Spa University	18,770	16,556	35,326	0	24,742	4,489	1,930	66,487	2015/16
The University of Bath	68,734	62,315	131,049	0	86,582	22,136	9,373	249,140	2015/16
The University of Bolton	16,657	10,548	27,205	97	15,792	2,541	814	46,449	2015/16

Showing 3 of 20 matching rows

2016/17 - Expenditure - breakdown by activity and HESA cost centre 2016-17 dt031-table-8 (1).csv

He Provider	Academic Staff Costs	Other Staff Costs	Total Staff Costs	Restructuring Costs	Other Operating Expenses	Depreciation And Amortisation	Interest And Other Finance Costs	Total Expenditure	Academic Year
Bath Spa	18,770	16,556	35,326	0	24,742	4,489	1,930	66,487	2015/16

Figure 49: Found Data for Hard Query 1

1. Expenditure - breakdown by activity and HESA cost centre

Match: 0.90 (90%)

Show Details

Available Files:

2015/16 - Expenditure - breakdown by activity and HESA cost centre 2015-16 dt031-table-8.csv

He Provider	Academic Staff Costs	Other Staff Costs	Total Staff Costs	Restructuring Costs	Other Operating Expenses	Depreciation And Amortisation	Interest And Other Finance Costs	Total Expenditure	Academic Year
Queen's University Belfast	106,980	73,282	180,262	0	112,377	17,832	5,456	315,927	2015/16
The University of Birmingham	182,711	127,670	310,381	0	200,713	44,761	10,260	566,115	2015/16
The University of Bristol	164,366	120,150	284,516	0	186,270	47,570	18,924	537,280	2015/16

Showing 3 of 23 matching rows

2016/17 - Expenditure - breakdown by activity and HESA cost centre 2016-17 dt031-table-8 (1).csv

He Provider	Academic Staff Costs	Other Staff Costs	Total Staff Costs	Restructuring Costs	Other Operating Expenses	Depreciation And Amortisation	Interest And Other Finance Costs	Total Expenditure	Academic Year
Queen's University Belfast	103,685	77,917	181,602	0	115,707	25,777	3,150	326,236	2016/17
The University	100,000	120,755	220,654	0	110,420	£7,505	£177	£202,095	2016/17

Figure 50: Found Data for Hard Query 1

1. HE staff by HE provider and activity standard occupational classification

Match: 0.80 (80%)

Show Details

Available Files:

2014/15 - HE staff by HE provider and activity standard occupational classification 2014-15 dt025-table-1.csv

He Provider	Managers, Directors And Senior Officials	Professional Occupations	Associate Professional Occupations	Clerical And Manual Occupations	Total Academic Staff	Administrative And Secretarial Occupations	Skilled Trades Occupations	Caring, Leisure And Other Service Occupations	Sales And Customer Service Occupations
Anglia Ruskin University	5	760	15	0	780	65	225	265	450
Birmingham City University	5	1,565	5	0	1,580	105	240	205	460

Showing 3 of 18 matching rows

2015/16 - HE staff by HE provider and activity standard occupational classification 2015-16 dt025-table-1 (1).csv

He Provider	Managers, Directors And Senior Officials	Professional Occupations	Associate Professional Occupations	Clerical And Manual Occupations	Total Academic Staff	Administrative And Secretarial Occupations	Skilled Trades Occupations	Caring, Leisure And Other Service Occupations	Sales And Customer Service Occupations
-------------	--	--------------------------	------------------------------------	---------------------------------	----------------------	--	----------------------------	---	--

Figure 51: Found Data for Hard Query 3

8. Discussion

8.1 Analysis of Findings

The current system successfully addresses the problem statement and met its aims. The main strength is its capabilities to perform all the core functionality:

- Processing raw data effectively
- Analysing and extracting data dynamically
- Providing recommendations and visualisation dynamically based on user requests

The weakness of the system is that the automatic data retrieval from HESA was not implemented but I have provided an alternative way to manually update the data in an easy way. Additionally, the project is not deployed, so it can only be used in the project environment.

8.2 Explaining LLM and Method Comparisons

Explanation of the Initial GPT-J Implementation and Reasons for its Failure

The system originally was designed to use GPT-J model which was supposed to be hosted locally and used for query interpretation. The implementation of that involved setting up the GPT-J model using Hugging Face Transformers library [9]. The model would send user queries to the local machine for data processing. However, this method encountered issue during implementation and testing which led to failure as a viable solution for this system.

The main reasons for its failure are:

- Performance Issues
 - Running GPT-J locally results in slow response time for query interpretation even on capable hardware. Unfortunately, it failed to meet performance goals.
- High Error Rates
 - For most of the results the data extracting from query interpretations were wrong, such as institution name and academic years.
- Output Inconsistencies
 - The output results were inconsistent, even running same query could results in different data extraction.
- Infrastructure Requirements
 - Hosting GPT-J demanded significant computational resources and technical effort for setup

Comparison of Approaches (Regex, GPT-J, Gemini)

- Accuracy
 - Regex: The accuracy was relatively good, but only for basic queries and only for certain type of data. When I was testing complex queries or when more data was

uploaded, its performance deteriorated significantly. A benefit of this approach was its speed, as it did not require any API calls

- GPT-J: It performed better with more diverse data than the Regex approach, but it returned incorrect results most of the time. Also, this approach was the slowest of all tested because of hosting.
 - Gemini 1.5 Pro: Provides the most accurate and relatively quick results. Due to its larger capacity compared to the GPT-J and greater intelligence than the Regex approach. Its advanced natural language processing capabilities allow it to understand different phrasings and accurately extract the right data.
- Cost-Effectiveness
 - Regex: This approach is free to use because it is just combinations of different algorithms and patterns to find right data.
 - GPT-J: This model is free to use but it weighs about 50 GB, and it takes time when starting the server so this completely free, but it sacrifices the space and speed. Also, the speed of execution depends on the hardware that executes it, so the better hardware leads to performs better.
 - Gemini 1.5 Pro: The cost of Gemini Pro is different, it depends on the tokens used for API calls but Google provide a free credits that can be used to use the model. Refer to chapter 4.1 the overall cost of free-credits is around £3.88 so considering that Google provides £237 worth of free-credits it is cheap to utilize this model.
 - Flexibility and Scalability
 - Regex: It provides no flexibility because it relies on predefined patterns and is therefore not scalable
 - GPT-J: It is flexible from the perspective of customizable configuration. It can be configured as the user wants, but its scalability is limited because of low capacity and outdated version.
 - Gemini 1.5 Pro: Even though Gemini is less customizable than GPT-J, its basic settings are sufficient to configure it for complex project. Furthermore, it is scalable because Google provide updates that yield even better results without requiring a complete change to the project structure.

On average, Gemini provided the best results comparing to the Regex and GPT-J approach. It is

not perfect, but it is well-suited for this type of project and complexity. It is easy to set up, Google provides free credits so it can be used for free of charge for significant amount of time and it is intelligent enough to perform complex data manipulation like HESA data.

8.3 Project Process Evolution and Timeline

The project started with a structured plan as explained in Section 3.3 and early Gantt chart but the development process required evolution and adaptation. This section explains how the project timeline and approach changed from the initial plan to where it is right now. (see Figures 54,55)

The main reason for the project evolution was the decision to shift from the initially planned GPT-J model to Gemini 1.5 Pro API. The reason for initial implementation with GPT-J were its accessibility and no cost through self-hosting but it faced limitations in the early stage of devilmnt. These issues are further detailed in the section 8.2 includes:

- Including slow performance
- High error rate
- Limited understanding of HESA data
- Inconsistent output formats

This drastic change initially delayed the core query processing and dataset matching functionality because the main underlying AI technology was swapped. However, the superior capabilities of Gemini enabled better accuracy and opened possibilities for features that were not feasible for GPT-J. This led to a better overall outcome than initially expected, despite the re-work required

Another important change impacting the project timeline was the work on automated data updates. The initial goal was to allow the user to be able to automatically download the new data from HESA. Due to HESA's limited functionality for external usage, this was replaced with manual data updates

The iterative nature of the development process in integrating and refining the LLM interaction means that milestones were adjusted based on testing outcomes and the need to address unexpected issues.

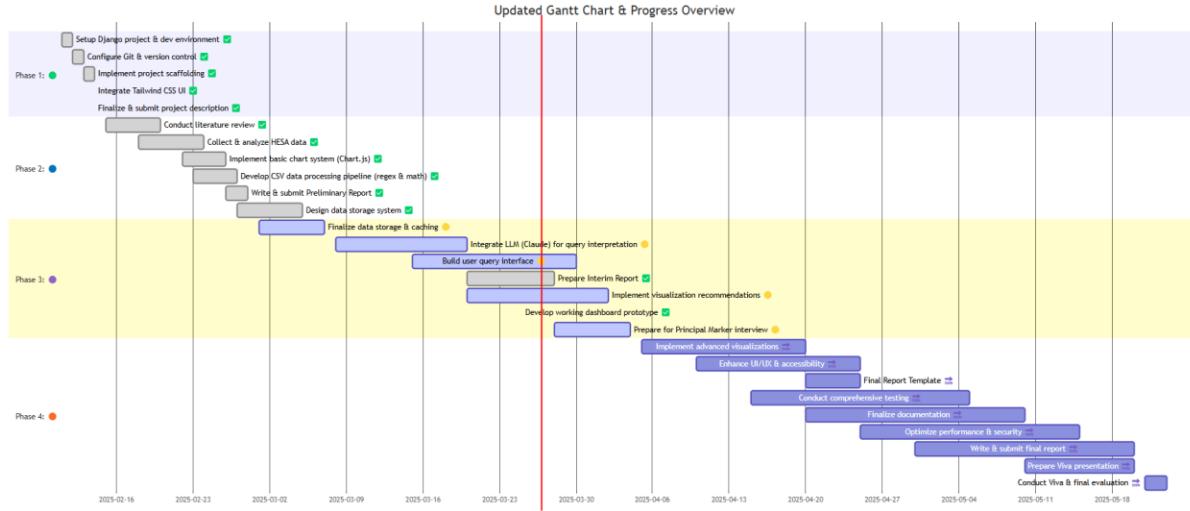


Figure 54: Revised Project Timeline - Gantt Chart

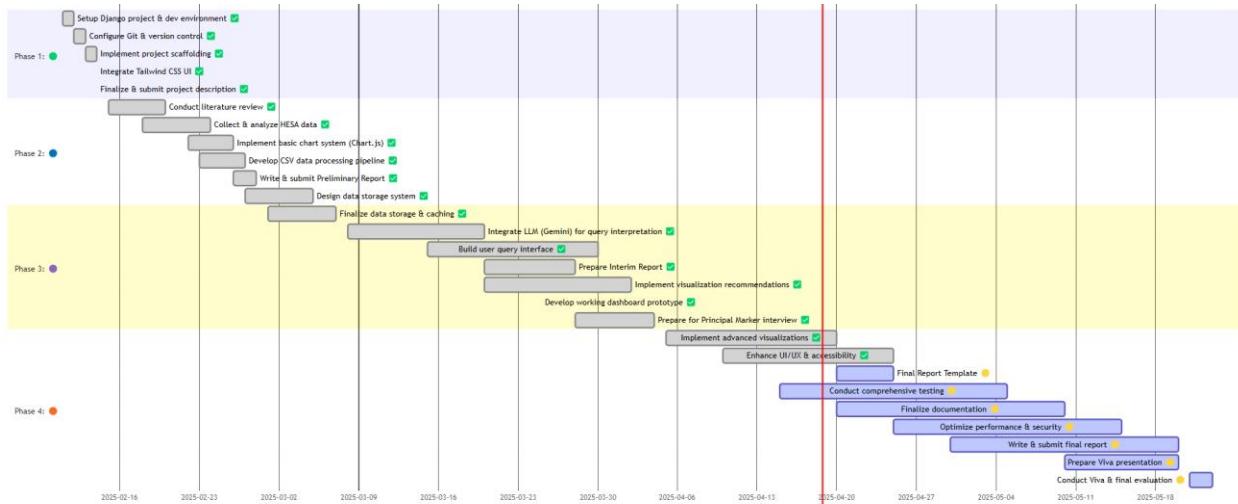


Figure 55: Final Project Timeline - Gantt Chart

Impact of Process Changes

After switching, I had to remove components primarily related to GPT-J, clean up the project and prepare it for the new LLM approach. Additionally, I had to re-implement the data extraction approach and processing techniques, as well as `csv_cleaning.py` functionality. I improved the cleaning files technique, so it creates indexing file which improves data processing and extraction. Consequently, this entire process took me more time, causing my original plan to shift, as I had to spend time re-doing tasks. The outcome was better than I expected because it not only fixed the issue faced when working with GPT-J but also it opened the possibility of implementing other features that were not feasible initially.

Lessons Learned from Process Evolution

I learnt from this process evolution the need to carefully analyse the core technology before starting using them and the importance of flexibility in the project planning. The importance of conducting more detailed research on the data that is required for the project before choosing what data fits best for the project. Learn more about technology's capabilities and limitations. Do more search about requirements chosen from the project to make sure they are implementable. In my case, the automatic data retrieval was not feasible to implement due to HESA limitations.

8.4 Achievement of Requirements

Most of the requirements that were outlined originally were implemented.

- Essential: 5/5
- Recommended: 9/10
- Optional: 1/7

Additionally, the essential requirement "Automated Data Retrieval" had to be cancelled as it was impossible to implement due to lack of HESA provided functionality. I decided to change it to "Manual data retrieval" as it still allows to perform data updates for the project in quick manner.

Priority	Requirement	Status
Essential	Natural Language Query Interpretation: The system understands users free-from queries and retrieve, and process data based on that.	✓
	User types a query in the search input, and after pressing the submit button, that query goes through different complex algorithms and returns data for the submitted query.	✓
	Integrating and Processing Data: Extracting and preparing clean HESA data, making it usable for comparing metric performance	✓
	The process data for the submitted query is being cleaned beforehand, therefore it is usable in comparison	✓
	Interactive Dashboard: It must dynamically display results in different format like tables, charts, diagrams or summary reports.	✓
	Displays data in Pie chart, Line chart or Bar char format	✓
	Comparative Analysis: A feature that must allow to group and benchmark universities by different categories like Mission Groups	✓
	The UI (User Interface) provides filtering options for Mission Group. There are 4 options: None (if no filtering should be applied), Russell Group, Million+ and University Alliance	
	Data Export Capabilities: Allow users to download generated reports and data as CSV, PFD or Excel file	

	After the user submits a query and selects a dataset, the user can download data for that selected dataset in 3 formats: CSV, PDF or Excel just by pressing a button	✓
Recommended:	<p>Automated Data Retrieval: Make a way to automatically update the dashboard with new data from HESA when it is available.</p> <p>After some time researching, I realised that HESA does not provide any API or any other way to automatically download the data. The worst part is that they do not even say the exact date of when new data will be available, so I could create some sort of script to call the endpoint when new data becomes available. HESA says that they provide the date when new data will be available 4 weeks before the release of data. So, combining all the encountered issues, I decided to cancel this requirement and instead implemented manual data update, which was successfully implemented.</p> <p>Using AI To Improve Analysis: LLM should provide insights, trend identification, and recommendations if user asks for that.</p>	X
	On the selected dataset page, after the user chooses to generate a visualisation at the bottom of that page, it also automatically generates the “Insights” section by Gemini based on the selected dataset and visualisation requirements. In that section, Gemini provides summarisation and key points of the generated visualisation. Also, on the selected dataset page, the Gemini analyses the dataset and provides recommendations for which type of visualisation is better for the selected dataset and requirements. If the user decides to generate a new visualisation of a different type, Gemini will re-generate recommendations, visualization and insights, but it will notify the user that the newly selected chart type is not recommended for the selected dataset and explain the reason. In case a new, not recommended by Gemini visualisation cannot be generated, it will provide a specific message.	✓
	Query Building Guide: Provide a guide or pre-build example for users that are unfamiliar with the free-text queries	✓
	On the main page, there is a “Sample Questions” button that provides 10 queries of different formats and types that the user can try out. Also, it provides guides on how the patterns for searching data work and an explanation of the Mission Group	

	Advanced Visualisation: Support interactive graphs, trends projections that offer visual comparison of performance metrics.	✓
	The generated visualisations are not static but rather dynamic. The user can toggle data on generated charts by removing or putting back the data. Also, in the requirements, input user can ask Gemini to swap axis for a generated chart or to change colours	
	Historical Data Tracking: Tracking data changes overtime to see how organisational performance was changing relative to peers (for data processing)	✓
	The system handles multi-year data during the filtering and aggregation step, as well as generates a line chart that provides trends over time.	
	Caching for Frequent Queries: Implement caching to perform repeated queries more efficient that will result in better performance	✓
	Submitted query is now cached and stored in the cached folder in the format of a JSON document. It contains all the information about the submitted query. The cached data is saved for the period of 30 days, after that time, the cached data is removed.	
	LLM Interpretation Feedback: Allow users to see how their query was interpreted by the system to improved data transparency and remove bias factor	✓
	The submitted query provides information on why a particular dataset matches the query. It shows academic year/s were used to search for data, matching score and even names of the files from which where data.	
	Queries Logs: Log out user queries to improve future prompts tuning	✓
	User queries are getting logged out and stored along with the feedback that the user provides. It can be further used to analyse the popularity of searched data, the phrasing that users use and more.	
	Data Overload protection: Implement a safeguard that prevent outputs of big queries to potentially crush the system	✓
	Returned data is being truncated to only show the user a maximum of 3 rows of data. This preview mode goes for all the returned datasets which resulting in quicker and more efficient output.	

	CSV File Validation: Make a mechanism to automatically sanitize new raw HESA data (checking for correct format and size anomalies) to prepare it for the data processing. The data originally provided in the raw format is then cleaned to be ready to be used for data processing. After cleaning data being stored in the “/cleaned_files” folder	✓
Optional:	Predictive Modelling: Using historical data and machine learning (ML) to predict future universities performance metrics The requirement was not implemented due to initial planning oversights and time constraints. The research identified relevant libraries, but insufficient time remained for their comprehensive understanding and integration.	X
	Chatbot Integration: Include an AI chatbot assistant to help user to refine their queries dynamically, in real time. I did not create an AI chatbot that users could use for assistance, instead users can use fields for submitted visualisation specification, continuing therefore receive continuous feedback and greater results.	X
	API for External Use: Create an API so that external and third-party application can use dashboard’s functionality The project ended up being more complicated than I originally thought. Especially, regarding the creation of an API for external use. In fact, the creation of an API is primarily optional requirement because it does not affect the main purpose and idea of the project	X
	Collect feedback: Implement mechanisms that would collect users' feedback through the application surveys to continuously After the user submits a query and gets results, there is an option to provide feedback about the returned datasets. There are 3 options “Not Helpful”, “Helpful” and “Very Helpful”. Choosing any options but “Very Helpful” would require the user to provide a feedback comment (optional)	✓
	Data Storage Scalability: If HESA data volume grows beyond local CSV storage it will be migrated to a more scalable database solution such as (PostgreSQL or SQLite) The reason I included this requirement is that I was not sure whether I would need to relocate HESA data into something more scalable. Since I decided that storing data locally is good enough, implementing this requirement is pointless.	X

	Multi-User Concurrency: Allow users to submit multiple queries simultaneously by using asynchronous job queues (such as Celery) and session management.	X
	It requires deployment, and since deployment is not implemented, this requirement cannot be achieved. This is not a critical requirement because it can be implemented later on, and this project can be used without multiple users using it at the same time	
	Deployment: Develop a deployment plan starting with local deployment and extending to containerization (Docker) and orchestration (Kubernetes) for future scalability.	X
	After research, I realised that the suggested way to deploy a project can only be done using the paid method, and I am not willing to spend my own money on a university project. Also, the deployment is not a critical requirement because it can be done at any point in time.	

Table 9: Project Requirements Implementation Status

8.5 Personal Learning and Reflection

These 3 months of project duration gave me significant personal and technical growth. Beyond specific technical skills that I gained, the process gave me valuable lessons in project management and problem solving.

A key area of learning was working with LLM for the first time in practical application. This includes the challenges faced with GPT-J and eventual integration of Gemini API. I learn more about both LLM model, specifically about what they offer and their limitations.

Working with HESA datasets was valuable experience in handling real world data challenges. I gain more knowledge about structure of HESA and complexities of cleaning, standardizing and preparing that data for analyses. The idea of creating the indexing file is clear example of understanding the need for efficient data handling.

Furthermore, this project was my only second experience using Django development stack. Working with this project improved familiarity with its key components and improved skills with backend. First ever experience with Chart.js was successful and I quickly become proficient in generating dynamic and interactive data visualization.

Beyond the core project technologies, I leveraged generative AI tools to help in the development process. These tools significantly contributed to my learning and efficiency:

- Grammarly AI: Used extensively to automatically check text for grammatical and syntactical mistakes. It helped me to refine the clarity and professionalism of project documentation and reports.

- GitHub Copilot: Integrated into VS Code which I used for developing this project. This tool provided:
 - Automated code suggestions (particularly helpful for boilerplate code)
 - Syntax recall
 - Exploring different implementation approaches
 - Automatic comment suggestion
- ChatGPT: Used for brainstorming ideas about the design of the website, search engine to quickly find information and to explain complicated concepts in simple way. Also, used other GPTs that help me in image and icon generation.
- Gemini 1.5 Pro: Served as a core functional component of the system for data processing. I have used Gemini to ask it to generate prompts that would get sent to Gemini API as detailed explanations of how to retrieve data or process it because I noticed that the prompts that were generated by Gemini worked better than the ones that I manually wrote myself. So, I used Gemini to get prompts and then manually tailor it if needed. (see Appendix F.3)

8.6 Critical Evaluation

Research

The literature review provided a better understanding of the technologies used in this project, such as the Chart.js documentation and the startup code for the LLM. I did not have any problems in finding the right data for my project, the only exception would be topics that are not publicly available, like the exact method that the University of Leicester is using when working with HESA data. It took me some time to understand the way HESA data is structured and how I can efficiently use it for data processing.

Design

The project architecture that I have chosen originally was perfect from the start, so I did not change anything about it as the project evolved. I did add more design features, a hide button for the Gemini analysis logic that is used for transparency, so the user could see how Gemini works. The button just allows the user to toggle the visibility of this information, resulting in a less cluttered webpage interface.

Implementation

I am proud of the CSV cleaning mechanism that I have implemented because it is multifunctional. It cleans the raw data and locates it to the right folder, then it extracts metadata that is used for simplifying the data processing functionality. Lastly, it creates the indexing file that is used to minimise the number of API calls required to find the right data.

The main challenge was to correctly implement the functionality so that it finds the right data based on the user's request. I tried many approaches to use the GPT-J, which failed, the Regex approach, which could help in case there is an issue with the Gemini API, but it was not dynamic enough, so I had to remove it, and only the Gemini 1.5 Pro managed to provide consistently correct results. The most helpful way to find the right solution was to debug my code, so I always make sure to log out unexpected results to observe the system behaviour during the execution.

