

# Uni-Mate: A Retrieval-Augmented Generation System to Provide High School Students with Accurate Academic Guidance

Samuele Mazzei<sup>1</sup>, Lorenzo Zambotto<sup>1</sup>, Gabriele Tealdo<sup>1</sup>, Alberto Macagno<sup>1</sup> and Alessio Palmero Aprosio<sup>1,\*</sup>

<sup>1</sup>Department of Psychology and Cognitive Science, University of Trento, Corso Bettini 84, Rovereto, Italy

## Abstract

This paper introduces the development and evaluation of a Retrieval-Augmented Generation (RAG) system designed to assist prospective students in navigating university options. The system provides accurate academic guidance by retrieving and synthesizing information on undergraduate and single-cycle master's degree programs, as well as library resources, from the University of Trento and the University of Verona. The RAG pipeline utilizes a streamlined toolchain, incorporating a Markdown parser for efficient data handling and the Llama3-8b-8192 Large Language Model (LLM) for query processing. The system's performance was assessed through both automated evaluation, using the Llama3-70b LLM as a reference, and blinded human evaluation. The results demonstrate the system's potential for providing relevant and accurate information to students. The evaluation also highlighted areas for further development, including enhanced retrieval mechanisms and expanded LLM testing. Future work aims to broaden the system's scope to include more degree levels and universities, ultimately creating a comprehensive platform to support students in their academic decision-making journey.

## Keywords

Retrieval-Augmented Generation (RAG), Natural Language Processing (NLP), Large Language Models (LLMs), Dataset Creation, Academic Guidance

## 1. Introduction

Choosing a university path is one of the most complex and significant decisions for students nearing the end of high school. This, combined with the overwhelming amount of new information encountered when browsing various and often inconsistent university websites, creates confusion and a sense of being lost, leading to wasted time and uncertainty. These challenges stem from both the dispersion of available information and the lack of intuitive tools to guide students through the decision-making process.

We deal with this problem by creating a platform called Uni-Mate (formerly referred to as *MyVision* and later renamed to better align with startup branding goals, offering a more appealing name for potential users and investors). The system aims to integrate an AI-powered chatbot that provides relevant information about partner universities and online counseling services within a single interface.

A survey, conducted among 183 students from the Department of Psychology and Cognitive Science and the School of Innovation between October and November 2024, was instrumental in identifying a significant need among students for improved online educational guidance and revealed significant challenges faced by students in choosing their academic paths. A striking 74% reported at least one major difficulty in the orientation process. The most common issues included a lack of clear and comparable information across courses and institutions (43%), uncertainty regarding personal interests and aptitudes (38%), and confusion about the differences among European universities (29%). Additionally, limited access to insights from alumni was also noted (17%). When seeking guidance, students primarily relied on official university websites (65%) and personal networks such as parents or friends (58%), while only 21% consulted academic counselors. Moreover, fewer than 10% found digital comparison tools to be truly effective.

The data also highlights a strong interest in innovative orientation tools. Notably, 81% of respondents expressed a willingness to use a platform like Uni-Mate, which would feature personalized course matching algorithms and structured reviews from former students. Furthermore, 67% indicated a readiness to pay for such a service if it proved to be effective. These results point to a clear gap in the current academic orientation offerings, which are seen as fragmented, non-interactive, and lacking personalization. There is a strong latent demand for com-

*CLiC-it 2025: Eleventh Italian Conference on Computational Linguistics, September 24 – 26, 2025, Cagliari, Italy*

\*Corresponding author.

<sup>†</sup>These authors contributed equally.

✉ samuele.mazzei@studenti.unitn.it (S. Mazzei);  
lorenzo.zambotto@studenti.unitn.it (L. Zambotto);  
gabriele.tealdo@studenti.unitn.it (G. Tealdo);  
alberto.macagno@studenti.unitn.it (A. Macagno);  
a.palmeroaprosio@unitn.it (A. Palmero Aprosio)



© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

prehensive digital solutions that provide personalized guidance, real-life experiences, and comparative tools to support students in making well-informed educational decisions.

Following this initial validation, the team submitted MyVision as a proposal to DigiEduHack,<sup>1</sup> a European innovation challenge promoted by the European Commission and aimed at fostering technological advancements in the field of education. During the event, the concept was further developed and ultimately awarded first place in the Expert Category, the most competitive and high-level track of the competition.<sup>2</sup>

As a result, the team has been invited to present MyVision at the DigiEduHack Award Ceremony, scheduled for June 24th in Brussels, as part of the Digital Education Stakeholder Forum 2025 — a major annual event organized by the European Commission to promote dialogue and policy development in digital education.