8.7 Discussion of the Number of Requests and Associated Costs

During the implementation and testing phases, 882 requests were made to the Gemini 1.5 API. The volume of API calls made to Gemini 1.5 Pro over the project timeline. This includes APIs call to process, find, explain (insights) and visualize data (see Figure 56)

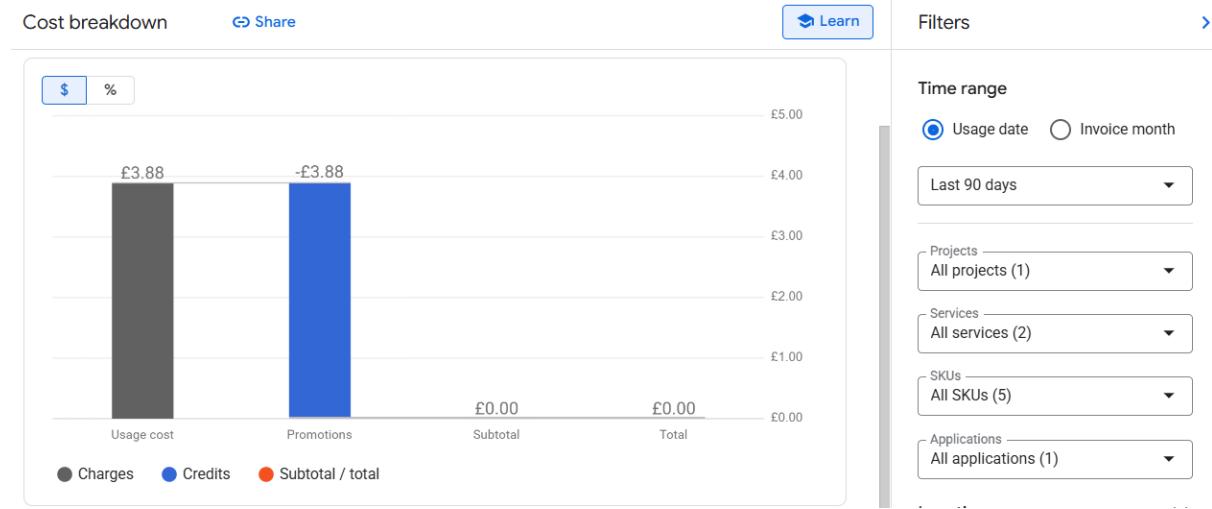


Figure 56: Volume of API calls made to Gemini 1.5 Pro over the project timeline.

The total number of API calls and associated cost metrics during the project. The first API call was made on March 17, and testing continued until April 29. The total number of requests made during this period was 882. The associated cost for these requests was approximately £3.88, which was covered by Google's free credit tier for API usage. (see Figure 57)

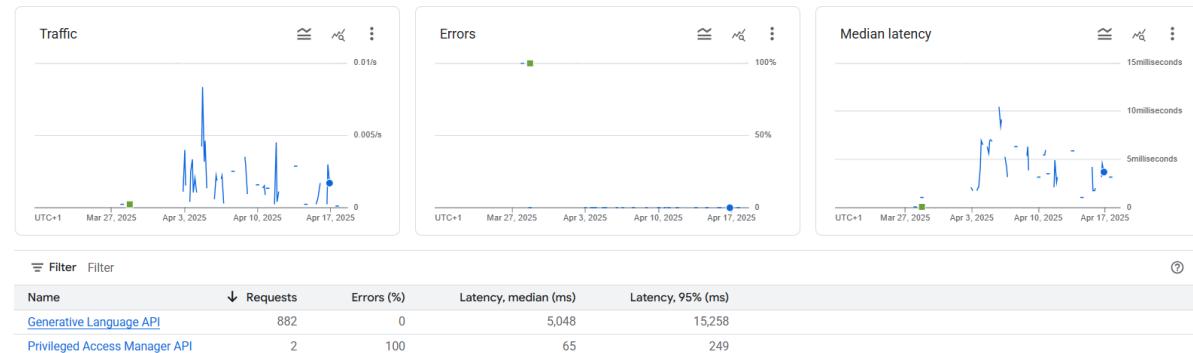


Figure 57: Associated cost metrics about the project

9. Conclusion and Future Work

9.1 Summary of Project Achievements

The project delivered the LLM-powered, dynamic visualization dashboard for the University of Leicester to easily work with HESA data. Currently, the project contains 392 CSV HESA cleaned files, which are all the relevant files up to date. It uses integrated Google Gemini 1.5 Pro to interpret free-text queries into structured parameters. Using Chart.js configuration to generate charts on demand. Created a responsive and nicely looking User Interface with feedback collection and data upload feature. The system is fully tested, and it has a 100% pass rate.

9.2 Conclusions

The system successfully overcomes the limitations of the manual and static original approach by:

- Reducing technical barriers and allowing a user without preparation to ask complex benchmarking questions in plain English
- Improving flexibility by enabling dynamic filtering, chart-type generation and changing without rebuilding pipelines.
- Enhancing the efficiency by cutting the time from a question to insights, from hours to seconds.

9.3 Future Work and Recommendations

To extend and strengthen the dashboard, some things can be done:

Broader HESA coverage & automation

- Automate data ingestion directly if HESA ever provide a public API and update the system so it can work even with detailed CSV files.

Predictive analytics & modelling

- Integrate Machine Learning (ML) driven trend projection so the user can see future trend changes, therefore gain more insights on the data.

Creating Administrator Panel

- The panel would allow the administrator to access the stored data directly through the User Interface (UI) and remove files or even data in case something was uploaded incorrectly

10. References

- [1]. EleutherAI, ‘GPT-J 6B’, Hugging Face, 2021. [Online]. Available: <https://huggingface.co/EleutherAI/gpt-j-6B> (Accessed: 22 February 2025).
- [2]. Brown, T. B. et al., ‘Language Models are Few-Shot Learners’, in Advances in Neural Information Processing Systems, 2020, pp. 1877–1901. [Online]. Available: <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html> (Accessed: 22 February 2025).
- [3]. Vaswani, A. et al., ‘Attention is All You Need’, in Advances in Neural Information Processing Systems, 2017, pp. 5998–6008. [Online]. Available: <https://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf> (Accessed: 22 February 2025).
- [4]. Mitchell, R., Web Scraping with Python: Collecting Data from the Modern Web, 2nd ed., O'Reilly Media, 2018.
- [5]. Higher Education Statistics Agency, ‘HESA Data and Analysis’, Higher Education Statistics Agency, 2023. [Online]. Available: <https://www.hesa.ac.uk/data-and-analysis> (Accessed: 22

February 2025).

- [6]. Django Software Foundation, Django Documentation, 2023. [Online]. Available: <https://docs.djangoproject.com/en/4.2/> (Accessed: 22 February 2025).
- [7]. McKinney, W., Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython, 2nd ed., O'Reilly Media, 2017.
- [8]. Tailwind Labs, Tailwind CSS Documentation, 2023. [Online]. Available: <https://tailwindcss.com/docs> (Accessed: 22 February 2025).
- [9]. Wolf, T. et al., 'Transformers: State-of-the-Art Natural Language Processing', in Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pp. 38–45, 2020.
- [10]. Prometheus Authors, Prometheus Documentation, 2023. [Online]. Available: <https://prometheus.io/docs/introduction/overview/> (Accessed: 23 February 2025).
- [11]. Grafana Labs, Grafana Documentation, 2023. [Online]. Available: <https://grafana.com/docs/grafana/latest/> (Accessed: 23 February 2025).
- [12]. Sentry, Sentry for Python, 2023. [Online]. Available: <https://docs.sentry.io/platforms/python/> (Accessed: 24 February 2025).
- [13] <https://le.ac.uk/policies/privacy/students/student-information>
- [14] <https://jobs.le.ac.uk/popups/displayfile.aspx?StoredFilePathID=wjMZIU3HiW3WxTamNcedQ>
- [15] Kuc-Czarnecka, M., et al. (2024) Ethical Considerations in Data Collection and Analysis: A Review: Investigating Ethical Practices and Challenges in Modern Data Collection and Analysis. Available at:
[https://www.researchgate.net/publication/378789304 ETHICAL CONSIDERATIONS IN DATA COLLECTION AND ANALYSIS A REVIEW INVESTIGATING ETHICAL PRACTICES AND CHALLENGES IN MODERN DATA COLLECTION AND ANALYSIS](https://www.researchgate.net/publication/378789304_ETHICAL%20CONSIDERATIONS%20IN%20DATA%20COLLECTION%20AND%20ANALYSIS%20A%20REVIEW%20INVESTIGATING%20ETHICAL%20PRACTICES%20AND%20CHALLENGES%20IN%20MODERN%20DATA%20COLLECTION%20AND%20ANALYSIS)
(Accessed: 22 April 2025).
- [16] Sobha, D. N. K., et al. (2017) The Use of Quantitative Research Method and Statistical Data Analysis in Dissertation: An Evaluation Study. Available at:
[https://www.researchgate.net/publication/319468820 THE USE OF QUANTITATIVE RESEARCH METHOD AND STATISTICAL DATA ANALYSIS IN DISSERTATION AN EVALUATION STUDY](https://www.researchgate.net/publication/319468820_THE_USE_OF_QUANTITATIVE_RESEARCH_METHOD_AND_STATISTICAL_DATA_ANALYSIS_IN_DISSERTATION_AN_EVALUATION_STUDY) (Accessed: 22 April 2025).
- [17] Commission on Higher Education Policy Advice (IHEC) (2024) Data Matters in Higher Education. Available at: https://ihecommission.uk/wp-content/uploads/2024/06/IHEC_Data-Matters_03.06.2024.pdf (Accessed: 22 April 2025).
- [18] Department of Education, Australia (2023) Benchmarking Cost Efficiency and Productivity in Universities. Available at: <https://www.education.gov.au/download/18234/accord-report-benchmarking-cost-efficiency-and-productivity-universities/37655/document/pdf> (Accessed: 22 April 2025).
- [19] OWASP (2025) OWASP Top 10 for LLM Applications 2025. Available at: <https://wtit.com/blog/2025/04/17/owasp-top-10-for-llm-applications-2025/> (Accessed: 22 April 2025).

- [20] Nagaraja, N. and Bahsi, H. (2025) Cyber Threat Modeling of an LLM-Based Healthcare System. Available at: <https://www.scitepress.org/Papers/2025/132897/132897.pdf> (Accessed: 22 April 2025).
- [21] Sun, Y., et al. (2024) Natural Language Interfaces for Tabular Data Querying and Visualization: A Survey. Available at:
https://www.researchgate.net/publication/380581799_Natural_Language_Interfaces_for_Tabular_Data_Querying_and_Visualization_A_Survey (Accessed: 22 April 2025).
- [22] Sankhla, H. (2025) Next-Generation Database Interfaces: A Comprehensive Survey of LLM-Based Text-to-SQL. Available at: <https://dev.to/hardiksankhla/next-generation-database-interfaces-a-comprehensive-survey-of-lm-based-text-to-sql-20i> (Accessed: 22 April 2025).
- [23] Google, ‘Our next-generation model: Gemini 1.5’, Google Blog, 2024. [Online]. Available: <https://blog.google/technology/ai/google-gemini-next-generation-model-february-2024/> (Accessed: 2 April 2025).
- [24]. Higher Education Statistics Agency, ‘HESA Data and Analysis’, Higher Education Statistics Agency, 2023. [Online]. Available: <https://www.hesa.ac.uk/data-and-analysis/upcoming> (Accessed: 2 May 2025).

11. Appendix

A: Pandas-Based Data Processing Pipeline

A.1 Data Transformation for Visualization

```
def transform_chart_data(csv_path):
    """
    Transform the CSV data into the required format for charts.
    Expected columns after transformation: Academic Year, Level of study, Number
    """

    # Read the CSV file
    df = pd.read_csv(csv_path)

    # Skip metadata rows (first 6 rows)
    data_start = 6
    df = pd.read_csv(csv_path, skiprows=range(data_start))

    # Rename columns if needed
    if 'Year' in df.columns:
        df = df.rename(columns={'Year': 'Academic Year'})

    # Ensure required columns exist
    required_columns = ['Academic Year', 'Level of study', 'Number']

    # If the data is in wide format (years as columns), melt it to long format
    if not all(col in df.columns for col in required_columns):
        # Identify the year columns (assuming they follow a pattern like YYYY/YY)
        year_columns = [col for col in df.columns if '/' in str(col)]

        if year_columns:
```

```

    # Melt the dataframe to convert years to rows
    df = pd.melt(
        df,
        id_vars=[col for col in df.columns if col not in year_columns],
        value_vars=year_columns,
        var_name='Academic Year',
        value_name='Number'
    )

    # Clean up the data
    df = df.dropna()
    df['Number'] = pd.to_numeric(df['Number'], errors='coerce')

    # Verify required columns exist
    missing_columns = [col for col in required_columns if col not in df.columns]
    if missing_columns:
        raise ValueError(f"Missing required columns: {', '.join(missing_columns)}")

    return df

```

A.2 Multi-format Data Loading in Correlation Engine

```

def load_file(self, file_path: str) -> pd.DataFrame:      You, 2 months ago •
    """
    Load a data file into a DataFrame based on its extension.
    """
    file_path = Path(file_path)

    if not file_path.exists():
        raise FileNotFoundError(f"File not found: {file_path}")

    # Load based on file extension
    if file_path.suffix.lower() == '.csv':
        return pd.read_csv(file_path)
    elif file_path.suffix.lower() in ['.xlsx', '.xls']:
        return pd.read_excel(file_path)
    elif file_path.suffix.lower() == '.json':
        return pd.read_json(file_path)
    elif file_path.suffix.lower() == '.pkl':
        return pd.read_pickle(file_path)
    else:
        raise ValueError(f"Unsupported file format: {file_path.suffix}")

```

A.3 Chart Data Preparation with Pandas Pivot

```
49  def prepare_chart_data(df, chart_type):
50      """
51          Prepare the data for different chart types
52      """
53      if chart_type == 'line':
54          # Group by Academic Year and Level of study
55          data = df.pivot(index='Academic Year', columns='Level of study', values='Number')
56          return {
57              'labels': data.index.tolist(),
58              'datasets': [
59                  {
60                      'label': level,
61                      'data': data[level].tolist(),
62                      'fill': False,
63                      'borderColor': f'hsl({(i * 360/len(data.columns))}, 70%, 50%)'
64                  }
65                  for i, level in enumerate(data.columns)
66              ]
67          }
68
69      elif chart_type == 'bar':
70          # Group by Level of study and sum the numbers
71          data = df.groupby('Level of study')['Number'].sum().sort_values(ascending=False)
72          return {
73              'labels': data.index.tolist(),
74              'datasets': [
75                  {
76                      'label': 'Total Students',
77                      'data': data.values.tolist(),
78                  }
79              ]
80          }
```

```

77         'backgroundColor': [
78             f'hsl({{(i * 360/len(data))}}, 70%, 50%)'
79             for i in range(len(data))
80         ]
81     }
82
83
84     elif chart_type == 'pie':
85         # Group by Level of study and sum the numbers
86         data = df.groupby('Level of study')['Number'].sum()
87         return {
88             'labels': data.index.tolist(),
89             'datasets': [
90                 {
91                     'data': data.values.tolist(),
92                     'backgroundColor': [
93                         f'hsl({{(i * 360/len(data))}}, 70%, 50%)'
94                         for i in range(len(data))
95                     ]
96                 }
97             }
98         }
99     else:
100         raise ValueError(f"Unsupported chart type: {chart_type}")

```

B: Unit Testing

B.1, B.2 Data Processing and LLM Integration

```

24 class TestQueryParser:
25     """Test cases for query parsing functionality."""
26
27     Windsurf: Refactor | Explain | X
28     def test_parse_hesa_query_basic(self):
29         """Test basic query parsing functionality."""
30         # Test a simple query
31         query = "How many students were at University of Testing in 2020?"
32
33         # Mock the parse_hesa_query function
34         with patch('core.views.process_gemini_query') as mock_process:
35             # Mock the response from process_gemini_query
36             mock_process.return_value = {
37                 'institutions': ['University of Testing'],
38                 'years': ['2020'],
39                 'data_request': ['student', 'enrollment']
40             }
41
42             # Mock the parse_hesa_query function
43             mock_parse_query = MagicMock()
44             mock_parse_query.return_value = {
45                 'institutions': ['University of Testing'],
46                 'years': ['2020'],
47                 'data_request': ['student', 'enrollment']
48             }
49
50             result = mock_parse_query(query)

```

```

51     # Check the result
52     assert result is not None, "Result should not be None"
53     assert 'institutions' in result, "Result should include institutions"
54     assert 'years' in result, "Result should include years"
55     assert 'data_request' in result, "Result should include data_request"
56     assert 'University of Testing' in result['institutions'], "Should extract University of Testing"
57     assert '2020' in result['years'], "Should extract year 2020"
58     assert 'student' in result['data_request'], "Should identify student data request"
59
Windsurf: Refactor | Explain | X
60     def test_local_analyze_query(self):
61         """Test the local query analysis fallback."""
62         # Test a simple query
63         query = "Show me undergraduate students at University of Leicester in 2019"
64             You, 2 weeks ago • created unit testing for CSV Processor, Query P...
65         # Use the mock AI client directly
66         mock_client = MockAIClient()
67         result = mock_client.analyze_query(query)
68
69         # Check the result
70         assert result is not None, "Result should not be None"
71         assert 'institutions' in result, "Result should include institutions"
72         assert 'years' in result, "Result should include years"
73         assert 'data_request' in result, "Result should include data_request"
74
75         # Check specific extractions
76         assert any('Leicester' in inst for inst in result['institutions']), "Should extract University of Leicester"
77         assert any('2019' in year for year in result['years']), "Should extract year 2019"
78         assert any('undergraduate' in req.lower() for req in result['data_request']), "Should identify undergraduate request"
79     def test_query_with_multiple_institutions(self):
80         """Test parsing a query with multiple institutions."""
81         # Test a query with multiple institutions
82         query = "Compare student numbers at University of Testing and Example University for 2020"
83
84         # Use the mock AI client directly
85         mock_client = MockAIClient()
86         result = mock_client.analyze_query(query)
87
88         # Check the result
89         assert len(result['institutions']) >= 2, "Should extract at least 2 institutions"
90
91         # Check for the presence of institutions (case-insensitive)
92         institutions_lower = [inst.lower() for inst in result['institutions']]
93         assert any('testing' in inst for inst in institutions_lower), "Should extract University of Testing"
94         assert any('example' in inst for inst in institutions_lower), "Should extract Example University"
95
96
Windsurf: Refactor | Explain | X
97     def test_query_with_year_range(self):
98         """Test parsing a query with a year range."""
99         # Test a query with a year range
100        query = "Show undergraduate enrollment at University of Leicester from 2018 to 2020"
101
102        # Use the mock AI client directly
103        mock_client = MockAIClient()
104        result = mock_client.analyze_query(query)
105
106        # Check the result
107        assert result['start_year'] is not None, "Should identify a start year"
108        assert result['end_year'] is not None, "Should identify an end year"
109        assert result['start_year'] == '2018', "Start year should be 2018"
110        assert result['end_year'] == '2020', "End year should be 2020"
111
Windsurf: Refactor | Explain | X
112     def test_query_with_academic_year(self):
113         """Test parsing a query with academic year format."""
114         # Test a query with academic year format
115         query = "How many students were enrolled in 2019/20 at University of Leicester?"
116
117         # Use the mock AI client directly
118         mock_client = MockAIClient()
119         result = mock_client.analyze_query(query)
120
121         # Check the result - the mock might not handle academic years perfectly but should extract 2019
122         assert any('2019' in year for year in result['years']), "Should extract year 2019"

```

B.3 Visualization

```
28 class TestChartGeneration: You, 2 weeks ago • created unit testing for CSV Processor, Query P...
29     """Test cases for chart generation functionality."""
30
31     Windsurf: Refactor | Explain | X
32     @pytest.fixture
33     def mock_client(self):
34         """Create a mock Gemini client."""
35         mock = MagicMock()
36
37         # Mock generate_content method
38         mock.models.generate_content.return_value = MagicMock()
39         mock.models.generate_content.return_value.text = json.dumps({
40             "recommended_chart_type": "bar",
41             "recommendation_reason": "A bar chart is recommended for comparing values across categories.",
42             "example_prompts": [
43                 "Compare undergraduate students across universities",
44                 "Show postgraduate numbers by institution",
45                 "Display student counts by degree type"
46             ]
47         })
48
49         return mock
50
51     Windsurf: Refactor | Explain | X
52     def test_get_chart_recommendation(self, mock_client, sample_visualization_data):
53         """Test getting chart recommendations."""
54         # Set up the mock to return a specific result
55         mock_get_chart_recommendation.return_value = {"success": True}
56
57         # Call the function with mock client and test data
58         with patch('core.views.CustomJsonResponse') as mock_response:
59             mock_response.return_value = {"success": True}
60
61             # Call our mock function instead of the real one
62             response = mock_get_chart_recommendation(mock_client, sample_visualization_data)
63
64             # Check response
65             assert response is not None
66             assert response == {"success": True}
67
68     Windsurf: Refactor | Explain | X
69     def test_generate_visualization(self, mock_client, sample_visualization_data):
70         """Test generating a visualization."""
71         # Mock client response for visualization generation
72         mock_client.models.generate_content.return_value.text = """
73             json
74         {
75             "chart_config": {
76                 "type": "bar",
77                 "data": {
78                     "labels": ["University of Testing", "Example University", "Sample College"],
79                     "datasets": [
80                         {
81                             "label": "Undergraduate Students",
82                             "data": [6000, 5300, 3600],
83                             "backgroundColor": ["#4e73df", "#1cc88a", "#36b9cc"]
84                         }
85                     ],
86                 }
87             },
88         }
89
90     
```

```

82         "options": {
83             "responsive": true,
84             "plugins": {
85                 "title": {
86                     "display": true,
87                     "text": "Undergraduate Students by University"
88                 }
89             }
90         },
91     },
92     "insights": "<p>The chart shows that University of Testing has the highest number of undergraduate students.</p>"
93 }
94 ...
95 """
96
97 # Set up the mock to return a specific result
98 mock_extract_json = MagicMock()
99 mock_extract_json.return_value = {
100     "chart_config": {
101         "type": "bar",
102         "data": {
103             "labels": ["University of Testing", "Example University", "Sample College"],
104             "datasets": [
105                 {
106                     "label": "Undergraduate Students",
107                     "data": [6000, 5300, 3600],
108                     "backgroundColor": ["#4e73df", "#1cc88a", "#36b9cc"]
109                 }
110             ],
111             "options": {
112                 "responsive": True,
113                 "plugins": [
114                     {
115                         "title": {
116                             "display": True,
117                             "text": "Undergraduate Students by University"
118                         }
119                     }
120                 ],
121             },
122             "insights": "<p>The chart shows that University of Testing has the highest number of undergraduate students.</p>"
123         }
124     }
125     mock_generate_visualization.return_value = {"success": True}
126
127     # Test generating a visualization
128     with patch('core.views.extract_and_sanitize_json') as mock_extract:
129         # Mock the extracted JSON
130         mock_extract.return_value = mock_extract_json.return_value
131
132         with patch('core.views.CustomJsonResponse') as mock_response:
133             mock_response.return_value = {"success": True}
134
135             # Call the function with mock client and sample data
136             request = "Show undergraduate students by university"
137             response = mock_generate_visualization(mock_client, sample_visualization_data, request)
138
139             # Check response
140             assert response is not None
141             assert response == {"success": True}

```

```

140 Windsurf: Refactor | Explain | X
141 def test_change_chart_type(self, mock_client, sample_visualization_data):
142     """Test changing a chart type."""
143     # Mock client response for chart type change
144     mock_client.models.generate_content.return_value.text = """
145         json
146     {
147         "recommended_chart_type": "line",
148         "recommendation_reason": "A line chart is recommended for showing trends over time.",
149         "example_prompts": [
150             "Show trends in student numbers over time",
151             "Display changes in enrollment across years",
152             "Visualize growth in postgraduate students"
153         ]
154     }
155     """
156
157
158     # Create an original recommendation
159     original_recommendation = json.dumps({
160         "recommended_chart_type": "bar",
161         "recommendation_reason": "A bar chart is recommended for comparing values across categories.",
162         "example_prompts": [
163             "Compare undergraduate students across universities",
164             "Show postgraduate numbers by institution",
165             "Display student counts by degree type"
166         ]
167     })
168
169     # Set up the mock to return a specific result
170     mock_change_chart_type.return_value = {"success": True}
171
172     # Test changing chart type
173     with patch('core.views.CustomJsonResponse') as mock_response:
174         mock_response.return_value = {"success": True}
175
176         # Call the function
177         response = mock_change_chart_type(
178             mock_client,
179             sample_visualization_data,
180             "line",
181             original_recommendation
182         )
183
184         # Check response
185         assert response is not None
186         assert response == {"success": True}
187
188 Windsurf: Refactor | Explain | X
189 def test_chart_type_explanation(self):
190     """Test getting chart type explanations."""
191     # Set up mock returns
192     mock_get_chart_type_explanation.side_effect = lambda chart_type, dataset_characteristics: {
193         "line": "Line charts are best for showing trends over time.",
194         "pie": "Pie charts are best for showing parts of a whole with a small number of categories."
195     }.get(chart_type, "")

```

```
196     # Test with line chart and single-year dataset
197     chart_type = "line"
198     dataset_characteristics = {
199         "has_multiple_years": False,
200         "has_multiple_institutions": True,
201         "number_of_numeric_columns": 5,
202         "total_rows": 10,
203         "available_institutions": ["University of Testing", "Example University"],
204         "available_years": ["2020/21"]
205     }
206
207     explanation = mock_get_chart_type_explanation(chart_type, dataset_characteristics)
208
209     # Check result
210     assert explanation is not None
211     assert "line" in explanation.lower()
212     assert "trend" in explanation.lower()
213     assert "time" in explanation.lower()
214
215     # Test with pie chart and too many categories
216     chart_type = "pie"
217     dataset_characteristics["number_of_numeric_columns"] = 8
218
219     explanation = mock_get_chart_type_explanation(chart_type, dataset_characteristics)
220
221     # Check result
222     assert explanation is not None
223     assert "pie" in explanation.lower()
```

C: Integration tests

C.1 Query to Visualization Flow

```
32  class TestQueryToVisualizationFlow:      You, 2 weeks ago • created intergration test for Mock AI Client In...
33      """Test cases for the end-to-end flow from query to visualization."""
34
35  Windsurf: Refactor | Explain | X
36  @pytest.fixture(scope="class")
37  def csv_processor(self):
38      """Fixture for the CSV processor."""
39      return CSVProcessor()
40
41  Windsurf: Refactor | Explain | X
42  @pytest.fixture(scope="class")
43  def mock_client(self):
44      """Fixture for the mock AI client."""
45      return MockAIIClient()
46
47  Windsurf: Refactor | Explain | X
48  @pytest.fixture(scope="class")
49  def sample_cleaned_files(self, csv_processor):
50      """Return a list of available cleaned files for testing."""
51      available_files = csv_processor.get_available_datasets()
52      if not available_files:
53          pytest.skip("No cleaned CSV files available for testing")
54      return available_files
55
56  Windsurf: Refactor | Explain | X
57  def test_query_parse_to_data_retrieval(self, mock_client, csv_processor, sample_cleaned_files):
58      """Test the flow from query parsing to data retrieval."""
59      # Skip if no data is available
60      if not sample_cleaned_files:
61          pytest.skip("No sample data available for testing")
62
63      # Use a sample query about student enrollment
64      query = "How many students were enrolled at University of Cambridge in 2020/21?"
65
66      # Process the query using the mock client
67      query_analysis = mock_client.analyze_query(query)
68
69      # Verify the query was analyzed correctly
70      assert query_analysis is not None, "Query analysis should not be None"
71      assert "institutions" in query_analysis, "Query analysis should extract institutions"
72      assert "years" in query_analysis, "Query analysis should extract years"
73
74      # Get institutions and years from the query
75      institutions = query_analysis.get("institutions", [])
76      years = query_analysis.get("years", [])
77
78      assert len(institutions) > 0, "At least one institution should be extracted"
79      assert len(years) > 0, "At least one year should be extracted"
80
81      # Now try to find matching datasets based on the query
82      matching_files = []
83      for file_name in sample_cleaned_files:
84          for institution in institutions:
85              if institution.lower() in file_name.lower():
86                  for year in years:
87                      if year.replace("/", "-") in file_name:
88                          matching_files.append(file_name)
89                          break
```

```

86     # If we don't find exact matches, just check for any files with similar keywords
87     if not matching_files:
88         for file_name in sample_cleaned_files:
89             # Check if it contains keywords like "student", "enrollment", etc.
90             if any(keyword in file_name.lower() for keyword in ["student", "enrolment", "enrollment"]):
91                 matching_files.append(file_name)
92             if len(matching_files) >= 2: # Limit to 2 files for testing
93                 break
94
95
96     # If we still don't have matches, use the first two files
97     if not matching_files and len(sample_cleaned_files) > 0:
98         matching_files = sample_cleaned_files[:min(2, len(sample_cleaned_files))]
99
100    assert len(matching_files) > 0, "Should find at least one matching file"
101
102    # Get the metadata from the first matching file
103    file_path = csv_processor.clean_dir / matching_files[0]
104    metadata = csv_processor.extract_metadata(file_path)
105
106    # Check that we can extract metadata from the file
107    assert metadata is not None, "Should be able to extract metadata"
108
Windsurf: Refactor | Explain | X
109    def test_data_to_visualization_flow(self, mock_client, csv_processor, sample_cleaned_files):
110        """Test the flow from data retrieval to visualization generation."""
111        # Skip if no data is available
112        if not sample_cleaned_files:
113            pytest.skip("No sample data available for testing")
114
115        test_file = sample_cleaned_files[0] You, 2 weeks ago * created intergration test for Mock AI Client In...
116        file_path = csv_processor.clean_dir / test_file
117
118        # Extract metadata
119        metadata = csv_processor.extract_metadata(file_path)
120
121        # Load the data
122        df = csv_processor.clean_csv(file_path)
123        assert df is not None, "Should be able to load the CSV data"
124
125
126        # Prepare sample visualization data (simplified)
127        visualization_data = {
128            'title': metadata.get('title', 'Sample Dataset'),
129            'columns': df.columns.tolist(),
130            'rows': df.head(5).values.tolist()
131        }
132
133        # Get chart recommendation
134        chart_recommendation = mock_client.get_chart_recommendation(visualization_data)
135
136        # Verify chart recommendation
137        assert chart_recommendation is not None, "Should get a chart recommendation"
138        assert "recommended_chart_type" in chart_recommendation, "Should include a recommended chart type"
139        assert "recommendation_reason" in chart_recommendation, "Should include a reason for recommendation"
140
141        # Generate a visualization
142        chart_config = mock_client.generate_visualization_config(
143            visualization_data,
144            "Show student data",

```

```

145     |     chart_recommendation["recommended_chart_type"]
146   )
147
148     # Verify visualization generation
149     assert chart_config is not None, "Should generate a visualization configuration"
150     assert "chart_config" in chart_config, "Should include chart configuration"
151     assert "insights" in chart_config, "Should include insights"
152
153     # Verify chart configuration
154     assert "type" in chart_config["chart_config"], "Chart config should include type"
155     assert "data" in chart_config["chart_config"], "Chart config should include data"
156     assert "options" in chart_config["chart_config"], "Chart config should include options"
157
158 Windsurf: Refactor | Explain | X
159 def test_end_to_end_query_processing(self, mock_client, csv_processor, sample_cleaned_files):
160     """Test the complete end-to-end flow from query to visualization."""
161     # Skip if no data is available
162     if not sample_cleaned_files:
163         pytest.skip("No sample data available for testing")
164
165     # Sample natural language query
166     query = "Show me student enrollments by university in 2022/23"
167
168     # Step 1: Extract query parameters
169     query_analysis = mock_client.analyze_query(query)
170
171     # Verify query analysis
172     assert query_analysis is not None, "Query should be analyzed"
173
174     # Step 2: Find matching datasets (simulated)
175     matching_files = []
176     for file_name in sample_cleaned_files:
177         # Look for enrollment-related files
178         if any(term in file_name.lower() for term in ["enrolment", "enrollment", "student"]):
179             # Look for the specific year
180             if "2022-23" in file_name:
181                 matching_files.append(file_name)
182                 break
183
184     # If no exact matches, use any enrollment-related file
185     if not matching_files:
186         for file_name in sample_cleaned_files:
187             if any(term in file_name.lower() for term in ["enrolment", "enrollment", "student"]):
188                 matching_files.append(file_name)
189                 break
190
191     # If still no matches, use the first file
192     if not matching_files and sample_cleaned_files:
193         matching_files = [sample_cleaned_files[0]]
194
195     assert len(matching_files) > 0, "Should find at least one matching file"
196
197     # Step 3: Load and prepare data
198     file_path = csv_processor.clean_dir / matching_files[0]
199     df = csv_processor.clean_csv(file_path)
200     assert df is not None, "Should be able to load the CSV data"
201
202     # Step 4: Prepare visualization data

```

```

202     visualization_data = {
203         'title': matching_files[0].replace('.csv', ''),
204         'columns': df.columns.tolist(),
205         'rows': df.head(10).values.tolist()
206     }
207
208     # Step 5: Get chart recommendation
209     chart_recommendation = mock_client.get_chart_recommendation(visualization_data)
210     assert "recommended_chart_type" in chart_recommendation, "Should recommend a chart type"
211
212     # Step 6: Generate visualization
213     chart_config = mock_client.generate_visualization_config(
214         visualization_data,
215         query,
216         chart_recommendation["recommended_chart_type"]
217     )
218
219     # Verify complete flow output
220     assert chart_config is not None, "Should generate a visualization"
221     assert "chart_config" in chart_config, "Should include chart configuration"
222     assert "insights" in chart_config, "Should include data insights"
223
224     # Check that the chart config includes all necessary components
225     chart_type = chart_config["chart_config"]["type"]
226     assert chart_type in ["bar", "line", "pie", "scatter", "doughnut"], f"Chart type should be valid, got {chart_type}"
227     assert "data" in chart_config["chart_config"], "Chart config should include data"
228     assert "labels" in chart_config["chart_config"]["data"], "Chart data should include labels"
229     assert "datasets" in chart_config["chart_config"]["data"], "Chart data should include datasets"

```

C.2 Mock AI Client Integration

```

19  class TestMockAIIClientIntegration:      You, 2 weeks ago * created intergration test for Mock AI Client In...
20      """Test cases for the MockAIclient in integration scenarios."""
21
22      Windsurf: Refactor | Explain | X
23      @pytest.fixture
24      def mock_client(self):
25          """Mock AI client fixture."""
26          return MockAIclient()
27
28      Windsurf: Refactor | Explain | X
29      def test_chart_recommendation(self, mock_client):
30          """Test that chart recommendations work as expected."""
31          # Create visualization data
32          visualization_data = {
33              'title': 'Student Enrollments 2020-21',
34              'columns': ['University', 'Undergraduate', 'Postgraduate', 'Total'],
35              'rows': [
36                  ['University A', '1000', '500', '1500'],
37                  ['University B', '2000', '1000', '3000'],
38                  ['University C', '3000', '1500', '4500']
39              ]
40          }
41
42          # Test the chart recommendation function
43          chart_recommendation = mock_client.get_chart_recommendation(visualization_data)
44
45          # Assert that the mock client returns the expected structure
46          assert "recommended_chart_type" in chart_recommendation, "Should include a recommended chart type"
47          assert "recommendation_reason" in chart_recommendation, "Should include a reason for the recommendation"
48          assert "example_prompts" in chart_recommendation, "Should include example prompts"
49      def test_visualization_generation(self, mock_client):
50          """Test visualization generation functionality."""
51          # Create visualization data
52          visualization_data = {
53              'title': 'Student Enrollments 2020-21',
54              'columns': ['University', 'Undergraduate', 'Postgraduate', 'Total'],
55              'rows': [
56                  ['University A', '1000', '500', '1500'],
57                  ['University B', '2000', '1000', '3000'],
58                  ['University C', '3000', '1500', '4500']
59              ]
60          }
61
62          # Generate a visualization config
63          chart_config = mock_client.generate_visualization_config(
64              visualization_data,
65              "Show student enrollments",
66              "bar"
67          )
68
69          # Verify the chart configuration structure
70          assert chart_config is not None, "Should generate a visualization configuration"
71          assert "chart_config" in chart_config, "Should include chart configuration"
72          assert "insights" in chart_config, "Should include insights"
73
74          # Basic test of chart configuration structure
75          assert "type" in chart_config["chart_config"], "Chart config should include type"
76          assert "data" in chart_config["chart_config"], "Chart config should include data"
77          assert "labels" in chart_config["chart_config"]["data"], "Chart data should include labels"

```

```
78     def test_chart_type_change(self, mock_client):
79         """Test changing chart type functionality."""
80         # Create visualization data
81         visualization_data = {
82             'title': 'Student Enrollments 2020-21',
83             'columns': ['University', 'Undergraduate', 'Postgraduate', 'Total'],
84             'rows': [
85                 ['University A', '1000', '500', '1500'],
86                 ['University B', '2000', '1000', '3000'],
87                 ['University C', '3000', '1500', '4500']
88             ]
89         }
90
91         # Generate a visualization with a different chart type
92         chart_config = mock_client.generate_visualization_config(
93             visualization_data,
94             "Show student enrollments",
95             "pie" # Use pie chart instead of default bar
96         )
97
98         # Verify the chart type was changed
99         assert chart_config["chart_config"]["type"] == "pie", "Chart type should be pie"
100        assert "data" in chart_config["chart_config"], "Chart config should include data"
```

D: Data Processing Pipeline Implementation

D.1 Data Ingestion Process

```
374     def process_single_file(self, file_name: str) -> bool:      You, 2 months ago • Using regular expression predefined query
375         """
376             Process a single file by name from the raw directory.
377         """
378         file_path = self.raw_dir / file_name
379
380         try:
381             logger.info(f"Processing single file: {file_name}")
382
383             # Validate file
384             if not self.validate_csv(file_path):
385                 return False
386
387             # Create new empty metadata
388             metadata = {}
389
390             # Read the first 20 lines to extract title and academic year
391             try:
392                 with open(file_path, 'r', encoding='utf-8', errors='ignore') as f:
393                     lines = [line.strip() for line in f.readlines()[:20]]
394
395                     # Extract title from the first line only
396                     first_line = lines[0] if lines else ""
397
398                     # Log the raw line for debugging
399                     logger.info(f"Raw first line: [{first_line}]")
400
401                     # Check if the line contains the title pattern
402                     if 'Title:' in first_line:
403                         # Get everything after 'Title:'
404                         title_text = first_line.split('Title:')[1].strip()
405                         logger.info(f"After split by 'Title': [{title_text}]")
406
407                         # Remove leading commas and quotes
408                         title_text = title_text.lstrip(',').strip('')
409                         logger.info(f"After cleaning commas/quotes: [{title_text}]")
410
411                         # Find the part after "Table XX - " if it exists
412                         table_match = re.search(r'Table\s+\d+\s*\-(.*\s*)', title_text)
413                         if table_match:
414                             # Extract just the part after the dash
415                             extracted_title = table_match.group(1).strip()
416                             # Remove any trailing quotes
417                             extracted_title = extracted_title.rstrip('')
418                             metadata['title'] = extracted_title
419                             logger.info(f"Extracted title: [{extracted_title}]")
420                         else:
421                             # If no table pattern, just use the whole text
422                             cleaned_title = title_text.strip()
423                             metadata['title'] = cleaned_title
424                             logger.info(f"No table pattern found, using cleaned title: [{cleaned_title}]")
425                         else:
426                             logger.warning(f"First line does not contain 'Title': [{first_line}]")
427
428                         # Extract academic year
429                         for line in lines:
430                             # Skip lines with "to" as they indicate a range of years (e.g., "Academic years 2020/21 to 2023/24")
431                             if "to" in line:
```

```

432         continue
433
434         # Look for patterns like ",Academic year,YYYY/YY" or "Filters:,Academic year,YYYY/YY"
435         year_match = re.search(r'(?:,|:)\s*Academic year\s*,\s*(20\d{2}/\d{2})', line)
436         if year_match:
437             metadata['academic_year'] = year_match.group(1).strip()
438             logger.info(f"Found academic year format: {metadata['academic_year']}")
439             break
440     except Exception as e:
441         logger.warning(f"Error extracting title/academic year: {str(e)}")
442
443     # If title not found, use "Unknown title"
444     if 'title' not in metadata:
445         logger.warning(f"Could not extract title from file content")
446         metadata['title'] = "Unknown title"
447
448     # If academic year not found, use "Unknown"
449     if 'academic_year' not in metadata:
450         logger.warning(f"Could not extract academic year from file content")
451
452     # Try to extract from filename before defaulting to Unknown
453     year_match = re.search(r'(\d{4})[.-](\d{2})', file_name)
454     if year_match:
455         metadata['academic_year'] = f"{year_match.group(1)}/{year_match.group(2)}"
456         logger.info(f"Extracted academic year from filename: {metadata['academic_year']}")
457     else:
458         metadata['academic_year'] = "Unknown"
459
460     # Clean data
461     cleaned_df = self.clean_csv(file_path)
462     if cleaned_df is None:
463         return False
464
465     # Extract title-based keywords
466     title_keywords = self.extract_keywords_from_title(metadata['title'])
467     metadata['keywords_title'] = title_keywords
468
469     # Extract column-based keywords
470     column_keywords = self.extract_keywords_from_columns(cleaned_df.columns)
471     metadata['keywords_columns'] = column_keywords
472
473     # Generate new filename based on title and academic year
474     original_filename = file_name
475     new_filename = self.generate_clean_filename(metadata['title'], metadata['academic_year'], original_filename)
476
477     # Store original filename in metadata
478     metadata['original_filename'] = original_filename
479     # Store new filename in metadata
480     metadata['new_filename'] = new_filename
481
482     # Create metadata line
483     metadata_line = f"#METADATA:{json.dumps(metadata)}\n"
484
485     # Save cleaned file with metadata using new filename
486     cleaned_path = self.clean_dir / new_filename
487
488     # Write metadata line first, then the data
489     with open(cleaned_path, 'w', encoding='utf-8') as f:
490         f.write(metadata_line)
491
492     # Then append the dataframe
493     cleaned_df.to_csv(cleaned_path, index=False, mode='a')
494
495     logger.info(f"Successfully processed {file_name} with metadata")
496     logger.info(f"Title: {metadata.get('title', 'Not found')}")
497     logger.info(f"Academic Year: {metadata.get('academic_year', 'Not found')}")
498     logger.info(f"New filename: {new_filename}")
499     logger.info(f"Title keywords: {title_keywords}")
500     logger.info(f"Column keywords: {column_keywords}")
501
502     return True
503
504 except Exception as e:
505     logger.error(f"Error processing {file_name}: {str(e)}")
506     return False

```

D.2 Manual Update Process

```
5110     @csrf_exempt
5111     def upload_data_view(request):
5112         """Render the upload data page."""
5113         return render(request, 'upload_data.html', {'show_form': True})
5114
5115     @csrf_exempt
5116     @require_http_methods(['POST'])
5117     def process_upload(request):
5118         """
5119
5120         # Define BASE_DIR
5121         BASE_DIR = Path(__file__).resolve().parent.parent
5122
5123         # Set up logging
5124         logger = logging.getLogger(__name__)
5125         logger.info("Processing file upload request")
5126
5127         # Prepare response data
5128         response_data = {
5129             'success': False,
5130             'error': None,
5131             'total_files': 0,
5132             'added_files_count': 0,
5133             'added_files': [],
5134             'error_files': [],
5135             'cleaning_status': 'Not started',
5136             'indexing_created': False
5137         }
5138
5139         try:
5140             # Get file processing option
5141             processing_option = request.POST.get('processing_option', 'add')
5142             logger.info(f"Processing option selected: {processing_option}")
5143
5144             # Check if files were uploaded
5145             if 'csv_files' not in request.FILES:
5146                 response_data['error'] = 'No files were uploaded'
5147                 return JsonResponse(response_data)
5148
5149             # Get uploaded files
5150             files = request.FILES.getlist('csv_files')
5151             logger.info(f"Received {len(files)} files for upload")
5152
5153             # Check if uploaded files exceed the maximum allowed number
5154             max_files = getattr(settings, 'DATA_UPLOAD_MAX_NUMBER_FILES', 20) # Default is 20
5155             if len(files) > max_files:
5156                 logger.error(f"Too many files: received {len(files)}, max allowed is {max_files}")
5157                 response_data['error'] = f"The number of files exceeded the maximum limit of {max_files}. Please upload"
5158                 return JsonResponse(response_data)
5159
5160             # Validate files are CSVs
5161             non_csv_files = [f.name for f in files if not f.name.lower().endswith('.csv')]
5162             if non_csv_files:
5163                 response_data['error'] = f"Only CSV files are allowed. Invalid files: {', '.join(non_csv_files)}"
5164                 return JsonResponse(response_data)
5165
5166             # Define paths
5167             raw_files_dir = RAW_FILES_DIR
5168             cleaned_files_dir = CLEANED_FILES_DIR
5169
5170             # Create directories if they don't exist
5171             raw_files_dir.mkdir(parents=True, exist_ok=True)
5172             cleaned_files_dir.mkdir(parents=True, exist_ok=True)
```

```

5193     if processing_option == 'overwrite':
5194         logger.info("Overwrite option selected - removing existing files")
5195         # Remove all files in raw_files directory
5196         for file_path in raw_files_dir.glob('*.*'):
5197             try:
5198                 file_path.unlink()
5199             except Exception as e:
5200                 logger.error(f"Error deleting raw file {file_path}: {str(e)}")
5201
5202         # Remove all files in cleaned_files directory
5203         for file_path in cleaned_files_dir.glob('*.*'):
5204             try:
5205                 file_path.unlink()
5206             except Exception as e:
5207                 logger.error(f"Error deleting cleaned file {file_path}: {str(e)}")
5208
5209         # Remove index file if it exists
5210         index_file_path = BASE_DIR / "hesa_files_index.json"
5211         if index_file_path.exists():
5212             try:
5213                 index_file_path.unlink()
5214                 logger.info("Removed existing index file")
5215             except Exception as e:
5216                 logger.error(f"Error deleting index file: {str(e)}")
5217         # For 'add' option, clean only cleaned files and index, but keep raw files
5218     elif processing_option == 'add':
5219         logger.info("Add files option selected - cleaning only processed files")
5220         # Remove all files in cleaned_files directory
5221         for file_path in cleaned_files_dir.glob('*.*'):
5222             try:
5223                 file_path.unlink()
5224                 logger.info(f"Removed cleaned file: {file_path.name}")
5225             except Exception as e:
5226                 logger.error(f"Error deleting cleaned file {file_path}: {str(e)}")
5227
5228         # Remove index file if it exists
5229         index_file_path = BASE_DIR / "hesa_files_index.json"
5230         if index_file_path.exists():
5231             try:
5232                 index_file_path.unlink()
5233                 logger.info("Removed existing index file")
5234             except Exception as e:
5235                 logger.error(f"Error deleting index file: {str(e)}")
5236
5237         # Process each uploaded file
5238         for file in files:
5239             file_name = file.name
5240             logger.info(f"Processing file: {file_name}")
5241
5242             # Check line count
5243             try:
5244                 # Read and count lines without loading entire file into memory
5245                 line_count = 0
5246                 for line in file.chunks():
5247                     line_count += line.count(b'\n')
5248
5249                 # Reset file pointer to beginning

```

```

5262     # Save file to raw_files directory
5263     try:
5264         destination_path = raw_files_dir / file_name
5265         with open(destination_path, 'wb+') as destination:
5266             for chunk in file.chunks():
5267                 destination.write(chunk)
5268
5269         response_data['added_files'].append(file_name)
5270         logger.info(f"Successfully saved {file_name} to raw files directory")
5271     except Exception as e:
5272         logger.error(f"Error saving {file_name}: {str(e)}")
5273         response_data['error_files'].append(f"{file_name} (error saving file)")
5274         continue
5275
5276     # Run CSV cleaning script if files were added
5277     if response_data['added_files']:
5278         response_data['added_files_count'] = len(response_data['added_files'])
5279
5280     try:
5281         # Initialize the CSV processor
5282         from core.data_processing.csv_processor import CSVProcessor
5283         processor = CSVProcessor()
5284
5285         # Process the files
5286         results = processor.process_all_files()
5287
5288         # Add the project root to the path if it's not already there
5289         project_root = str(BASE_DIR)
5290         if project_root not in sys.path:
5291             sys.path.insert(0, project_root)
5292
5293         # Import create_index_file
5294         from csv_cleaning import create_index_file
5295
5296         # Generate the index file
5297         index_created = create_index_file(processor, results)
5298         response_data['indexing_created'] = index_created
5299
5300         # Count total files in the system
5301         response_data['total_files'] = len(list(raw_files_dir.glob('*.*')))
5302
5303         # Set success flag
5304         response_data['success'] = True
5305
5306     except Exception as e:
5307         logger.error(f"Error during CSV cleaning process: {str(e)}")
5308         response_data['error'] = f"Error during processing: {str(e)}"
5309         response_data['cleaning_status'] = "Failed"
5310
5311     else:
5312         response_data['error'] = "No valid files were uploaded"
5313
5314 except Exception as e:
5315     logger.error(f"Unexpected error in file upload processing: {str(e)}")
5316     response_data['error'] = f"Unexpected error: {str(e)}"
5317
5318 return JsonResponse(response_data)

```

D.3 Metadata Extraction

```

75     def extract_metadata(self, file_path: Path) -> Dict[str, Any]:      You, 2 months ago • change the way data is being p
76         """Extract metadata from the first line of the file if it exists."""
77         try:
78             with open(file_path, 'r', encoding='utf-8', errors='ignore') as f:
79                 first_line = f.readline().strip()
80
81             if first_line.startswith('#METADATA:'):
82                 # Extract the JSON metadata
83                 metadata_json = first_line[9:] # Skip the #METADATA: prefix
84                 metadata = json.loads(metadata_json)
85                 return metadata
86             else:
87                 return {}
88         except Exception as e:
89             logger.warning(f"Error extracting metadata from {file_path}: {str(e)}")
90             return {}

```

D.4 Indexing

```
29 def create_index_file(processor, results):      You, 2 months ago + created index file for the AI approach. fix th...
30 """
31     Creates an index file at the root level containing metadata for all CSV files
32     in the cleaned_files directory.
33 """
34 logger = logging.getLogger(__name__)
35
36 # Define path for the index file
37 base_dir = Path(__file__).resolve().parent
38 index_file_path = base_dir / "hesa_files_index.json"
39
40 # Initialize index data structure
41 index_data = {
42     "hesa_files": [],
43     "updated_at": datetime.datetime.now().isoformat(),
44     "total_files": 0
45 }
46
47 # Process each cleaned file to extract metadata
48 for file_name in processor.get_available_datasets():
49     if file_name == '.gitkeep':
50         continue
51
52     cleaned_file_path = processor.clean_dir / file_name
53
54     try:
55         # Extract metadata from the cleaned file using the processor's method
56         metadata = processor.extract_metadata(cleaned_file_path)
57
58         if not metadata:
59             logger.warning(f"No metadata found using processor for {file_name}, trying direct extraction")
60             try:
61                 # Read the first line to extract metadata
62                 with open(cleaned_file_path, 'r', encoding='utf-8', errors='ignore') as f:
63                     first_line = f.readline().strip()
64
65                     if first_line.startswith('#METADATA:'):
66                         # Extract the JSON metadata
67                         metadata_json = first_line[9:] # Skip the #METADATA: prefix
68                         metadata = json.loads(metadata_json)
69                         logger.info(f"Successfully extracted metadata directly from file: {file_name}")
70                     else:
71                         logger.warning(f"File {file_name} doesn't have metadata line")
72                         metadata = {}
73             except Exception as e:
74                 logger.warning(f"Error extracting metadata from first line of {file_name}: {str(e)}")
75             metadata = {}
76
77         # If still no metadata or missing key fields, try to create minimal metadata
78         if not metadata or 'title' not in metadata or 'academic_year' not in metadata:
79             logger.warning(f"Incomplete metadata for {file_name}, creating minimal metadata")
80
81             # If title is missing, extract it from the raw file
82             if 'title' not in metadata:
83                 # Try to find the corresponding raw file
84                 raw_file_path = processor.raw_dir / file_name
```

```

87     if raw_file_path.exists():
88         try:
89             # Extract title using the same logic as in processor.process_all_files
90             with open(raw_file_path, 'r', encoding='utf-8', errors='ignore') as f:
91                 lines = [line.strip() for line in f.readlines()[:20]]
92
93             # Extract title from the first line
94             first_line = lines[0] if lines else ""
95
96             if 'Title:' in first_line:
97                 # Get everything after 'Title:'
98                 title_text = first_line.split('Title:')[1].strip()
99
100            # Remove leading commas and quotes
101            title_text = title_text.lstrip(',').strip('\"')
102
103            # Find the part after "Table XX - " if it exists
104            import re
105            table_match = re.search(r'Table\s+\d+\s*-.*', title_text)
106            if table_match:
107                # Extract just the part after the dash
108                extracted_title = table_match.group(1).strip()
109                # Remove any trailing quotes
110                extracted_title = extracted_title.rstrip('\"')
111                metadata['title'] = extracted_title
112            else:
113                # If no table pattern, just use the whole text
114                cleaned_title = title_text.strip()
115                metadata['title'] = cleaned_title
116
117                logger.info(f"Extracted title from raw file: {metadata['title']}"))
118            else:
119                # Fallback to cleaned filename
120                metadata['title'] = file_name.replace('.csv', '').split(' 20')[0]
121                logger.warning(f"Using fallback title from filename: {metadata['title']}"))
122            except Exception as e:
123                logger.warning(f"Error extracting title from raw file: {str(e)}")
124                # Fallback to cleaned filename
125                metadata['title'] = file_name.replace('.csv', '').split(' 20')[0]
126            else:
127                # Fallback to cleaned filename
128                metadata['title'] = file_name.replace('.csv', '').split(' 20')[0]
129                logger.warning(f"Raw file not found, using fallback title from filename: {metadata['title']}"))
130
131            # If academic year is missing, extract it from the raw file first, then try filename
132            if 'academic_year' not in metadata:
133                raw_file_path = processor.raw_dir / file_name
134
135            if raw_file_path.exists():
136                try:
137                    # Read the first 20 lines to look for academic year
138                    with open(raw_file_path, 'r', encoding='utf-8', errors='ignore') as f:
139                        lines = [line.strip() for line in f.readlines()[:20]]
140
141                    # Look for the exact pattern with comma at the start: ",Academic year,YYYY/YY"
142                    academic_year_found = False
143                    for line in lines:
144                        # Skip lines with "to" as they indicate a range of years

```



```
203
204     # Add reference (file name)
205     metadata['reference'] = file_name
206
207     # Add metadata to index
208     index_data["hesa_files"].append(metadata)
209     logger.info(f"Added metadata for {file_name} to index with title: '{metadata.get('title', 'No title')}'")
210     logger.info(f"Academic Year: {metadata.get('academic_year', 'Unknown')}")  

211
212 except Exception as e:
213     logger.error(f"Error processing {file_name} for index: {str(e)}")
214
215 # Update total files count
216 index_data["total_files"] = len(index_data["hesa_files"])
217
218 # Write index file
219 try:
220     with open(index_file_path, 'w', encoding='utf-8') as f:
221         json.dump(index_data, f, indent=2)
222     logger.info(f"Successfully created index file at {index_file_path}")
223     logger.info(f"Total files indexed: {index_data['total_files']}")
224     return True
225 except Exception as e:
226     logger.error(f"Error writing index file: {str(e)}")
227     return False
```

D.5 Cleaning Process

```
374     def process_single_file(self, file_name: str) -> bool:      You, 2 months ago * using regular expression predefined query i
375         """
376         Process a single file by name from the raw directory.
377         """
378         file_path = self.raw_dir / file_name
379
380         try:
381             logger.info(f"Processing single file: {file_name}")
382
383             # Validate file
384             if not self.validate_csv(file_path):
385                 return False
386
387             # Create new empty metadata
388             metadata = {}
389
390             # Read the first 20 lines to extract title and academic year
391             try:
392                 with open(file_path, 'r', encoding='utf-8', errors='ignore') as f:
393                     lines = [line.strip() for line in f.readlines()[:20]]
394
395                 # Extract title from the first line only
396                 first_line = lines[0] if lines else ""
397
398                 # Log the raw line for debugging
399                 logger.info(f"Raw first line: [{first_line}]")
400
401                 # Check if the line contains the title pattern
402                 if 'Title:' in first_line:
403                     # Get everything after 'Title:'
404                     title_text = first_line.split('Title:')[1].strip()
405                     logger.info(f"After split by 'Title': [{title_text}]")
406
407                     # Remove leading commas and quotes
408                     title_text = title_text.lstrip(',').strip('"')
409                     logger.info(f"After cleaning commas/quotes: [{title_text}]")
410
411                     # Find the part after "Table XX - " if it exists
412                     table_match = re.search(r'Table\s+\d+\s*-.*', title_text)
413                     if table_match:
414                         # Extract just the part after the dash
415                         extracted_title = table_match.group(1).strip()
416                         # Remove any trailing quotes
417                         extracted_title = extracted_title.rstrip('"')
418                         metadata['title'] = extracted_title
419                         logger.info(f"Extracted title: [{extracted_title}]")
420                     else:
421                         # If no table pattern, just use the whole text
422                         cleaned_title = title_text.strip()
423                         metadata['title'] = cleaned_title
424                         logger.info(f"No table pattern found, using cleaned title: [{cleaned_title}]")
425                 else:
426                     logger.warning(f"First line does not contain 'Title': [{first_line}]")
427
428                 # Extract academic year
429                 for line in lines:
430                     # Skip lines with "to" as they indicate a range of years (e.g., "Academic years 2020/21 to 2023/24")
```

```

431     if "to" in line:
432         continue
433
434         # Look for patterns like ",Academic year,YYYY/YY" or "Filters:,Academic year,YYYY/YY"
435         year_match = re.search(r'(? :,|:)\s*Academic year\s*,\s*(20\d{2})/\d{2})', line)
436         if year_match:
437             metadata['academic_year'] = year_match.group(1).strip()
438             logger.info(f"Found academic year format: {metadata['academic_year']}")
439             break
440     except Exception as e:
441         logger.warning(f"Error extracting title/academic year: {str(e)}")
442
443         # If title not found, use "Unknown title"
444         if 'title' not in metadata:
445             logger.warning(f"Could not extract title from file content")
446             metadata['title'] = "Unknown title"
447
448         # If academic year not found, use "Unknown"
449         if 'academic_year' not in metadata:
450             logger.warning(f"Could not extract academic year from file content")
451
452         # Try to extract from filename before defaulting to Unknown
453         year_match = re.search(r'(\d{4})[.-_](\d{2})', file_name)
454         if year_match:
455             metadata['academic_year'] = f'{year_match.group(1)}/{year_match.group(2)}'
456             logger.info(f"Extracted academic year from filename: {metadata['academic_year']}")
457         else:
458             metadata['academic_year'] = "Unknown"
459
460         # Clean data
461         cleaned_df = self.clean_csv(file_path)
462         if cleaned_df is None:
463             return False
464
465         # Extract title-based keywords
466         title_keywords = self.extract_keywords_from_title(metadata['title'])
467         metadata['keywords_title'] = title_keywords
468
469         # Extract column-based keywords
470         column_keywords = self.extract_keywords_from_columns(cleaned_df.columns)
471         metadata['keywords_columns'] = column_keywords
472
473         # Generate new filename based on title and academic year
474         original_filename = file_name
475         new_filename = self.generate_clean_filename(metadata['title'], metadata['academic_year'], original_filename)
476
477         # Store original filename in metadata
478         metadata['original_filename'] = original_filename
479         # Store new filename in metadata
480         metadata['new_filename'] = new_filename
481
482         # Create metadata line
483         metadata_line = f"#METADATA:{json.dumps(metadata)}\n"
484
485         # Save cleaned file with metadata using new filename
486         cleaned_path = self.clean_dir / new_filename
487
488         # Write metadata line first, then the data

```

```
489     with open(cleaned_path, 'w', encoding='utf-8') as f:
490         f.write(metadata_line)
491
492     # Then append the dataframe
493     cleaned_df.to_csv(cleaned_path, index=False, mode='a')
494
495     logger.info(f"Successfully processed {file_name} with metadata")
496     logger.info(f"Title: {metadata.get('title', 'Not found')}")
497     logger.info(f"Academic Year: {metadata.get('academic_year', 'Not found')} ")
498     logger.info(f"New filename: {new_filename}")
499     logger.info(f"Title keywords: {title_keywords}")
500     logger.info(f"Column keywords: {column_keywords}")
501
502     return True
503
504 except Exception as e:
505     logger.error(f"Error processing {file_name}: {str(e)}")
506     return False
```

D.6 CSVProcessor

```
16 class CSVProcessor:
17     Windsurf: Refactor | Explain | Generate Docstring | X
18     def __init__(self):
19         self.raw_dir = RAW_FILES_DIR
20         self.clean_dir = CLEANED_FILES_DIR
21
22         Windsurf: Refactor | Explain | X
23         def validate_csv(self, file_path: Path) -> bool:      You, 3 months ago • changed tailwind CND to direct instalation and
24             """
25             Validates if the CSV file is properly formatted and contains expected data.
26             """
27             try:
28                 df = pd.read_csv(file_path)
29
29                 # Check if file is empty
30                 if df.empty:
31                     logger.error(f"File {file_path} is empty")
32                     return False
33
34                 # Check for minimum required columns (customize based on your HESA data structure)
35                 required_columns = ['institution_id', 'year'] # Add your required columns
36                 missing_columns = [col for col in required_columns if col not in df.columns]
37
38                 if missing_columns:
39                     logger.error(f"Missing required columns: {missing_columns}")
40                     return False
41
42                 return True
43
44             except Exception as e:
45                 logger.error(f"Error validating {file_path}: {str(e)}")
46                 return False
47
48         Windsurf: Refactor | Explain | X
49         def clean_csv(self, file_path: Path) -> Optional[pd.DataFrame]:
50             """
51             Cleans the CSV data by handling missing values, removing duplicates, etc.
52             """
53             try:
54                 df = pd.read_csv(file_path)
55
56                 # Remove duplicate rows
57                 df = df.drop_duplicates()
58
59                 # Handle missing values (customize based on your needs)
60                 df = df.fillna({
61                     'numeric_columns': 0, # Fill numeric missing values with 0
62                     'string_columns': 'Unknown' # Fill string missing values with 'Unknown'
63                 })
64
65                 # Convert date columns to datetime (if any)
66                 date_columns = [] # Add your date column names
67                 for col in date_columns:
68                     df[col] = pd.to_datetime(df[col], errors='coerce')
69
70                 return df
71
72             except Exception as e:
```

```

71     logger.error(f"Error cleaning {file_path}: {str(e)}")
72     return None
73
74     Windsurf: Refactor | Explain | X
75     def process_all_files(self) -> Dict[str, bool]:
76         """
77             Process all CSV files in the raw directory and save cleaned versions.
78         """
79         results = {}
80
81         # Create cleaned directory if it doesn't exist
82         self.clean_dir.mkdir(parents=True, exist_ok=True)
83
84         for file_path in self.raw_dir.glob('*.*csv'):
85             try:
86                 # Validate file
87                 if not self.validate_csv(file_path):
88                     results[file_path.name] = False
89                     continue
90
91                 # Clean data
92                 cleaned_df = self.clean_csv(file_path)
93                 if cleaned_df is None:
94                     results[file_path.name] = False
95                     continue
96
97                 # Save cleaned file
98                 cleaned_path = self.clean_dir / file_path.name
99                 cleaned_df.to_csv(cleaned_path, index=False)
100                results[file_path.name] = True
101
102                logger.info(f"Successfully processed {file_path.name}")
103
104            except Exception as e:
105                logger.error(f"Error processing {file_path.name}: {str(e)}")
106                results[file_path.name] = False
107
108        return results
109
110     Windsurf: Refactor | Explain | X
111     def get_available_datasets(self) -> List[str]:
112         """
113             Returns a list of available cleaned dataset names.
114         """
115         return [f.name for f in self.clean_dir.glob('*.*csv')]

```

E: Query Processing and Interpretation

E.1 Transform free text into query parameters

```
36     def analyze_query(self, query: str) -> Dict[str, Any]:      You, 2 months ago • created indeing file for the AI approach.
37         """
38             Analyze a natural language query using the Gemini API.
39
40         """
41         import requests
42         import json
43         import re
44         import logging
45
46         logger = logging.getLogger(__name__)
47         logger.info(f"Analyzing query with Gemini API: {query}")
48
49         if not self.api_key:
50             logger.error("No Gemini API key provided")
51             raise ValueError("Gemini API key is required")
52
53         # Construct the prompt
54         system_prompt = """
55             You are a helpful assistant that extracts structured information from a user's query about higher education statistics
56
57             Extract the following information and output it in JSON format:
58             1. Institutions mentioned (e.g., "University of Oxford", "University of Leicester")
59             2. Years mentioned (e.g., "2020", "2020/21")
60             3. If a year range is given, identify the start_year and end_year
61             4. The type of data being requested (e.g., "student numbers", "enrollment", "graduates")
62
63             IMPORTANT: You need to detect and correct ONLY SPELLING typos in institution names and years.
64             For institutions, follow these strict rules:
65             - ONLY consider a term to have a typo if it has clear misspellings (e.g., "Univercity" → "University", "Lieicester" → "Leicester")
66             - DO NOT consider "Oxford University" vs "University of Oxford" as a typo - these are different ways to refer to the
67                 same institution
68             - DO NOT mark an institution as having a typo if it's spelled correctly but uses a different naming convention
69             - DO NOT transform a city name like "london" into "The University of London"
70             - PRESERVE the original terms from the query - if user says "london", keep it as "london" in the institutions list
71             - If there are no spelling errors, original_institutions should exactly match institutions, and has_institution_typos
72                 should be false
73
74             For years:
75             - ONLY mark years as having typos if they have clear numerical errors (e.g., "20025" → "2025")
76             - If no year typos exist, original_years should exactly match years, and has_year_typos should be false
77
78             Be careful about interpreting years in academic context:
79             - If the query contains phrases like "starting in [YEAR]" or "beginning in [YEAR]", interpret [YEAR] as the start of
80                 an academic year.
81             - If the query contains phrases like "end of [YEAR]" or "ending in [YEAR]", interpret [YEAR] as the end of an
82                 academic year.
83             - For ranges like "2016 to 2017", treat both as starting years of academic years.
84             - If there's no clarification, assume a year refers to the starting year of an academic year.
85             - For "past X years", calculate the years based on the current year (2025).
86
87             For the data_request field, use the SPECIFIC TERMS from the query whenever possible:
88             - If the query mentions "undergraduates", use ["undergraduate"] not ["student_enrollment"]
89             - If the query mentions "postgraduates", use ["postgraduate"] not ["student_enrollment"]
90             - If the query mentions "student numbers", use ["student_count"]
```

```

88     - If the query mentions "enrollment" or "enrolment", use ["enrollment"]
89     - If the query mentions "staff" or "teachers", use ["staff_data"]
90     - If the query mentions "research", use ["research_data"]
91     - NEVER generalize specific educational levels - keep them exactly as requested
92
93     Output format:
94     {
95         "institutions": ["University X", "london"],
96         "original_institutions": ["University X", "london"],
97         "has_institution_typos": true,
98         "years": ["2019/20", "2020/21"],
99         "original_years": ["20019", "2020"],
100        "has_year_typos": true,
101        "start_year": "2019",
102        "end_year": "2020",
103        "data_request": ["undergraduate", "graduation_rates"]
104    }
105
106    If no specific institutions are mentioned, return an empty list for institutions and original_institutions.
107    If no typos were detected, original_institutions should match institutions, and has_institution_typos should be false.
108    The same applies to years and original_years.
109    """
110
111    # The actual user query
112    user_prompt = f"Extract information from this query, paying special attention to academic year conventions and
113    potential typos in institution names and years: '{query}'"
114    ^
115
116    try:
117        # Gemini API endpoint - Updated to use the most current endpoint
118        url = "https://generativelanguage.googleapis.com/v1/models/gemini-1.5-pro:generateContent"
119
120        # Request payload
121        payload = {
122            "contents": [
123                {
124                    "role": "user",
125                    "parts": [
126                        {"text": system_prompt},
127                        {"text": user_prompt}
128                    ]
129                },
130                "generationConfig": {
131                    "temperature": 0.1,
132                    "topP": 0.95,
133                    "topK": 40,
134                    "maxOutputTokens": 2048
135                }
136            }
137
138            # Headers with API key
139            headers = {
140                "Content-Type": "application/json",
141                "x-goog-api-key": self.api_key
142            }

```

```
143     # Make the request
144     response = requests.post(url, json=payload, headers=headers)
145
146     # Check for successful response
147     if response.status_code != 200:
148         logger.error(f"Gemini API error: {response.status_code} - {response.text}")
149         raise Exception(f"Gemini API returned status code {response.status_code}: {response.text}")
150
151     # Parse the response
152     response_data = response.json()
153
154     # Extract text from response
155     if 'candidates' in response_data and len(response_data['candidates']) > 0:
156         if 'content' in response_data['candidates'][0]:
157             response_text = response_data['candidates'][0]['content']['parts'][0]['text']
158         else:
159             logger.error("No content in Gemini API response")
160             raise Exception("No content in Gemini API response")
161     else:
162         logger.error("No candidates in Gemini API response")
163         raise Exception("No candidates in Gemini API response")
164
165     # Find the JSON response within the text
166     json_match = re.search(r'\{.*\}', response_text, re.DOTALL)
167     if not json_match:
168         logger.error(f"Failed to extract JSON from response: {response_text}")
169         raise Exception("Failed to extract JSON from the Gemini API response")
170
```

```

170
171     # Parse the JSON data
172     result = json.loads(json_match.group(0))
173
174     # Default values for missing fields
175     result.setdefault('institutions', [])
176     result.setdefault('original_institutions', [])
177     result.setdefault('has_institution_typos', False)
178     result.setdefault('years', [])
179     result.setdefault('original_years', [])
180     result.setdefault('has_year_typos', False)
181     result.setdefault('data_request', ['general_data'])
182
183     # Always ensure University of Leicester is included if any institution is mentioned
184     if result['institutions'] and 'University of Leicester' not in result['institutions']:
185         result['institutions'].append('University of Leicester')
186
187     # Fix "in University" issue by filtering out partial names
188     result['institutions'] = [
189         inst for inst in result['institutions']
190         if len(inst.split()) > 1 # Must be at least two words
191         and inst.lower() != "in university" # Explicitly exclude "in university"
192     ]
193
194     # If we don't have typo information but have institutions, assume no typos
195     if 'original_institutions' not in result or not result['original_institutions']:
196         result['original_institutions'] = result['institutions'].copy()
197         result['has_institution_typos'] = False
198     if 'original_years' not in result or not result['original_years']:
199         result['original_years'] = result['years'].copy()
200         result['has_year_typos'] = False
201
202     # Handle special case for "past X years"
203     if 'start_year' not in result or 'end_year' not in result:
204         # Check if the query contains "past X years"
205         past_years_match = re.search(r'past\s+(\d+)\s+years?', query.lower())
206         if past_years_match:
207             num_years = int(past_years_match.group(1))
208             # Use current year from datetime
209             current_year = datetime.now().year
210             end_year = current_year
211             start_year = current_year - num_years
212
213             result['end_year'] = str(end_year)
214             result['start_year'] = str(start_year)
215
216             # Expand years array to include all years in the range
217             if not result['years']:
218                 for year in range(start_year, end_year + 1):
219                     # Only add academic year format, not plain year
220                     result['years'].append(f"{year}/{str(year+1)[2:4]}")
221
222             # Apply our academic year logic for special cases if Gemini didn't handle it well
223             self._process_academic_year_logic(query, result)
224
225             logger.info(f"Extracted information: {result}")
226
227     return result
228

```

E.2 Date Interpretation

```
234     def _process_academic_year_logic(self, query, result):      You, 2 months ago + changed the way dataset is created. it now
235         """
236             Process academic year logic based on context clues in the query.
237             This handles cases like:
238             - "starting in 2017" → 2017/18
239             - "end of 2017" → 2016/17
240             - "2016 to 2017" → 2016/17 - 2017/18 (both as starting years)
241
242         """
243         import re
244         import logging
245
246         logger = logging.getLogger(__name__)
247         query_lower = query.lower()
248
249         # Extract years with context
250         start_patterns = [
251             r'start(?:ing|s|ed)?\s+(?:in|from|at)?\s+(?:the\s+)?(?:year\s+)?(\d{4})',
252             r'begin(?:ning|s)?\s+(?:in|from|at)?\s+(?:the\s+)?(?:year\s+)?(\d{4})',
253             r'from\s+(?:the\s+)?(?:year\s+)?(\d{4})'
254         ]
255
256         end_patterns = [
257             r'end(?:ing|s|ed)?\s+(?:in|at|of)?\s+(?:the\s+)?(?:year\s+)?(\d{4})',
258             r'finish(?:ing|es|ed)?\s+(?:in|at)?\s+(?:the\s+)?(?:year\s+)?(\d{4})',
259             r'(?:(in|at)\s+(?:the\s+)?end\s+of\s+(?:the\s+)?(?:year\s+)?(\d{4})'
260         ]
261
262         # Handle range logic for years without context
263         if 'start_year' in result and 'end_year' in result and result['start_year'] is not None and result['end_year'] is not
264             None:
265                 # Range years should be treated as starting years of academic years
266                 start_year = result['start_year']
267                 end_year = result['end_year']
268
269                 try:
270                     # Clear existing years array to rebuild it with correct academic years
271                     result['years'] = []
272
273                     # Add academic years for the range
274                     for year in range(int(start_year), int(end_year) + 1):
275                         academic_year = f"{year}/{str(year+1)[2:4]}"
276                         result['years'].append(academic_year)
277
277                     logger.info(f"Processed year range: {start_year}-{end_year} as academic years: {', '.join(result['years'])}")
278                 except (ValueError, TypeError) as e:
279                     logger.error(f"Error processing year range: {e}. start_year={start_year}, end_year={end_year}")
280
281             return
282
283         # Explicitly defined start years
284         for pattern in start_patterns:
285             start_match = re.search(pattern, query_lower)
286             if start_match:
287                 year = start_match.group(1)
288                 try:
289                     academic_year = f"{year}/{str(int(year)+1)[2:4]}"
```

```

290             if academic_year not in result['years']:
291                 result['years'].append(academic_year)
292                 result['start_year'] = year
293                 logger.info(f"Processed starting year {year} as academic year {academic_year}")
294             except (ValueError, TypeError) as e:
295                 logger.error(f"Error processing starting year {year}: {e}")
296
297             # Explicitly defined end years
298             for pattern in end_patterns:
299                 end_match = re.search(pattern, query_lower)
300                 if end_match:
301                     year = end_match.group(1)
302                     try:
303                         previous_year = str(int(year) - 1)
304                         academic_year = f"{previous_year}/{year[2:4]}"
305
306                         if academic_year not in result['years']:
307                             result['years'].append(academic_year)
308                             result['end_year'] = year
309                             logger.info(f"Processed ending year {year} as academic year {academic_year}")
310                         except (ValueError, TypeError) as e:
311                             logger.error(f"Error processing ending year {year}: {e}")
312
313             # For single years without context, treat as starting years
314             plain_year_pattern = r'\b(\d{2})\b'
315             years_mentioned = re.findall(plain_year_pattern, query)
316
317             # Skip years that were already processed with context
318             already_processed = True
319             break
320
321         if not already_processed:
322             try:
323                 # No context, treat as starting year
324                 academic_year = f"{year}/{str(int(year)+1)[2:4]}"
325                 if academic_year not in result['years']:
326                     result['years'].append(academic_year)
327                     # If no explicit start_year is set, use this as default
328                     if 'start_year' not in result or result['start_year'] is None:
329                         result['start_year'] = year
330                     logger.info(f"Processed plain year {year} as academic year {academic_year}")
331
332                     # Remove the plain year from the years list if it exists
333                     if year in result['years']:
334                         result['years'].remove(year)
335                         logger.info(f"Removed plain year {year} from years list to avoid duplication")
336             except (ValueError, TypeError) as e:
337                 logger.error(f"Error processing plain year {year}: {e}")
338
339             # Make sure years are unique
340             result['years'] = list(set(result['years']))
341
342             # Filter out any plain years that might still be in the list from other sources
343             result['years'] = [year for year in result['years'] if '/' in year]
344

```

F: Visualization Generation

F.1 Chart Generation

```
3848 def generate_visualization(client, dataset_info, user_request, chart_type=None):
3849     """
3850     Generate a Chart.js visualization based on the dataset and user request
3851     """
3852     try:
3853         # Extract relevant information from dataset_info
3854         title = dataset_info.get('title', '')
3855         columns = dataset_info.get('columns', [])
3856         rows = dataset_info.get('rows', [])
3857         query = dataset_info.get('query', '')
3858         institutions = dataset_info.get('institutions', [])
3859         years = dataset_info.get('years', [])
3860         all_year_data = dataset_info.get('allYearData', [])
3861
3862         # Log information about the request
3863         logging.info(f"Visualization request for dataset: {title}")
3864         logging.info(f"User request: {user_request}")
3865         logging.info(f"Years detected: {years}")
3866         logging.info(f"Number of rows: {len(rows)}")
3867         logging.info(f"Number of all year data rows: {len(all_year_data) if all_year_data else 0}")
3868         logging.info(f"Requested chart type: {chart_type}")
3869
3870         # Extract institution names directly from the dataset
3871         dataset_institutions = []
3872         provider_col_idx = -1
3873
3874         if chart_type == 'line' and not dataset_characteristics["has_multiple_years"]:
3875             chart_compatibility_check = """
3876             IMPORTANT NOTE: The user has requested a line chart, but this dataset only contains data for a single academic
3877             year.
3878             Line charts are generally better suited for showing trends over multiple time periods.
3879
3880             You should:
3881             1. Try to generate a line chart as requested, adapting the data to work with a single year if possible
3882             2. Make sure the chart configuration is valid JavaScript that can be evaluated
3883             3. Explicitly mention the limitation of using a line chart with single-year data in your insights
3884             """
3885             chart_compatibility_issue = True
3886
3887         elif chart_type == 'pie' and len([col for col in columns if col not in ['HE provider', 'Academic Year']]) > 5:
3888             chart_compatibility_check = """
3889             IMPORTANT NOTE: The user has requested a pie chart with a dataset that has many data columns.
3890             Pie charts work best with a small number of categories (typically 5 or fewer).
3891
3892             You should:
3893             1. Generate a pie chart as requested, but focus on the most important categories
3894             2. Make sure the chart configuration is valid JavaScript that can be evaluated
3895             3. Mention the limitation of using pie charts with many categories in your insights
3896             """
3897             chart_compatibility_issue = True
3898
3899     
```

```

3943     if len(institutions) > 5 and 'compare' in user_request.lower() and any(inst.lower() in user_request.lower() for inst
3944         in institutions):
3945             too_many_institutions_warning = """
3946             IMPORTANT NOTE: The user has requested comparing multiple institutions, but displaying too many institutions
3947             on a single chart can make it cluttered and hard to read.
3948
3949             You should:
3950             1. Limit the visualization to the 3-5 most relevant institutions mentioned in the user's request
3951             2. If the user explicitly mentioned specific institutions, prioritize those
3952             3. Mention in the insights that you've limited the chart to the most relevant institutions for readability
3953             """
3954             chart_compatibility_check += too_many_institutions_warning
3955             too_many_institutions = True
3956             chart_compatibility_issue = True
3957
3958             # Ensure 'The University of Leicester' is always considered
3959             leicester_variants = ['university of leicester', 'the university of leicester', 'leicester university']
3960             leicester_included = False
3961             for inst in institutions:
3962                 if inst.lower() in leicester_variants:
3963                     leicester_included = True
3964                     break
3965
3966             if not leicester_included and 'leicester' not in '.join(institutions).lower():
3967                 institutions.append('The University of Leicester')
3968                 logging.info("Added University of Leicester to institutions list")
3969
3970             institution_analysis = ""
3971             if len(mentioned_institutions) > 0 and len(referenced_institutions) == 0:
3972                 # User mentioned institutions but none match our dataset
3973                 institution_analysis = f"""
3974                 IMPORTANT WARNING: The user has requested information about institutions that don't exist in the dataset.
3975                 The available institutions in this dataset are: {', '.join(institutions[:10])}{'...' if len(institutions) > 10
3976                 else ''}
3977
3978                 You MUST adapt the visualization to use only institutions that exist in the dataset and clearly explain in the
3979                 insights
3980                 that the requested institution(s) aren't available in this dataset.
3981                 """
3982
3983             elif any(university_term in user_request.lower() for university_term in ["university", "college", "school"]) and len
3984                 (referenced_institutions) == 0:
3985                 # Generic institution terms but no specific matches
3986                 institution_analysis = """
3987                     IMPORTANT WARNING: The user's request mentions educational institutions, but I couldn't determine which specific
3988                     institutions from the dataset they're interested in.
3989
3990                     Please create a visualization using the institutions available in the dataset, and clearly explain in the
3991                     insights
3992                     which institutions were selected and why.
3993                     """
3994
3995
3996
3997
3998
3999
4000
4001
4002
4003
4004

```

```

4056     axis_instructions = ""
4057     if "axis" in user_request.lower() or "axes" in user_request.lower():
4058         axis_instructions = """
4059             I've noticed the user has specified some axis preferences. Please follow these instructions carefully:
4060             1. If the user wants to swap X and Y axes, make sure to implement this exactly as requested
4061             2. If the user specifies what should be on X or Y axis, honor this request precisely
4062             3. For pie charts, note that traditional X/Y axes don't apply, but use the specified dimensions for data
4063                 selection
4064             """
4065
4066     # Add any compatibility warnings for the selected chart type
4067     chart_compatibility_note = ""
4068     if chart_type:
4069         chart_compatibility_warning = get_chart_type_explanation(chart_type, dataset_characteristics)
4070         if chart_compatibility_warning:
4071             chart_compatibility_note = f"\n\nIMPORTANT NOTE ABOUT {chart_type.upper()} CHARTS:
4072                 {chart_compatibility_warning}"
4073
4074     # Set instructions for requested chart type
4075     chart_type_instruction = ""
4076     if chart_type:
4077         chart_type_instruction = f"""
4078             IMPORTANT: You MUST generate a {chart_type} chart for this visualization. Do not suggest or use any other chart
4079                 type.
4080             The user has specifically requested a {chart_type} chart, so optimize the visualization for this chart type.
4081                 {chart_compatibility_check}
4082
4083             chart_type_instruction = ""
4084             if chart_type:
4085                 chart_type_instruction = f"""
4086                     IMPORTANT: You MUST generate a {chart_type} chart for this visualization. Do not suggest or use any other chart
4087                         type.
4088                     The user has specifically requested a {chart_type} chart, so optimize the visualization for this chart type.
4089                         {chart_compatibility_check}
4090
4091             For a {chart_type} chart specifically:
4092                 - If it's a 'pie' chart: Ensure the Chart.js configuration has type: 'pie' and data formatted for a pie chart
4093                 - If it's a 'line' chart: Ensure the Chart.js configuration has type: 'line' and datasets with x/y coordinates
4094                 - If it's a 'bar' chart: Ensure the Chart.js configuration has type: 'bar' with labels and dataset values
4095
4096             CRITICAL: The JSON must be valid JavaScript that can be evaluated directly.
4097             DO NOT use string literals with quotes, Object.assign, or template literals in the chart configuration.
4098             """
4099
4100     # Create years description
4101     years_description = f"{len(years)} academic year{'s' if len(years) > 1 else ''}"

```

```
4119 |     IMPORTANT INSTRUCTIONS:  
4120 |     1. Generate a Chart.js configuration for a {chart_type if chart_type else "suitable"} chart visualization.  
4121 |     2. STRICTLY USE ONLY the {chart_type.upper() if chart_type else "appropriate"} CHART TYPE for visualization, even if  
4122 |     another type might be better suited.  
4123 |     3. Ensure the visualization accurately represents the data and addresses the user's request.  
4124 |     4. ONLY use the exact institution names as they appear in the dataset.  
4125 |     5. If the user mentions institutions not in the dataset, adapt the visualization using only available institutions  
4126 |     and clearly explain this in the insights.  
4127 |     6. Generate clear and insightful analysis of the visualization.  
4128 |     7. If the user's request cannot be fulfilled with the available data, provide a clear explanation.  
4129 |     8. Make sure colors are distinct and the visualization is easy to interpret.  
4130 |     9. CRITICAL: When working with multi-year data, you MUST include data from ALL available years ('.join(years))  
4131 |     in your visualization and analysis.  
4132 |     10. Do NOT state that data for certain years is missing unless you've verified it's actually missing from the  
4133 |         dataset.  
4134 |     11. When a user requests information about "London College of Business Sciences" and "Empire College London  
4135 |         Limited", ensure you include data for BOTH 2015/16 and 2016/17.  
4136 |     12. For visualization requests involving specific institutions across years, show the values for EVERY year in the  
4137 |         dataset.  
4138 |  
4139 |     CRITICAL FOR TIME-BASED VISUALIZATIONS:  
4140 |     1. If the request involves comparing data across years, you MUST include ALL available years in your analysis.  
4141 |     2. For requests involving trends or changes over time, be sure to show the complete timeline from {years[0]} to  
4142 |     {years[-1] if len(years) > 1 else years[0]}.  
4143 |     3. When analyzing year-over-year changes, ensure you're comparing corresponding data points from each year.  
4144 |     4. For institutions across multiple years, you MUST include a data point for each year, clearly showing the  
4145 |         comparison.  
4146 |     5. The "Total" column values should be compared between years, showing how enrollment changed from one year to the  
4147 |         next.
```

F.2 Handles different chart type requests

```

332 def get_chart_data(request, chart_type):      You, 2 months ago * created charts of type pie, line and graph for ...
333     try:
334         # Path to the CSV file
335         csv_path = os.path.join(settings.BASE_DIR, 'data', 'cleaned_files', 'chart-1_cleaned.csv')
336
337         # Read the CSV file with pandas
338         try:
339             # First read the metadata to determine where data starts
340             with open(csv_path, 'r') as f:
341                 lines = f.readlines()
342
343             # Find where the actual data starts (after metadata)
344             data_start = 0
345             for i, line in enumerate(lines):
346                 if 'Academic Year' in line or 'Level of study' in line:
347                     data_start = i
348                     break
349
350             # Read the CSV file starting from the data
351             df = pd.read_csv(csv_path, skiprows=data_start)
352
353             # Clean up column names
354             df.columns = df.columns.str.strip()
355
356             # Ensure required columns exist
357             if 'Level of study' not in df.columns:
358                 raise ValueError("Column 'Level of study' not found in the data")
359             if 'Academic Year' not in df.columns:
360                 raise ValueError("Column 'Academic Year' not found in the data")
361             if 'Number' not in df.columns:
362                 raise ValueError("Column 'Number' not found in the data")
363
364             # Convert Number to numeric, replacing any non-numeric values with NaN
365             df['Number'] = pd.to_numeric(df['Number'], errors='coerce')
366
367             # Drop any rows with missing values
368             df = df.dropna(subset=['Number', 'Level of study', 'Academic Year'])
369
370     except Exception as e:
371         raise ValueError(f"Error reading CSV file: {str(e)}")
372
373     # Create figure (using Agg backend)
374     plt.figure(figsize=(10, 6))
375
376     if chart_type == 'line':
377         # Group by years and plot each level of study
378         for level in df['Level of study'].unique():
379             level_data = df[df['Level of study'] == level]
380             plt.plot(level_data['Academic Year'], level_data['Number'], label=level, marker='o')
381         plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left')
382         plt.xticks(rotation=45)
383         plt.title('Students by Level of Study Over Time')
384
385     elif chart_type == 'bar':
386         # Group by Level of study and sum the numbers
387         data = df.groupby('Level of study')['Number'].sum().sort_values(ascending=False)
388         plt.bar(range(len(data)), data.values)

```

```
398     else:
399         raise ValueError(f"Unsupported chart type: {chart_type}")
400
401     # Adjust layout
402     plt.tight_layout()
403
404     # Save to bytes buffer
405     buffer = io.BytesIO()
406     plt.savefig(buffer, format='png', bbox_inches='tight', dpi=300)
407     buffer.seek(0)
408     image_png = buffer.getvalue()
409
410     # Encode
411     graphic = base64.b64encode(image_png)
412     graphic = graphic.decode('utf-8')
413
414     return CustomJsonResponse({
415         'data': {
416             'image': f'data:image/png;base64,{graphic}'
417         }
418     })
419
420
421
422 except Exception as e:
423     return CustomJsonResponse({'error': str(e)}, status=500)
424
```

F.3 Gets AI recommendations for chart types

```
3522 def get_chart_recommendation(client, dataset_info):      You, last month • For the selected dataset Gemini gives recommend...
3523 """
3524     Get a chart type recommendation from Gemini based on the dataset
3525 """
3526 try:
3527     # Extract relevant information from dataset_info
3528     title = dataset_info.get('title', '')
3529     columns = dataset_info.get('columns', [])
3530     sample_rows = dataset_info.get('rows', [])[:5] # Use up to 5 rows as a sample
3531     query = dataset_info.get('query', '')
3532     institutions = dataset_info.get('institutions', [])
3533     years = dataset_info.get('years', [])
3534     all_year_data = dataset_info.get('allYearData', [])
3535
3536     # Log information about the data
3537     logging.info(f"Chart recommendation for dataset: {title}")
3538     logging.info(f"Years detected: {years}")
3539     logging.info(f"Number of institutions: {len(institutions)}")
3540     logging.info(f"Institutions: {institutions}")
3541
3542     # Extract institution names directly from the dataset
3543     dataset_institutions = []
3544     if sample_rows and len(sample_rows) > 0:
3545         # Find the institution/provider column index
3546         provider_col_idx = -1
3547         for i, col in enumerate(columns):
3548             if col.lower() in ['he provider', 'institution', 'provider', 'university']:
3549                 prompt = f"""
3550 I have a dataset titled "{title}" with the following columns:
3551 {', '.join(columns)}
3552
3553 Here are some sample rows from the dataset:
3554 {format_sample_rows(columns, sample_rows)}
3555
3556 The institutions in this dataset are: {', '.join(institutions)}
3557 The academic years in the dataset are: {', '.join(years)}
3558
3559 IMPORTANT INSTRUCTIONS:
3560 1. Analyze this specific dataset and determine the most appropriate chart type based solely on the data
3561 characteristics.
3562 2. ONLY use the EXACT institution names from the list I provided - do not modify, abbreviate or combine them.
3563 3. Do not mention or suggest institutions that aren't in the provided dataset list.
3564 4. Consider the following chart theory principles:
3565     - Line charts are best for showing trends and changes over time (multiple years/time periods)
3566     - Bar charts are best for comparing data across different categories or groups
3567     - Pie charts are best for showing proportions of parts to a whole (typically at a single point in time)
3568 5. Ignore the user's original query completely and focus only on what visualizations would be most valuable for this
3569 dataset.
3570
3571 Please provide:
3572 1. The recommended chart type (bar, line, or pie) that would be most appropriate for this specific dataset
3573 2. A clear explanation of why this chart type is recommended based on the dataset characteristics
3574 3. Three specific example visualization requests that:
3575     - Only use institutions and years that are available in this dataset
3576     - Must use the EXACT institution names as listed in the dataset - never abbreviate them
3577     - Would work well with the recommended chart type
3578     - Are realistic and can be fulfilled with the available data
```

F.4 Dynamically changing the visualization

```
4536 def change_chart_type(client, dataset_info, chart_type, original_recommendation, current_request=''):      You, last month •
4537 """
4538     Generate a new chart recommendation and examples for a different chart type
4539 """
4540
4541     try:
4542         # Extract relevant information from dataset_info
4543         title = dataset_info.get('title', '')
4544         columns = dataset_info.get('columns', [])
4545         sample_rows = dataset_info.get('rows', [])[:5] # Use up to 5 rows as a sample
4546         query = dataset_info.get('query', '')
4547         institutions = dataset_info.get('institutions', [])
4548         years = dataset_info.get('years', [])
4549         all_year_data = dataset_info.get('allYearData', [])
4550
4551         # Log information about the request
4552         logging.info(f"Chart type change request for dataset: {title}")
4553         logging.info(f"Changing to chart type: {chart_type}")
4554         logging.info(f"Original recommendation: {original_recommendation}")
4555         logging.info(f"Current visualization request: {current_request}")
4556
4557         # Special case checks for incompatible chart types
4558         if chart_type == 'pie' and current_request and ('axis' in current_request.lower() or 'axes' in current_request.lower()
4559             () or 'swap' in current_request.lower()):
4560             compatibility_warning = "Pie charts don't have traditional X and Y axes, so axis swapping or specific axis
4561             assignments don't apply to this chart type."
4562         elif chart_type == 'line' and not dataset_characteristics['has_multiple_years'] and 'year' in current_request.lower():
4563             compatibility_warning = "Line charts work best with multiple time periods. Your request involves a time-based
4564             comparison but the dataset might not have enough time periods to show meaningful trends."
4565
4566         # If we have a compatibility warning, include it in the response
4567         compatibility_note = ""
4568         if compatibility_warning:
4569             compatibility_note = f"<div class='text-amber-600 mb-2 font-medium'><strong>Note:</strong>
4570             {compatibility_warning}</div>"
4571
4572         # Format information about multi-year data if available
4573         multi_year_info = ""
4574         if years and len(years) > 0:
4575             multi_year_info = f"\nThis dataset contains data for multiple academic years: {', '.join(years)}."
4576             multi_year_info += "\nConsider suggesting visualizations that could compare trends across these years or focus
4577             on specific years."
4578
4579     return {
4580         'chart_type': chart_type,
4581         'recommendation': original_recommendation,
4582         'note': compatibility_note,
4583         'multi_year_info': multi_year_info
4584     }
```

```

4690
4691 prompt = f"""
4692 I have a dataset titled "{title}" with the following columns:
4693     ', '.join(columns)}
4694
4695     Here are some sample rows from the dataset:
4696     {format_sample_rows(columns, sample_rows)}
4697
4698     The institutions in this dataset are: ', '.join(institutions)}
4699     The academic years in the dataset are: ', '.join(years)}
4700
4701     IMPORTANT INSTRUCTIONS:
4702     1. I need examples SPECIFICALLY for a {chart_type.upper()} CHART. Every example MUST work with this chart type.
4703     2. ONLY use the EXACT institution names from the list I provided - do not modify, abbreviate or combine them.
4704     3. If a {chart_type} chart isn't ideal for some aspects of the data, still create examples that WILL WORK with a
4705         {chart_type} chart.
4706     4. Your examples will be used as-is and rendered as a {chart_type} chart, so they MUST be appropriate for this
4707         chart type.
4708     5. Focus on what {chart_type} charts do best:
4709         - LINE CHARTS: Show trends over time/years
4710         - BAR CHARTS: Compare values across categories/institutions
4711         - PIE CHARTS: Show proportions of a whole for a single period
4712
4713     Please provide:
4714     1. A brief explanation of how a {chart_type} chart could be used with this specific dataset, noting its
4715         strengths and limitations.
4716     2. 3 example visualization requests that:
4717         - Are SPECIFICALLY designed to work with a {chart_type} chart
4718         - ONLY use institutions and years that are actually available in this dataset (as listed above)
4719         - Are realistic and can be fulfilled with the available data
4720
4721 prompt = f"""
4722 I have a dataset titled "{title}" with the following columns:
4723     ', '.join(columns)}
4724
4725     Here are some sample rows from the dataset:
4726     {format_sample_rows(columns, sample_rows)}
4727
4728     The institutions in this dataset are: ', '.join(institutions)}
4729     The academic years in the dataset are: ', '.join(years)}
4730
4731     IMPORTANT INSTRUCTIONS:
4732     1. I need examples SPECIFICALLY for a {chart_type.upper()} CHART. Every example MUST work with this chart type.
4733     2. ONLY use the EXACT institution names from the list I provided - do not modify, abbreviate or combine them.
4734     3. If a {chart_type} chart isn't ideal for some aspects of the data, still create examples that WILL WORK with a
4735         {chart_type} chart.
4736     4. Your examples will be used as-is and rendered as a {chart_type} chart, so they MUST be appropriate for this
4737         chart type.
4738     5. Focus on what {chart_type} charts do best:
4739         - LINE CHARTS: Show trends over time/years
4740         - BAR CHARTS: Compare values across categories/institutions
4741         - PIE CHARTS: Show proportions of a whole for a single period
4742
4743     Please provide:
4744     1. A brief explanation of how a {chart_type} chart could be used with this specific dataset, noting its
4745         strengths and limitations.
4746     2. 3 example visualization requests that:
4747         - Are SPECIFICALLY designed to work with a {chart_type} chart
4748         - ONLY use institutions and years that are actually available in this dataset (as listed above)
4749         - Are realistic and can be fulfilled with the available data

```

G: Front-End Rendering

G.1 Rendering Visualization

```
1332 function renderChart(chartConfig) {      You, last month • For the selected dataset Gemini gives recommend...
1333     // Get the canvas element for rendering
1334     const canvas = document.getElementById('visualizationChart');
1335     const ctx = canvas.getContext('2d');
1336
1337     try {
1338         // If we have a previous chart, destroy it properly
1339         if (window.currentChart) {
1340             window.currentChart.destroy();
1341             window.currentChart = null;
1342         }
1343
1344         // In case the chart wasn't properly destroyed, clear the canvas manually
1345         canvas.width = canvas.width; // This effectively clears the canvas
1346
1347         // Create a new chart with the provided configuration
1348         window.currentChart = new Chart(ctx, chartConfig);
1349
1350         // Ensure chart is visible
1351         document.getElementById('visualizationContainer').classList.remove('hidden');
1352         document.getElementById('visualizationError').classList.add('hidden');
1353
1354     } catch (error) {
1355         console.error("Error rendering chart:", error);
1356
1357         // Display error to user with friendly message based on chart type
1358         let errorMessage = "The chart couldn't be displayed with your current settings.";
1359
1360         if (chartConfig && chartConfig.type) {
1361             errorMessage = `Your ${chartConfig.type} chart couldn't be created with the selected data.`;
1362             errorMessage += "<br>Try selecting a different chart type or modifying your request.";
1363         }
1364
1365         document.getElementById('visualizationErrorMessage').innerHTML = errorMessage;
1366         document.getElementById('visualizationError').classList.remove('hidden');
1367         document.getElementById('visualizationContainer').classList.add('hidden');
1368     }
1369 }
```

G.2 Interaction Control

```
1112 function generateVisualization(datasetInfo, request, chartType) {      You, last month • so progress with charts but still n
1113     const loadingEl = document.getElementById('visualizationLoading');
1114     const errorEl = document.getElementById('visualizationError');
1115     const errorMsgEl = document.getElementById('visualizationErrorMessage');
1116     const alternativesEl = document.getElementById('alternativeSuggestions');
1117     const suggestionsListEl = document.getElementById('suggestionsList');
1118     const visualizationEl = document.getElementById('visualizationContainer');
1119     const insightsEl = document.getElementById('insightsContent');
1120
1121     // Store the current request for future use
1122     currentRequest = request;
1123
1124     // Hide previous results
1125     errorEl.classList.add('hidden');
1126     alternativesEl.classList.add('hidden');
1127     visualizationEl.classList.add('hidden');
1128
1129     // Show loading state
1130     loadingEl.classList.remove('hidden');
1131
1132     // Prepare data for API call
1133     const requestData = {
1134         action: 'generate_visualization',
1135         dataset_info: datasetInfo,
1136         user_request: request
1137     };
1138
1139     // If a specific chart type was requested, include it
1140     if (chartType) {
1141         requestData.chart_type = chartType;
1142
1143         // Set active chart button to reflect the selected type
1144         setActiveChartButton(chartType);
1145
1146         // Add a confirmation that this chart type is being enforced
1147         console.log(`Enforcing chart type: ${chartType} for visualization`);
1148
1149         // If the chart type wasn't specifically mentioned in the request, add info
1150         if (!request.toLowerCase().includes(chartType.toLowerCase())) {
1151             loadingEl.innerHTML = `<span class="text-sm text-blue-600">Generating a ${chartType} chart as requested...</span>`;
1152         }
1153     }
1154
1155     // Call backend API for Gemini visualization generation
1156     fetch('/visualization_api/', {
1157         method: 'POST',
1158         headers: {
1159             'Content-Type': 'application/json',
1160             'X-CSRFToken': getCsrfToken()
1161         },
1162         body: JSON.stringify(requestData)
1163     })
1164     .then(response => response.json())
1165     .then(data => {
```

```

1166     // Hide loading state
1167     loadingEl.classList.add('hidden');
1168
1169     // Store visualization data for feedback
1170     window.lastVisualizationData = {
1171       chart_type: chartType || 'auto',
1172       request: request,
1173       success: data.success,
1174       has_compatibility_warning: data.has_compatibility_warning || false,
1175       timestamp: new Date().toISOString()
1176     };
1177
1178   if (data.success) {
1179     // If successful, render the chart
1180     try {
1181       console.log("Chart config received:", data.chart_config);
1182
1183       // Safely evaluate the chart configuration
1184       let chartConfig;
1185       try {
1186         // First attempt: directly evaluate the config after cleaning it
1187         const cleanConfig = data.chart_config
1188           .replace(/Object\.\assign/g, '') // Remove Object.assign calls
1189           .replace(/['"]\s*\+\s*['"]|['"]\s*+|[']\+\s*['']|\+\s*['']/g, '') // Fix string concatenation
1190           .replace(/(['])?([a-zA-Z0-9_]+)(['])?:/g, '"$2":'); // Normalize property names to double quotes
1191
1192         chartConfig = eval('(' + cleanConfig + ')');
1193       } catch (evalError) {
1194         console.error("Error in first evaluation attempt:", evalError);
1195
1196         try {
1197           // Replace problematic quotes and syntax
1198           let sanitizedConfig = data.chart_config
1199             .replace(/Object\.\assign/g, '') // Remove Object.assign calls
1200             .replace(/(['"])?([a-zA-Z0-9_]+)(['])?:/g, '"$2":') // Ensure property names are quoted
1201             .replace(/:\s*(['"]*)\1/g, ':"$2"') // Ensure property values use double quotes
1202             .replace(/,(s*[\r\n])/g, '$1'); // Remove trailing commas
1203
1204             chartConfig = JSON.parse(sanitizedConfig);
1205         } catch (jsonError) {
1206           console.error("Error in second evaluation attempt:", jsonError);
1207
1208           // Third attempt: try to fix common issues with template literals and Object references
1209           try {
1210             let fixedConfig = data.chart_config
1211               .replace(/Object/g, '') // Remove 'Object'
1212               .replace(/\.\assign/g, '') // Remove '.assign'
1213               .replace(/`([^\`]*`)/g, '"$1"') // Replace template literals with double quotes
1214               .replace(/\$\{([^\}]*\}\)/g, '') // Remove template expressions
1215               .replace(/\\\v/g, ''); // Remove escaped characters
1216
1217
1218             if (!fixedConfig.startsWith('{')) {
1219               fixedConfig = '{' + fixedConfig;
1220             }
1221             if (!fixedConfig.endsWith('}')) {
1222               fixedConfig = fixedConfig + '}';
1223             }
1224           }
1225         }
1226       }
1227     }
1228   }
1229
1230   // Render the visualization
1231   render();

```

```
1233 // Only try to render if we have a valid chart config
1234 if (chartConfig && typeof chartConfig === 'object') {
1235     // STRICTLY enforce the selected chart type
1236     if (chartType) {
1237         console.log(`Enforcing chart type: ${chartType} (was: ${chartConfig.type})`);
1238         // Force the selected chart type to ensure it matches what the user selected
1239         chartConfig.type = chartType;
1240     }
1241
1242     // Add compatibility warning if the chart configuration suggests a different type
1243     if (chartType && data.has_compatibility_warning) {
1244         console.log("Chart has compatibility warning");
1245         // Check if the warning is already in the insights to avoid duplication
1246         if (data.insights && !data.insights.includes(data.compatibility_warning)) {
1247             insightsEl.innerHTML =
1248                 `<div class="bg-yellow-100 border-1-4 border-yellow-500 text-yellow-700 p-4 mb-4">
1249                 <p class="font-bold">Note about ${chartType} chart:</p>
1250                 <p>${data.compatibility_warning} || ^A ${chartType} chart may not be the optimal choice
1251             </div>
1252             ` + (data.insights || "");
1253         } else {
1254             insightsEl.innerHTML = data.insights || "";
1255         }
1256     } else {
1257         insightsEl.innerHTML = data.insights || "";
1258     }
1259 }
```

G.3 Changing Chart Type

```

1002 function changeChartType(chartType) {      You, last month * so progress with charts but still not fully wor...
1003     // Don't do anything if we haven't received the original recommendation yet
1004     if (!originalRecommendation) {
1005         console.warn("Cannot change chart type before receiving original recommendation");
1006         return;
1007     }
1008
1009     const loadingEl = document.getElementById('visualizationLoading');
1010     const recommendationEl = document.getElementById('chartRecommendation');
1011     const reasonEl = document.getElementById('recommendationReason');
1012     const chartTypeEl = document.getElementById('recommendedChartType');
1013
1014     // Show loading state
1015     loadingEl.classList.remove('hidden');
1016
1017     // Get the current visualization request if any
1018     currentRequest = document.getElementById('visualizationRequest').value || '';
1019
1020     // Prepare data for API call
1021     const requestData = {
1022         action: 'change_chart_type',
1023         dataset_info: window.datasetInfo, // Global variable set during page load
1024         chart_type: chartType,
1025         original_recommendation: JSON.stringify(originalRecommendation),
1026         current_request: currentRequest
1027     };
1028
1029     console.log("Sending chart type change request:", requestData);
1030     // Call backend API
1031     fetch('/visualization_api/', {
1032         method: 'POST',
1033         headers: {
1034             'Content-Type': 'application/json',
1035             'X-CSRFToken': getCsrfToken()
1036         },
1037         body: JSON.stringify(requestData)
1038     })
1039     .then(response => response.json())
1040     .then(data => {
1041         // Hide loading state
1042         loadingEl.classList.add('hidden');
1043         console.log("Chart type change response:", data);
1044
1045         if (data.success) {
1046             // Use our new chart type info update function to handle all UI updates
1047             updateChartTypeInfo(data);
1048
1049             // If there is a current visualization, regenerate it with the new chart type
1050             if (currentRequest && currentRequest.trim() !== '') {
1051                 // Clear previous visualization
1052                 if (window.currentChart) {
1053                     window.currentChart.destroy();
1054                     window.currentChart = null;
1055                 }
1056
1057                 // Hide visualization container

```

```

1059     document.getElementById('visualizationContainer').classList.add('hidden');
1060
1061     // Regenerate the visualization with the new chart type
1062     generateVisualization(window.datasetInfo, currentRequest, chartType);
1063   } else {
1064     // If there's no current request, just clear any existing visualization
1065     if (window.currentChart) {
1066       window.currentChart.destroy();
1067       window.currentChart = null;
1068     }
1069     document.getElementById('visualizationContainer').classList.add('hidden');
1070     document.getElementById('visualizationError').classList.add('hidden');
1071   }
1072 } else {
1073   showError(data.error || "Failed to change chart type");
1074 }
1075
1076 })
1077 .catch(error => {
1078   loadingEl.classList.add('hidden');
1079   showError("Network error: " + error.message);
1080   console.error("Error changing chart type:", error);
1081 });
1082

```

H: Interactive Visualization Controls

H.1 Back-end connection

```

3480 @csrf_exempt
3481 def visualization_api(request):
3482     """
3483     API endpoint for visualization recommendations and generation
3484     """
3485     if request.method != 'POST':
3486         return JsonResponse({'success': False, 'error': 'Only POST method is allowed'}, status=405)
3487
3488     try:
3489         data = json.loads(request.body)
3490         action = data.get('action')
3491         dataset_info = data.get('dataset_info')
3492
3493         if not action or not dataset_info:
3494             return JsonResponse({'success': False, 'error': 'Missing required parameters'}, status=400)
3495
3496         client = get_llm_client()
3497
3498         if action == 'get_recommendation':
3499             return get_chart_recommendation(client, dataset_info)
3500         elif action == 'change_chart_type':
3501             chart_type = data.get('chart_type')
3502             original_recommendation = data.get('original_recommendation')
3503             current_request = data.get('current_request', '')
3504             if not chart_type:
3505                 return JsonResponse({'success': False, 'error': 'Missing chart type'}, status=400)
3506             return change_chart_type(client, dataset_info, chart_type, original_recommendation, current_request)
3507         elif action == 'generate_visualization':
3508             user_request = data.get('user_request')
3509             chart_type = data.get('chart_type', None) # Get the requested chart type if provided
3510             if not user_request:
3511                 return JsonResponse({'success': False, 'error': 'Missing user request'}, status=400)
3512             return generate_visualization(client, dataset_info, user_request, chart_type)
3513         else:
3514             return JsonResponse({'success': False, 'error': f'Unknown action: {action}'}, status=400)
3515
3516     except json.JSONDecodeError:
3517         return JsonResponse({'success': False, 'error': 'Invalid JSON'}, status=400)
3518     except Exception as e:
3519         logging.exception("Error in visualization API: %s", str(e))
3520         return JsonResponse({'success': False, 'error': f'Server error: {str(e)}'}, status=500)

```

H.2 Chart Type Explanation

```
4802 def get_chart_type_explanation(chart_type, dataset_characteristics):
4803     """
4804     Get a specific explanation for the requested chart type based on dataset characteristics
4805     """
4806     if chart_type == 'line':
4807         if not dataset_characteristics.get("has_multiple_years", False):
4808             return "Line charts are most effective for showing trends over time. Since this dataset only contains data for a single year, you should focus on comparing values across different categories or institutions rather than showing trends over time."
4809     else:
4810         return "Line charts work best for showing trends over time or continuous data. Focus on how values change across the available years for one or more institutions."
4811
4812     elif chart_type == 'bar':
4813         if dataset_characteristics.get("has_multiple_institutions", False) and dataset_characteristics.get("has_multiple_years", False):
4814             return "Bar charts are excellent for comparing values across categories. You can compare institutions side by side, or show change over time by grouping bars by year."
4815         else:
4816             return "Bar charts are ideal for comparing values across categories. Focus on comparing different metrics or institutions within the available data."
4817
4818     elif chart_type == 'pie':
4819         if dataset_characteristics.get("has_multiple_institutions", False) and len(dataset_characteristics.get("available_institutions", [])) > 5:
4820             return "Pie charts work best with a small number of categories (5 or fewer) and when showing parts of a whole. Consider focusing on just a few key institutions or categories to maintain clarity."
4821         else:
4822             return "Pie charts are best for showing proportions or percentages of a whole. Focus on the distribution of a single metric across categories for a specific time period rather than comparisons over time."
4823
4824     return None| You, last month * so progress with charts but still not fully wor...
4825
4826 # Configure logging to prevent duplicate handlers
4827 logger = logging.getLogger('core.views')
4828 # Remove all handlers to avoid duplicates
4829 if logger.handlers:
4830     for handler in logger.handlers:
4831         logger.removeHandler(handler)
4832 # Add a single handler
4833 handler = logging.StreamHandler()
4834 formatter = logging.Formatter('%(asctime)s - %(name)s - %(levelname)s - %(message)s')
4835 handler.setFormatter(formatter)
4836 logger.addHandler(handler)
4837 logger.propagate = False # Prevent propagation to root logger
4838
```

I: Dashboard Interface

I.1 Display Query Results

```
281 |     function displayQueryAnalysis(data) {
282 |         console.log('Displaying query analysis:', data);
283 |
284 |         // Store the data for later use
285 |         window.currentQueryData = data;
286 |
287 |         if (!aiQueryResults) {
288 |             console.error('Query results container not found');
289 |             return;
290 |         }
291 |
292 |         // Check if there was an API error
293 |         if (data.api_error) {
294 |             aiQueryResults.innerHTML =
295 |                 <div class="bg-red-100 border-1-4 border-red-500 text-red-700 p-4 mt-4" role="alert">
296 |                     <p class="font-bold">Gemini API Error</p>
297 |                     <p>${data.api_error}</p>
298 |                     <p class="mt-2 text-sm">Using fallback analysis instead. For accurate entity extraction, please ensure
299 |                         the Gemini API key is configured correctly.</p>
300 |                 </div>
301 |
302 |             // We'll still show the fallback results below the error message
303 |         }
304 |
305 |         // Create HTML for the query analysis with collapsible section
306 |         let resultsHTML =
307 |             <div class="bg-white shadow-md rounded-lg p-6 mt-8">
308 |                 <div class="flex justify-between items-center mb-4">
309 |                     <h3 class="text-lg font-semibold">
310 |                         Query Analysis ${data.using_mock ? '(Mock AI)' : 'by Gemini AI'}
311 |                         ${data.using_mock ?
312 |                             '<span class="text-xs font-normal bg-yellow-100 text-yellow-800 px-2 py-1 rounded ml-2">Using
313 |                             offline AI simulation</span>' :
314 |                             '<span class="text-xs font-normal bg-green-100 text-green-800 px-2 py-1 rounded ml-2">Using
315 |                             Google Gemini API</span>'}
316 |                         ${data.cached === true ?
317 |                             '<span class="text-xs font-normal bg-purple-100 text-purple-800 px-2 py-1 rounded ml-2">⚡
318 |                             Retrieved from cache</span>' : ''}
319 |                     </h3>
320 |                     <button id="queryDetailsToggle" class="bg-blue-100 hover:bg-blue-200 text-blue-800 px-3 py-1 rounded
321 |                         text-sm transition-colors duration-150 ease-in-out">
322 |                         <span class="show-text">Show Details</span>
323 |                         <span class="hide-text hidden">Hide Details</span>
324 |                     </button>
325 |                 </div>
326 |
327 |                 <div class="mb-4">
328 |                     <p class="font-medium text-gray-700">Your query:</p>
329 |                     <p class="text-blue-800 bg-blue-50 p-2 rounded">${data.query}</p>
330 |                 </div>
331 |
332 |                 <div id="queryDetailsContent" class="hidden">
```

```

345 `      ${data.has_institution_typos || data.has_year_typos ? ``}
346   <div class="mb-4 bg-blue-50 p-3 rounded border border-blue-200">
347     <p class="font-medium text-gray-700">Corrected query:</p>
348     <p class="text-blue-800 p-2">${generateCorrectedQuery(data)}</p>
349   </div>
350   ` : ''
351
352   <div class="grid grid-cols-1 md:grid-cols-3 gap-4">
353     <div class="border rounded p-3">
354       <p class="font-medium text-gray-700">Institutions:</p>
355       <ul class="list-disc list-inside">
356         ${data.institutions && data.institutions.length > 0 ?
357           data.institutions.map(inst => `<li>${inst}</li>`).join('') :
358           '<li class="text-gray-500">None specified</li>'}
359       </ul>
360     </div>
361
362     <div class="border rounded p-3">
363       <p class="font-medium text-gray-700">Years:</p>
364       <ul class="list-disc list-inside">
365         ${data.years && data.years.length > 0 ?
366           data.years.map(year => `<li>${year}</li>`).join('') :
367           '<li class="text-gray-500">All available years</li>'}
368       </ul>
369     </div>
370
371     <div class="border rounded p-3">
372       <p class="font-medium text-gray-700">Year Range:</p>
373       <p class="${data.start_year ? 'font-semibold text-blue-700' : 'text-gray-500'}">
374         ${yearRangeDisplay}
375       </p>
376     </div>
377   </div>
378
379   ${data.mission_group ? `

380   <div class="mt-4 border rounded p-3 bg-blue-50">
381     <p class="font-medium text-gray-700">Mission Group Filter:</p>
382     <p class="font-semibold text-blue-700">${data.mission_group}</p>
383     <p class="text-sm mt-1">Including data for ${data.mission_group_institutions?.length || 0} institutions from this group</p>
384   </div>
385   ` : ''}

386
387   <div class="mt-4 border rounded p-3">
388     <p class="font-medium text-gray-700">Data requested:</p>
389     <p class="font-semibold ${dataRequestDisplay === 'General information' ? 'text-gray-500' : 'text-blue-700'}">
390       ${dataRequestDisplay}
391     </p>
392   </div>
393
394   <div class="mt-4 text-sm text-gray-600 border-t pt-2">
395     <p>Entity extraction performed by ${data.using_mock ? `local AI simulation (Gemini API unavailable)` : `Google's Gemini AI`}</p>
396     ${data.cached === true ?
397       '<p>Results retrieved from cache. Last processed: ' + new Date().toLocaleString() + '</p>' : ''}
398   </div>
399
400 // Add feedback component at the top
401 resultsHTML = `
402   <div id="feedbackContainer" class="bg-yellow-50 border border-yellow-200 rounded-lg p-6 mb-6 mt-6">
403     <h4 class="text-lg font-semibold mb-2">How helpful were these results?</h4>
404     <div class="flex flex-wrap gap-2">
405       <button id="veryHelpfulBtn" class="feedback-btn bg-green-100 hover:bg-green-200 text-green-800 px-4 rounded-md" data-value="Very Helpful">Very Helpful</button>
406       <button id="helpfulBtn" class="feedback-btn bg-blue-100 hover:bg-blue-200 text-blue-800 px-4 py-2 rounded-md" data-value="Helpful">Helpful</button>
407       <button id="notHelpfulBtn" class="feedback-btn bg-red-100 hover:bg-red-200 text-red-800 px-4 py-2 rounded-md" data-value="Not Helpful">Not Helpful</button>
408     </div>
409   </div>
410   ` + resultsHTML;
411
412
413
414
415
416

```

I.2 Reviewable Dataset

```
Windsurf: Refactor | Explain | X
720  function renderFilePreview(dataset, datasetIndex) {      You, last month • important stage ...
721    console.log(`Rendering file previews for dataset ${datasetIndex + 1}: ${dataset.title || 'Untitled'}`);
722
723    const filesContainer = document.getElementById(`files-container-${datasetIndex}`);
724    if (!filesContainer) {
725      console.error(`Files container not found for dataset ${datasetIndex}`);
726      return;
727    }
728
729    if (!dataset.matches || dataset.matches.length === 0) {
730      filesContainer.innerHTML = `
731        <div class="text-sm text-gray-500 p-3 bg-gray-50 rounded">
732          No files available
733        </div>
734      `;
735      return;
736    }
737
738    let filesHTML = '';
739
740    dataset.matches.forEach(match => {
741      console.log(`Processing match for dataset ${datasetIndex + 1}:`, match);
742
743      filesHTML += `
744        <div class="mt-4 mb-4 border rounded p-3 bg-white">
745          <div class="text-sm font-medium mb-2">${match.academic_year || 'Unknown Year'} - ${match.reference || 'Unknown Reference'}</div>
746
747          if (match.preview && match.preview.columns && match.preview.data) {
748            console.log(`Processing preview data for ${match.reference}`, {
749              columns: match.preview.columns.length,
750              data: match.preview.data.length
751            });
752
753            try {
754              // Clean metadata artifacts from columns and data
755              const { columns: uniqueColumns, data: previewData } = cleanMetadataArtifacts(
756                match.preview.columns,
757                match.preview.data
758              );
759
760              filesHTML += `
761                <div class="overflow-x-auto max-h-[300px] table-container border rounded">
762                  <table class="min-w-full border-collapse table-auto text-sm">
763                    <thead>
764                      <tr>
765                        ${uniqueColumns.map(column =>
766                          `<th class="px-4 py-2 border-b border-gray-300 text-left text-sm font-medium sticky top-0 bg-white z-10">${column}</th>`).
767                          .join('')}
768                      </tr>
769                    </thead>
770                    <tbody>
771                      ${previewData.length > 0 ?
772                        previewData.map(row => `<tr>
773
```

```

775           `${row.map(cell =>
776             `<td class="px-4 py-2 border-b border-gray-300 text-sm">${cell}</td>`  

777           ).join('')}`  

778         `).join('') :  

779         `<tr><td colspan="${uniqueColumns.length}" class="px-4 py-2 text-center  

780           border-gray-500">No data available or no matching institutions found</td></tr>  

781       }  

782     `</tbody>  

783   `</table>  

784 `</div>  

785 ${match.preview.has_more ?  

786   `<div class="text-xs border-gray-500 mt-2"  

787     Showing ${previewData.length} of ${match.preview.matched_rows || previewData.length}  

788     matching rows  

789   ` :  

790 }  

791 `;  

792 } catch (error) {  

793   console.error('Error processing preview data:', error);  

794   filesHTML +=  

795     `<div class="text-sm border-red-500 p-3 bg-red-50 rounded">  

796       Error processing preview data: ${error.message}  

797     `;  

798 }  

799 } else if (match.preview && match.preview.error) {  

800   filesHTML +=  

801     `<div class="text-sm border-red-500 p-3 bg-red-50 rounded">  

802       Error loading preview: ${match.preview.error}  

803     `;  

804 } else {  

805   filesHTML +=  

806     `<div class="text-sm border-gray-500 p-3 bg-gray-50 rounded">  

807       Preview not available  

808     `;  

809 }  

810 `;  

811 `;  

812 }  

813 `;  

814 filesHTML += `</div>`;  

815 `});  

816 `filesContainer.innerHTML = filesHTML;  

817 `};  

818 `}

```

I.3 Manage loading states

```
821     function showLoading(message = 'Processing your query...') {
822       console.log('Showing loading indicator');
823       let loadingIndicator = document.getElementById('aiLoadingIndicator');
824       if (!loadingIndicator) {
825         // Create loading indicator if it doesn't exist
826         loadingIndicator = document.createElement('div');
827         loadingIndicator.id = 'aiLoadingIndicator';
828         loadingIndicator.className = 'fixed top-0 left-0 w-full h-full flex items-center justify-center bg-gray-800 bg-opacity-20 rounded-lg shadow-lg flex-col';
829         loadingIndicator.innerHTML = `
830           <div class="backdrop-blur-md bg-gray-800 bg-opacity-20 p-8 rounded-lg flex flex-col items-center">
831             <div class="w-16 h-16 mb-4">
832               <svg class="animate-spin w-full h-full text-blue-500" xmlns="http://www.w3.org/2000/svg" fill="none">
833                 <circle class="opacity-25" cx="12" cy="12" r="10" stroke="currentColor" stroke-width="4"></circle>
834                 <path class="opacity-75" fill="currentColor" d="M4 12a8 8 0 018-8V0c5.373 0 5.373 0 12h4zm2 5.208 0 0" />
835               </svg>
836             </div>
837             <h2 class="text-xl font-semibold mb-2 text-white">Analyzing Data</h2>
838             <p class="text-gray-200 text-center">${message}</p>
839           </div>
840         `;
841         document.body.appendChild(loadingIndicator);
842         console.log('Loading indicator created and added to page');
843       } else {
844         const messageElement = loadingIndicator.querySelector('p');
845         if (messageElement) {
846           messageElement.textContent = message;
847         }
848         // Show existing loading indicator
849         loadingIndicator.classList.remove('hidden');
```

I.4 Query Input

```

88  if (aiSearchBtn && aiQueryInput) {
89    aiSearchBtn.addEventListener('click', function() {      You, 2 months ago • created indeing file for the AI approach.
90      const query = aiQueryInput.value.trim();
91
92      console.log('AI Search button clicked with query:', query);
93
94      if (!query) {
95        alert('Please enter a query');
96        return;
97      }
98
99      // Record the start time of the query processing
100     window.queryStartTime = performance.now();
101     console.log('Query processing started at:', window.queryStartTime);
102
103     // Show loading state
104     showLoading('Analyzing your query with Gemini AI...');
105
106     // Prepare max matches parameter
107     const maxMatches = aiMaxMatches ? aiMaxMatches.value : 3;
108
109     // Get selected mission group
110     const selectedMissionGroup = document.querySelector('input[name="missionGroup"]:checked').value;
111     console.log('Selected mission group:', selectedMissionGroup);
112
113     // Call the Gemini API endpoint
114     fetch('/process_gemini_query', {
115       method: 'POST',
116       headers: {
117         'Content-Type': 'application/json',
118         'X-CSRFToken': getCsrfToken(),
119       },
120       body: JSON.stringify({
121         query,
122         max_matches: maxMatches,
123         mission_group: selectedMissionGroup === 'none' ? null : selectedMissionGroup,
124         start_time: window.queryStartTime
125       })
126     })
127     .then(response => {
128       console.log('Received response with status:', response.status);
129
130       if (!response.ok) {
131         return response.text().then(text => {
132           console.error('Error response body:', text);
133           throw new Error('Network response was not ok: ' + response.status);
134         });
135       }
136       return response.json();
137     })
138     .then(data => {
139       console.log('Received data from API:', data);
140       // Calculate the end time and total duration
141       const endTime = performance.now();
142       const duration = endTime - window.queryStartTime;
143       console.log(`Query completed in ${duration.toFixed(2)} ms`);
144
145       // Send timing information to the server

```

```

146
147
148
149
150     recordQueryTiming(query, duration, data.matching_datasets ? data.matching_datasets.length : 0, selectedMissions);
151
152     hideLoading();
153
154     if (data.status === 'error') {
155         console.error('Error from backend:', data.error);
156         showError(data.error);
157         return;
158     }
159
160     // Display the query analysis
161     displayQueryAnalysis(data);
162
163     .catch(error => {
164         console.error('Fetch error:', error);
165         hideloading();
166         showError('An error occurred while processing your query. Please try again.');
167     });
168 };
169

```

I.5: Exporting data

```

618 document.addEventListener('DOMContentLoaded', function() {
619     // CSV download buttons
620     document.querySelectorAll('.download-csv-btn').forEach(button => {
621         button.addEventListener('click', function(e) {
622             // Prevent event from propagating to document level handlers
623             e.stopPropagation();
624
625             const tableId = this.getAttribute('data-table-id');
626             const tableIndex = tableId.split('-').pop();
627             const filename = `dataset_${tableIndex}.csv`;
628             downloadTableAsCSV(tableId, filename);
629
630             // Prevent default button behavior
631             e.preventDefault();
632         });
633     });
634
635     // PDF download buttons
636     document.querySelectorAll('.download-pdf-btn').forEach(button => {
637         button.addEventListener('click', function(e) {
638             // Prevent event from propagating to document level handlers
639             e.stopPropagation();
640
641             const tableId = this.getAttribute('data-table-id');
642             const tableIndex = tableId.split('-').pop();
643             const filename = `dataset_${tableIndex}.pdf`;
644             downloadTableAsPDF(tableId, filename);
645
646             // Prevent default button behavior
647             e.preventDefault();
648         });
649     });
650
651     // Excel download buttons
652     document.querySelectorAll('.download-excel-btn').forEach(button => {
653         button.addEventListener('click', function(e) {
654             // Prevent event from propagating to document level handlers
655             e.stopPropagation();
656
657             const tableId = this.getAttribute('data-table-id');
658             const tableIndex = tableId.split('-').pop();
659             const filename = `dataset_${tableIndex}.xlsx`;
660             downloadTableAsExcel(tableId, filename);
661
662             // Prevent default button behavior
663             e.preventDefault();
664         });
665     });
666 });

```

I.6 Prepare Data for Clean Display

```
211  function cleanMetadataArtifacts(columns, tableData) {      You, last month * important stage ...
212    console.log('Cleaning metadata artifacts from columns:', columns);
213
214    // Handle case where no columns are provided
215    if (!columns || columns.length === 0) {
216      console.warn('No columns provided to cleanMetadataArtifacts');
217      return { columns: [], data: tableData || [] };
218    }
219
220    // Remove any metadata string from column headers
221    const cleanColumns = columns.map(col => {
222      // Remove #METADATA and any JSON that follows it
223      if (typeof col === 'string') {
224        // Handle metadata in column headers
225        // Test Regex...
226        const metadataRemoved = col.replace(/#METADATA:.*?}/g, '').trim();
227
228        // Make column names look nicer (capitalize first letter of each word)
229        // Test Regex...
230        return metadataRemoved.split(/\s+/).map(word =>
231          | word.charAt(0).toUpperCase() + word.slice(1).toLowerCase()
232          ).join(' ');
233      }
234      return col;
235    });
236
237    console.log('Columns after metadata removal:', cleanColumns);
238
239    // Remove duplicate "Academic Year" columns
240    const uniqueColumns = [];
241    const seenColumns = new Set();
242
243    cleanColumns.forEach(col => {
244      if (!seenColumns.has(col)) {
245        uniqueColumns.push(col);
246        seenColumns.add(col);
247      } else {
248        console.log(`Duplicate column removed: ${col}`);
249      }
250    });
251
252    console.log('Final unique columns:', uniqueColumns);
253
254    // Clean row data - remove any metadata strings and ensure correct length
255    const cleanData = (tableData || []).map(row => {
256      // Handle metadata in row data
257      const cleanRow = row.map(cell => {
258        if (typeof cell === 'string') {
259          // Test Regex...
260          return cell.replace(/#METADATA:.*?}/g, '').trim();
261        }
262        return cell;
263      });
264
265      // Adjust row length to match columns
266      if (cleanRow.length > uniqueColumns.length) {
267        console.log(`Trimming row from ${cleanRow.length} to ${uniqueColumns.length} columns`);
268        return cleanRow.slice(0, uniqueColumns.length);
269      } else if (cleanRow.length < uniqueColumns.length) {
270        console.log(`Padding row from ${cleanRow.length} to ${uniqueColumns.length} columns`);
271        const paddedRow = [...cleanRow];
272        while (paddedRow.length < uniqueColumns.length) {
273          paddedRow.push('');
274        }
275        return paddedRow;
276      }
277      return cleanRow;
278    });
279
280    return { columns: uniqueColumns, data: cleanData };
281  }
```

I.7 Dynamic dataset preview toggle

```
toggleBtn.addEventListener('click', function() {      You, 4 ✓
    // Toggle content visibility
    contentDiv.classList.toggle('hidden');

    // Toggle button text
    showText.classList.toggle('hidden');
    hideText.classList.toggle('hidden');

    // Add smooth animation (if CSS is defined for it)
    if (!contentDiv.classList.contains('hidden')) {
        contentDiv.classList.add('animate-fade-in');
        setTimeout(() => {
            contentDiv.classList.remove('animate-fade-in');
        }, 300);
    }
});
```

J: Feedback System

```
1229 |     function setupFeedbackButtons(query) {      You, 4 weeks ago * working on feedback logging ...
1230 |       const feedbackButtons = document.querySelectorAll('.feedback-btn');
1231 |
1232 |       feedbackButtons.forEach(button => {
1233 |         button.addEventListener('click', function() {
1234 |           const feedback = this.textContent.trim();
1235 |
1236 |           // Disable all buttons to prevent multiple selections
1237 |           feedbackButtons.forEach(btn => btn.disabled = true);
1238 |
1239 |           // If "Helpful" or "Not Helpful" is clicked, show comment textarea
1240 |           if (feedback === "Helpful" || feedback === "Not Helpful") {
1241 |             // Create comment section
1242 |             const commentSection = document.createElement('div');
1243 |             commentSection.className = 'mt-4';
1244 |             commentSection.innerHTML = `
1245 |               <p class="mb-2 font-medium text-gray-700">How could we improve these results?</p>
1246 |               <textarea id="feedbackComment" class="w-full px-3 py-2 border border-gray-300 rounded-md focus:outline-none focus:ring-2 focus:ring-blue-500"
1247 |                 rows="3" placeholder="Your comments will help us improve our search results (optional)"></textarea>
1248 |               <div class="mt-2 flex justify-end">
1249 |                 <button id="submitFeedbackBtn" class="bg-blue-600 hover:bg-blue-700 text-white px-4 py-2 rounded-md">
1250 |                   Submit Feedback
1251 |                 </button>
1252 |               </div>
1253 |             `;
1254 |
1255 |             // Add to container
1256 |             document.getElementById('feedbackContainer').appendChild(commentSection);
1257 |
1258 |             // Enable the selected button with visual indication
1259 |             this.disabled = false;
1260 |             this.classList.add('ring-2', 'ring-blue-500');
1261 |
1262 |             // Handle submit button click
1263 |             document.getElementById('submitFeedbackBtn').addEventListener('click', function() {
1264 |               const commentText = document.getElementById('feedbackComment').value;
1265 |               submitFeedback(feedback, query, commentText);
1266 |             });
1267 |           } else {
1268 |             // If "Very Helpful" is clicked, submit immediately without comment
1269 |             // Show loading state in the button
1270 |             const originalText = this.textContent;
1271 |             this.innerHTML = `
1272 |               <div class="flex items-center">
1273 |                 <svg class="animate-spin -ml-1 mr-2 h-4 w-4 text-current" xmlns="http://www.w3.org/2000/svg"
1274 |                   fill="none" viewBox="0 0 24 24">
1275 |                   <circle class="opacity-25" cx="12" cy="12" r="10" stroke="currentColor" stroke-width="4"></circle>
1276 |                   <path class="opacity-75" fill="currentColor" d="M4 12a8 8 0 018-8V0C5.373 0 5.373 0
1277 |                     12h4zm2 5.291A7.962 7.962 0 0114 12H0c0 3.042 1.135 5.824 3 7.938l3-2.647z"></path>
1278 |                 </svg>
1279 |                 Submitting...
1280 |               </div>
1281 |             `;
1282 |           }
1283 |         }
1284 |       }
1285 |     }
1286 |   }
1287 | }
```

```

1286 // Function to submit feedback to the server
Windsurf: Refactor | Explain | X
1287 function submitFeedback(feedback, query, comment) {
1288     // If submitting with comment, show loading state in the submit button
1289     const submitBtn = document.getElementById('submitFeedbackBtn');
1290     if (submitBtn) {
1291         submitBtn.innerHTML =
1292             `<div class="flex items-center">
1293                 <svg class="animate-spin -ml-1 mr-2 h-4 w-4 text-white" xmlns="http://www.w3.org/2000/svg"
1294                 fill="none" viewBox="0 0 24 24">
1295                     <circle class="opacity-25" cx="12" cy="12" r="10" stroke="currentColor" stroke-width="4"/>
1296                     <circle>
1297                     <path class="opacity-75" fill="currentColor" d="M4 12a8 8 0 018-8V0C5.373 0 0 5.373 0 12h4zm2 5.
1298                         291A7.962 7.962 0 014 12H0c0 3.042 1.135 5.824 3 7.938l3-2.647z"></path>
1299                 </svg>
1300                 Submitting...
1301             </div>
1302         `;
1303         submitBtn.disabled = true;
1304     }
1305
1306     // Collect dataset information
1307     let datasetInfo = {};
1308     const selectedDataset = document.querySelector('.dataset-item.selected');
1309     if (selectedDataset) {
1310         const datasetId = selectedDataset.getAttribute('data-dataset-id');
1311         const datasetTitle = selectedDataset.querySelector('.dataset-title')?.textContent;
1312         const datasetYears = selectedDataset.querySelector('.dataset-years')?.textContent;
1313
1314         datasetInfo = {
1315             id: datasetId,
1316             title: datasetTitle,
1317             years: datasetYears
1318         };
1319
1320         // Add any dataset data attributes
1321         if (window.currentDatasetInfo) {
1322             datasetInfo.details = window.currentDatasetInfo;
1323         }
1324
1325     // Collect visualization information if available
1326     let visualizationData = {};
1327     if (window.lastVisualizationData) {
1328         visualizationData = window.lastVisualizationData;
1329     } else {

```

```

1342 // Send feedback to the backend
1343 fetch('/save_feedback', {
1344   method: 'POST',
1345   headers: {
1346     'Content-Type': 'application/json',
1347     'X-CSRFToken': getCsrfToken(),
1348   },
1349   body: JSON.stringify({
1350     query: query,
1351     feedback: feedback,
1352     comment: comment,
1353     dataset_info: datasetInfo,
1354     visualization_data: visualizationData,
1355     window_query_data: window.currentQueryData || {}
1356   })
1357 })
1358 .then(response => {
1359   if (!response.ok) {
1360     throw new Error(`Network response was not ok: ${response.status}`);
1361   }
1362   return response.json();
1363 })
1364 .then(data => {
1365   // Replace the entire feedback container with a thank you message
1366   const feedbackContainer = document.getElementById('feedbackContainer');
1367   if (feedbackContainer) {
1368     feedbackContainer.innerHTML =
1369       `<div class="text-green-700">
1370         <p class="font-medium">Thanks for your feedback!</p>
1371         <p class="text-sm">Your input helps us improve our search results.</p>
1372       </div>`;
1373   }
1374
1375   // Store visualization data for feedback
1376   const chartTypeButtons = document.querySelectorAll('.chart-type-button');
1377   const currentChartType = Array.from(chartTypeButtons).find(btn => btn.classList.contains('selected'))?.getAttribute('data-chart-type');
1378
1379   window.lastVisualizationData = {
1380     chart_type: visualizationData?.chart_type || currentChartType,
1381     request: document.getElementById('visualizationRequest')?.value || '',
1382     success: data.success,
1383     has_compatibility_warning: data.has_compatibility_warning || false
1384   };
1385 }
1386 )
1387 .catch(error => {
1388   console.error('Error submitting feedback:', error);
1389   // Show error message
1390   const feedbackContainer = document.getElementById('feedbackContainer');
1391   if (feedbackContainer) {
1392     feedbackContainer.innerHTML += `
1393       <div class="text-red-600 mt-2">
1394         <p>Failed to submit feedback. Please try again.</p>
1395       </div>`;
1396   }
1397 })

```

K: Caching

```
11  class QueryCache:      You, last month * implemented caching for queries, cahed being st...
12  """
13  A simple file-based cache for query results with expiration handling.
14  Caches query results in JSON format in the data/cache directory.
15  """
16
17  Windsurf: Refactor | Explain | X
18  def __init__(self):
19      """Initialize the query cache with the cache directory path."""
20      # Get base directory from settings or use default
21      base_dir = getattr(settings, 'BASE_DIR', None)
22      if base_dir:
23          self.cache_dir = Path(base_dir) / 'data' / 'cache'
24      else:
25          # Fallback to a relative path if BASE_DIR is not defined
26          self.cache_dir = Path(__file__).resolve().parent.parent.parent.parent / 'data' / 'cache'
27
28      # Ensure the cache directory exists
29      os.makedirs(self.cache_dir, exist_ok=True)
30      logger.info(f"Query cache initialized with directory: {self.cache_dir}")
31
32  Windsurf: Refactor | Explain | X
33  def _generate_cache_key(self, query):
34      """
35      Generate a unique cache key for a query.
36
37      """
38      # Create an MD5 hash of the query string to use as the filename
39      # This ensures unique filenames and handles special characters
40      hash_obj = hashlib.md5(query.encode('utf-8'))
41      return hash_obj.hexdigest()
42
43  Windsurf: Refactor | Explain | X
44  def _get_cache_file_path(self, cache_key):
45      """
46      Get the full path to a cache file.
47
48      """
49      return self.cache_dir / f"{cache_key}.json"
50
51  Windsurf: Refactor | Explain | X
52  def get(self, query):
53      """
54      Retrieve a cached query result if it exists and hasn't expired.
55
56      """
57      cache_key = self._generate_cache_key(query)
58      cache_file = self._get_cache_file_path(cache_key)
59
60      if not cache_file.exists():
61          logger.debug(f"Cache miss for query: {query[:50]}...")
62          return None
63
64      try:
65          with open(cache_file, 'r', encoding='utf-8') as f:
66              cached_data = json.load(f)
67
68              # Check if the cache has expired (30 days)
69              cached_time = datetime.datetime.fromisoformat(cached_data.get('cached_at', ''))
70              expiration = datetime.timedelta(days=30)
71
72              if datetime.datetime.now() - cached_time > expiration:
73                  logger.debug(f"Cache expired for query: {query[:50]}...")
74                  # Remove the expired cache file
75                  os.remove(cache_file)
76                  return None
77
78      logger.debug(f"Cache hit for query: {query[:50]}...")
```

```

75     # Update access count and last accessed time
76     cached_data['access_count'] = cached_data.get('access_count', 0) + 1
77     cached_data['last_accessed'] = datetime.datetime.now().isoformat()
78
79     # Write the updated metadata back to the cache
80     with open(cache_file, 'w', encoding='utf-8') as f:
81         json.dump(cached_data, f, ensure_ascii=False, indent=2)
82
83     return cached_data.get('result')
84
85 except (json.JSONDecodeError, KeyError, ValueError) as e:
86     logger.error(f"Error reading cache file {cache_file}: {str(e)}")
87     # Remove invalid cache file
88     if cache_file.exists():
89         os.remove(cache_file)
90     return None
91
92
93 Windsurf: Refactor | Explain | X
94 def set(self, query, result):
95     """
96     Cache a query result.
97
98     """
99     cache_key = self._generate_cache_key(query)
100    cache_file = self._get_cache_file_path(cache_key)
101
102    try:
103        # Create cache data structure with metadata
104        cache_data = {
105            'query': query,
106            'result': result,
107            'cached_at': datetime.datetime.now().isoformat(),
108            'access_count': 0,
109            'last_accessed': None
110        }
111
112        # Write to cache file
113        with open(cache_file, 'w', encoding='utf-8') as f:
114            json.dump(cache_data, f, ensure_ascii=False, indent=2)
115
116        logger.debug(f"Cached result for query: {query[:50]}...")
117        return True
118
119    except Exception as e:
120        logger.error(f"Error caching query result: {str(e)}")
121        return False
122
123 Windsurf: Refactor | Explain | X
124 def clear(self, query=None):
125     """
126     Clear specific cache entry or all cache entries.
127
128     """
129     if query:
130         # Clear specific cache entry
131         cache_key = self._generate_cache_key(query)

```

```

131     cache_file = self._get_cache_file_path(cache_key)
132
133     if cache_file.exists():
134         os.remove(cache_file)
135         logger.info(f"Cleared cache for query: {query[:50]}...")
136         return 1
137     return 0
138 else:
139     # Clear all cache entries
140     count = 0
141     for cache_file in self.cache_dir.glob('*json'):
142         os.remove(cache_file)
143         count += 1
144
145     logger.info(f"Cleared {count} entries from the query cache")
146     return count
147
148 Windsurf: Refactor | Explain | X
149 def get_cache_stats(self):
150     """
151     Get statistics about the current cache.
152     """
153     stats = {
154         'total_entries': 0,
155         'total_size_bytes': 0,
156         'oldest_entry': None,
157         'newest_entry': None,
158         'most_accessed': None,
159         'most_accessed_count': 0
160     }
161
162     # Get all cache files
163     cache_files = list(self.cache_dir.glob('*json'))
164     stats['total_entries'] = len(cache_files)
165
166     oldest_time = None
167     newest_time = None
168
169     for cache_file in cache_files:
170         # Get file size
171         file_size = os.path.getsize(cache_file)
172         stats['total_size_bytes'] += file_size
173
174         try:
175             with open(cache_file, 'r', encoding='utf-8') as f:
176                 cache_data = json.load(f)
177
178                 cached_time = datetime.datetime.fromisoformat(cache_data.get('cached_at', ''))
179
180                 # Track oldest and newest entries
181                 if oldest_time is None or cached_time < oldest_time:
182                     oldest_time = cached_time
183                     stats['oldest_entry'] = {
184                         'query': cache_data.get('query', '')[:50] + '...',
185                         'cached_at': cached_time.isoformat()
186                     }
187
188                 if newest_time is None or cached_time > newest_time:

```

```
189         newest_time = cached_time
190         stats['newest_entry'] = {
191             'query': cache_data.get('query', '')[:50] + '....',
192             'cached_at': cached_time.isoformat()
193         }
194
195     # Track most accessed entry
196     access_count = cache_data.get('access_count', 0)
197     if access_count > stats['most_accessed_count']:
198         stats['most_accessed_count'] = access_count
199         stats['most_accessed'] = {
200             'query': cache_data.get('query', '')[:50] + '....',
201             'access_count': access_count
202         }
203
204     except (json.JSONDecodeError, KeyError) as e:
205         logger.warning(f"Error reading cache file {cache_file}: {str(e)}")
206
207     # Convert total size to a more readable format
208     if stats['total_size_bytes'] > 1024 * 1024:
209         stats['total_size'] = f'{stats["total_size_bytes"] / (1024 * 1024):.2f} MB'
210     else:
211         stats['total_size'] = f'{stats["total_size_bytes"] / 1024:.2f} KB'
212
213 return stats
```

L: Gemini Client

```
26  class GeminiClient:      You, 2 months ago • created indeing file for the AI approch. fix th...
27      """Client for interacting with Google's Gemini API for natural language processing."""
28
29      Windsurf: Refactor | Explain | X
30      def __init__(self, api_key: Optional[str] = None):
31          """Initialize the Gemini client with API key."""
32          # Try to get API key from environment variable if not provided
33          self.api_key = api_key or os.environ.get('GEMINI_API_KEY')
34          if not self.api_key:
35              logger.warning("No Gemini API key provided. API calls will not work.")
36
37      Windsurf: Refactor | Explain | X
38      def analyze_query(self, query: str) -> Dict[str, Any]:
39          """
40              Analyze a natural language query using the Gemini API.
41
42          """
43          import requests
44          import json
45          import re
46          import logging
47
48          logger = logging.getLogger(__name__)
49          logger.info(f"Analyzing query with Gemini API: {query}")
50
51          if not self.api_key:
52              logger.error("No Gemini API key provided")
53              raise ValueError("Gemini API key is required")
54
55          # Construct the prompt
56          system_prompt = """
57          You are a helpful assistant that extracts structured information from a user's query about higher education
58          statistics.
59
60          Extract the following information and output it in JSON format:
61          1. Institutions mentioned (e.g., "University of Oxford", "University of Leicester")
62          2. Years mentioned (e.g., "2020", "2020/21")
63          3. If a year range is given, identify the start_year and end_year
64          4. The type of data being requested (e.g., "student numbers", "enrollment", "graduates")
65
66          IMPORTANT: You need to detect and correct ONLY SPELLING typos in institution names and years.
67
68          For institutions, follow these strict rules:
69          - ONLY consider a term to have a typo if it has clear misspellings (e.g., "Univercity" → "University", "Lieicester" →
70          "Leicester")
71          - DO NOT consider "Oxford University" vs "University of Oxford" as a typo - these are different ways to refer to the
72          same institution
73          - DO NOT mark an institution as having a typo if it's spelled correctly but uses a different naming convention
74          - DO NOT transform a city name like "london" into "The University of London"
75          - PRESERVE the original terms from the query - if user says "london", keep it as "london" in the institutions list
76          - If there are no spelling errors, original_institutions should exactly match institutions, and has_institution_typos
77          should be false      You, last month • important stage ...
78
79          For years:
80          - ONLY mark years as having typos if they have clear numerical errors (e.g., "20025" → "2025")
81          - If no year typos exist, original_years should exactly match years, and has_year_typos should be false
```

```

77     Be careful about interpreting years in academic context:
78     - If the query contains phrases like "starting in [YEAR]" or "beginning in [YEAR]", interpret [YEAR] as the start of
79     an academic year.
80     - If the query contains phrases like "end of [YEAR]" or "ending in [YEAR]", interpret [YEAR] as the end of an
81     academic year.
82     - For ranges like "2016 to 2017", treat both as starting years of academic years.
83     - If there's no clarification, assume a year refers to the starting year of an academic year.
84     - For "past X years", calculate the years based on the current year (2025).
85
86     For the data_request field, use the SPECIFIC TERMS from the query whenever possible:
87     - If the query mentions "undergraduates", use ["undergraduate"] not ["student_enrollment"]
88     - If the query mentions "postgraduates", use ["postgraduate"] not ["student_enrollment"]
89     - If the query mentions "student numbers", use ["student_count"]
90     - If the query mentions "enrollment" or "enrolment", use ["enrollment"]
91     - If the query mentions "staff" or "teachers", use ["staff_data"]
92     - If the query mentions "research", use ["research_data"]
93     - NEVER generalize specific educational levels - keep them exactly as requested
94
95     Output format:
96     {
97         "institutions": ["University X", "london"],
98         "original_institutions": ["Universcity X", "london"],
99         "has_institution_typos": true,
100        "years": ["2019/20", "2020/21"],
101        "original_years": ["20019", "2020"],
102        "has_year_typos": true,
103        "start_year": "2019",
104        "end_year": "2020",
105        "data_request": ["undergraduate", "graduation_rates"]
106
107    If no specific institutions are mentioned, return an empty list for institutions and original_institutions.
108    If no typos were detected, original_institutions should match institutions, and has_institution_typos should be false.
109    The same applies to years and original_years.
110    """
111
112    # The actual user query
113    user_prompt = f"Extract information from this query, paying special attention to academic year conventions and
114    potential typos in institution names and years: '{query}'"
115
116    try:
117        # Gemini API endpoint - Updated to use the most current endpoint
118        url = "https://generativelanguage.googleapis.com/v1/models/gemini-1.5-pro:generateContent"
119
120        # Request payload
121        payload = {
122            "contents": [
123                {
124                    "role": "user",
125                    "parts": [
126                        {"text": system_prompt},
127                        {"text": user_prompt}
128                    ]
129                },
130                "generationConfig": {
131                    "temperature": 0.1,
132                    "topP": 0.95,
133

```

```

132         "topK": 40,
133         "maxOutputTokens": 2048
134     }
135 }
136
137     # Headers with API key
138     headers = {
139         "Content-Type": "application/json",
140         "x-goog-api-key": self.api_key
141     }
142
143     # Make the request
144     response = requests.post(url, json=payload, headers=headers)
145
146     # Check for successful response
147     if response.status_code != 200:
148         logger.error(f"Gemini API error: {response.status_code} - {response.text}")
149         raise Exception(f"Gemini API returned status code {response.status_code}: {response.text}")
150
151     # Parse the response.
152     response_data = response.json()
153
154
155     def _process_academic_year_logic(self, query, result):
156         """
157             Process academic year logic based on context clues in the query.
158             This handles cases like:
159             - "starting in 2017" → 2017/18
160             - "end of 2017" → 2016/17
161             - "2016 to 2017" → 2016/17 - 2017/18 (both as starting years)
162
163         """
164
165         import re
166         import logging
167
168         logger = logging.getLogger(__name__)
169         query_lower = query.lower()
170
171         # Extract years with context
172         start_patterns = [
173             r'start(?:ing|s|ed)?\s*(?:in|from|at)?\s*(?:the\s+)?(?:year\s+)?(\d{4})',
174             r'begin(?:ning|s)?\s*(?:in|from|at)?\s*(?:the\s+)?(?:year\s+)?(\d{4})',
175             r'from\s*(?:the\s+)?(?:year\s+)?(\d{4})'
176         ]
177
178         end_patterns = [
179             r'end(?:ing|s|ed)?\s*(?:in|at|of)?\s*(?:the\s+)?(?:year\s+)?(\d{4})',
180             r'finish(?:ing|es|ed)?\s*(?:in|at)?\s*(?:the\s+)?(?:year\s+)?(\d{4})',
181             r'^(?:in|at)\s*(?:the\s+)?end\s*of\s*(?:the\s+)?(?:year\s+)?(\d{4})'
182         ]
183
184
185         # Handle range logic for years without context
186         if 'start_year' in result and 'end_year' in result and result['start_year'] is not None and result['end_year'] is not
187             None:
188             # Range years should be treated as starting years of academic years
189             start_year = result['start_year']
190             end_year = result['end_year']
191
192             try:
193                 # Clear existing years array to rebuild it with correct academic years
194                 result['years'] = []
195
196                 # Add academic years for the range
197                 for year in range(int(start_year), int(end_year) + 1):
198                     academic_year = f'{year}/{str(year+1)[2:4]}'
199                     result['years'].append(academic_year)
200
201
202                 logger.info(f"Processed year range: {start_year}-{end_year} as academic years: {', '.join(result['years'])}")
203             except (ValueError, TypeError) as e:
204                 logger.error(f"Error processing year range: {e}. start_year={start_year}, end_year={end_year}")
205
206             return
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282

```

```

351     def get_completion(self, prompt):
352         """
353             Get a completion from the Gemini API based on the provided prompt.
354
355         """
356         import requests
357         import logging
358
359         logger = logging.getLogger(__name__)
360         logger.info("Sending completion request to Gemini API")
361
362         if not self.api_key:
363             logger.error("No Gemini API key provided")
364             raise ValueError("Gemini API key is required")
365
366         try:
367             # Updated Gemini API endpoint
368             url = "https://generativelanguage.googleapis.com/v1/models/gemini-1.5-pro:generateContent"
369
370             # Request payload
371             payload = {
372                 "contents": [
373                     {
374                         "role": "user",
375                         "parts": [
376                             {"text": prompt}
377                         ]
378                     }
379                 ],
380                 "generationConfig": {
381                     "temperature": 0.2,
382                     "topP": 0.95,
383                     "topK": 40,
384                     "maxOutputTokens": 4096
385                 }
386             }
387
388             # Headers with API key
389             headers = {
390                 "Content-Type": "application/json",
391                 "x-goog-api-key": self.api_key
392             }
393
394             # Make the request
395             response = requests.post(url, json=payload, headers=headers)
396
397             # Check for successful response
398             if response.status_code != 200:
399                 logger.error(f"Gemini API error: {response.status_code} - {response.text}")
400                 raise Exception(f"Gemini API returned status code {response.status_code}: {response.text}")
401
402             # Parse the response
403             response_data = response.json()
404

```

```
405     # Extract text from response
406     if 'candidates' in response_data and len(response_data['candidates']) > 0:
407         if 'content' in response_data['candidates'][0]:
408             response_text = response_data['candidates'][0]['content']['parts'][0]['text']
409             return response_text
410         else:
411             logger.error("No content in Gemini API response")
412             raise Exception("No content in Gemini API response")
413     else:
414         logger.error("No candidates in Gemini API response")
415         raise Exception("No candidates in Gemini API response")
416
417     except Exception as e:
418         logger.error(f"Error getting completion from Gemini API: {str(e)}")
419         raise
420
421 Windsurf: Refactor | Explain | X
422 def generate_text(self, prompt):
423     """
424     Generate text using the Gemini API - alias for get_completion for compatibility
425     """
426     return self.get_completion(prompt)
427
```

M: Timing Queries

M.1 Recording the time of query processing

```
5335     @csrf_exempt
5336     @require_http_methods(["POST"])
5337     def record_query_timing(request):      You, 6 days ago * uploading all the cleaning csv files nad final ...
5338         """
5339         Record timing information for query processing.
5340         Data is stored in a JSON file at the root level.
5341         """
5342         import json
5343         import os
5344         import logging
5345         import traceback
5346         from pathlib import Path
5347         from django.conf import settings
5348
5349         logger = logging.getLogger(__name__)
5350         logger.info("== RECORDING QUERY TIMING ==")
5351
5352     try:
5353         # Parse the data from the request
5354         data = json.loads(request.body.decode('utf-8'))
5355         query = data.get('query', '')
5356         duration_ms = data.get('duration_ms', 0)
5357         # Store the requested match count, not the returned match count
5358         requested_match_count = data.get('requested_match_count', 0)
5359         mission_group = data.get('mission_group')
5360
5361         # Validate incoming data
5362         if not query:
5363             logger.warning("Empty query received in timing data")
5364             return JsonResponse({'status': 'error', 'error': 'Empty query'}, status=400)
5365
5366         if not isinstance(duration_ms, (int, float)) or duration_ms <= 0:
5367             logger.warning(f"Invalid duration value: {duration_ms}")
5368             return JsonResponse({'status': 'error', 'error': 'Invalid duration'}, status=400)
5369
5370         # Convert to integer milliseconds
5371         duration_ms_int = int(duration_ms)
5372
5373         # Format duration in seconds with 2 decimal places
5374         duration_sec = round(duration_ms / 1000, 2)
5375
5376         # Define the path to the timing data file - ensure it's in the project root
5377         timing_file_path = Path(settings.BASE_DIR) / 'timing_queries.json'
5378
5379         logger.info(f"Using timing file path: {timing_file_path}")
5380
5381         # Create a unique identifier for this query
5382         # A query is unique if it has the same text, mission group, and match count
5383         query_id = f"{query}|{mission_group or 'none'}|{requested_match_count}"
5384
5385         # Initialize with empty data if file doesn't exist
5386         if not os.path.exists(timing_file_path):
5387             logger.info(f"Creating new timing file at {timing_file_path}")
5388             with open(timing_file_path, 'w') as f:
```

```

5388     with open(timing_file_path, 'w') as f:
5389         json.dump({"queries": []}, f, indent=2)
5390
5391     # Read existing data with error handling
5392     try:
5393         with open(timing_file_path, 'r') as f:
5394             timing_data = json.load(f)
5395
5396             # Validate the structure
5397             if not isinstance(timing_data, dict) or 'queries' not in timing_data:
5398                 logger.warning("Invalid timing data structure, resetting file")
5399                 timing_data = {"queries": []}
5400             except (json.JSONDecodeError, FileNotFoundError) as e:
5401                 logger.error(f"Error reading timing file: {str(e)}")
5402             # Reset data if file is corrupted
5403             timing_data = {"queries": []}
5404
5405             # Find if this query already exists
5406             existing_entry = None
5407             for entry in timing_data.get('queries', []):
5408                 if (entry.get('query') == query and
5409                     entry.get('mission_group') == mission_group and
5410                     entry.get('match_count') == requested_match_count):
5411                     existing_entry = entry
5412                     break
5413
5414             # Update existing entry or add new one
5415             if existing_entry:
5416                 # Add timing to the existing entry
5417                 if 'durations_ms' not in existing_entry:
5418                     existing_entry['durations_ms'] = []
5419                 if 'durations_sec' not in existing_entry:
5420                     existing_entry['durations_sec'] = []
5421
5422                 existing_entry['durations_ms'].append(duration_ms_int)
5423                 existing_entry['durations_sec'].append(duration_sec)
5424             else:
5425                 # Create a new entry with a sequential ID
5426                 new_id = len(timing_data.get('queries', [])) + 1
5427                 new_entry = {
5428                     'id': new_id,
5429                     'query': query,
5430                     'mission_group': mission_group,
5431                     'match_count': requested_match_count,
5432                     'durations_ms': [duration_ms_int],
5433                     'durations_sec': [duration_sec]
5434                 }
5435                 timing_data['queries'].append(new_entry)
5436
5437             # Write the updated data back to the file
5438             try:
5439                 with open(timing_file_path, 'w') as f:
5440                     json.dump(timing_data, f, indent=2)
5441
5442                     logger.info(f"Recorded timing for query: '{query}'")
5443                     logger.info(f"Duration: {duration_ms_int} ms ({duration_sec} sec)")
5444             except Exception as write_error:
5445                 logger.error(f"Error writing timing data: {str(write_error)}")
5446                 return JsonResponse({'status': 'error', 'error': f'File write error: {str(write_error)}'}, status=500)
5447
5448             return JsonResponse({'status': 'success'})
5449
5450         except json.JSONDecodeError as e:
5451             logger.error(f"JSON decode error in timing data: {str(e)}")
5452             return JsonResponse({'status': 'error', 'error': f'Invalid JSON: {str(e)}'}, status=400)
5453         except Exception as e:
5454             logger.error(f"Error recording query timing: {str(e)}")
5455             logger.error(traceback.format_exc())
5456             return JsonResponse({'status': 'error', 'error': str(e)}, status=500)

```

M.2 Recorded time of processed queries

```
1  {
2    "queries": [
3      {
4        "id": 1,
5        "query": "How many postgraduates are in University of Cambridge in 2017?",
6        "mission_group": null,
7        "match_count": 2,
8        "durations_ms": [
9          6212,
10         59,
11         355
12       You, 6 days ago • uploading all the cleaning csv files nad final ...
13     ],
14     "durations_sec": [
15       6.21,
16       0.06,
17       0.36
18     ]
19   },
20   {
21     "id": 2,
22     "query": "What was the research income starting in 2017?",
23     "mission_group": null,
24     "match_count": 2,
25     "durations_ms": [
26       5482,
27       47
28     ],
29     "durations_sec": [
30       5.48,
31       0.05
32     ]
33   },
34   {
35     "id": 3,
36     "query": "Show carbon emissions for Coventry University ending in 2021",
37     "mission_group": null,
38     "match_count": 1,
39     "durations_ms": [
40       13173
41     ],
42     "durations_sec": [
43       13.17
44     ]
45   },
46   {
47     "id": 4,
48     "query": "Show carbon emissions for Coventry University in 2021",
49     "mission_group": null,
50     "match_count": 1,
51     "durations_ms": [
52       4292,
53       106
54     ],
55     "durations_sec": [
56       4.29,
57       0.11
58     ]
59   }
60 }
```

```
58 },
59 {
60   "id": 5,
61   "query": "How many full-time students lived in university accommodation at University of York from 2016 to 2018?",
62   "mission_group": null,
63   "match_count": 4,
64   "durations_ms": [
65     8665,
66     98
67   ],
68   "durations_sec": [
69     8.67,
70     0.1
71   ]
72 },
73 {
74   "id": 6,
75   "query": "Research funding for University of Manchester from 2017 to 2019",
76   "mission_group": null,
77   "match_count": 3,
78   "durations_ms": [
79     10309,
80     82
81   ],
82   "durations_sec": [
83     10.31,
84     0.08
85   ]
86 },
87 {
88   "id": 7,
89   "query": "Graduation rates at London institutions for the past 5 years",
90   "mission_group": null,
91   "match_count": 5,
92   "durations_ms": [
93     12935
94   ],
95   "durations_sec": [
96     12.94
97   ]
98 },
99 {
100   "id": 8,
101   "query": "Graduation rates at University of London institutions for the past 5 years",
102   "mission_group": null,
103   "match_count": 5,
104   "durations_ms": [
105     58301,
106     665,
107     181
108   ],
109   "durations_sec": [
110     58.3,
111     0.67,
112     0.18
113   ]
114 },
```