Notable precedents attempts to address these challenges exist in the academic guidance space. In the United States, ScholarMatch<sup>3</sup> focuses on helping first-generation, low-income students secure scholarships and complete their college education, addressing critical financial and support gaps. In contrast, a comparable all-in-one solution is lacking in Europe, where the challenges for students are less about tuition affordability and more about navigating a fragmented ecosystem of academic options. In the UK, Bonas MacFarlane<sup>4</sup> offers premium consulting for school and university placements, primarily targeting affluent families. These examples highlight both the proven demand for personalized academic support and the gap that MyVision seeks to fill in the EU context—by offering accessible, digital tools for orientation, comparison, reviews, and guidance all in one unified platform.

In Italy, UniverItaly<sup>5</sup> is the official portal developed by the Italian Ministry of University and Research to support students, both Italian and international, in navigating the higher education system in Italy. The web portal integrates a conversational assistant powered by large language models, which helps users navigate content and find relevant information interactively.

In this paper, we aim to lay the foundation for the development of our chatbot by focusing on the academic offerings of two Italian universities: the University of Trento (Unitn) and the University of Verona (Univr). These institutions were selected due to their geographical proximity and the presence of interdisciplinary and interuniversity courses, which offer significant opportunities for prospective students interested in studying in these ar-

eas. The system was developed with a specific focus on post-diploma university orientation, considering only bachelor's degree programs and single-cycle master's degrees. This approach addresses the needs of recent high school graduates by providing an innovative tool to explore available academic options in a simple and immediate way. Furthermore, we included information about the universities' libraries to provide new students with access to a valuable resource that can support their studies.

## 2. Related work

Numerous research groups and institutions have explored various strategies to support students in selecting the most suitable university.

For example, in a study [1], the researchers evaluated an educational app called GC Mobile and concluded that it enhanced the counseling process by leveraging technology to provide a scalable, accessible, and confidential platform for student guidance. As the authors noted, "The GC Mobile App allows students to see a counselor anytime and from any location without having to visit them in the office."

Another study [2] developed an AI-powered academic guidance and counseling system with the primary objective of supporting high school seniors in navigating the college application process and selecting suitable academic paths and universities for tertiary education. It also aimed to address the shortage of human resources in traditional counseling by providing an accessible, convenient, and time-saving alternative for students to obtain valuable insights without requiring face-to-face interaction or travel to gather university information.

Another approach was the creation of UniCompass [3], a platform designed to help students efficiently learn about and compare universities and departments, and to access diverse perspectives and shared experiences from peers. By consolidating information and providing structured guidance, UniCompass aims to save students time and support more informed academic and career decisions.

A similar application is "Major-Selection" [4], which functions as intelligent decision support software to assist students with major selection. It features a rule-based knowledge base containing information about university admission requirements and the skills and preferences relevant to various majors. This knowledge is derived from academic advisors and university guidelines.

In another study [5], the authors developed a web application that provides personalized recommendations and guidance to high school students. By using a questionnaire, the AI system builds a comprehensive profile of the student and delivers data-driven, customized guid-

<sup>1</sup><https://digieduhack.com/>

<sup>2</sup><https://digieduhack.com/news/digieduhack-2024-winners-announced-meet-the-innovators-shaping-digital-education>

<sup>3</sup><https://www.scholarmatch.org/>

<sup>4</sup><https://bonasmacfarlane.co.uk/en>

<sup>5</sup><https://www.universitaly.it/orientamento-universitario>

ance to support informed university and career decisions.

Lastly, myAlmaOrienta [6] was developed to support high school students in choosing a degree programme at the University of Bologna. It helps students navigate the selection process and identify programmes that match their skills and interests. The app was developed through a two-level co-design process involving both high school students (user-driven innovation) and university students (open innovation contest) to incorporate their needs and perspectives.

The chatbot involved in Uni-Mate uses Retrieval-Augmented Generation (RAG) to address the limitations inherent in traditional methods and standalone Large Language Models (LLMs) [7], such as limited context and possible hallucinations. Dieing et al. [8] describes a system for study program orientation that provides personalized recommendations using a Mixtral LLM paired with a RoBERTa embedding model. Their RAG approach retrieves data from a government website and achieves an average response accuracy above 0.75. Saha and Saha [9] reports that a GPT-3.5-based chatbot enhances support for international graduate students by combining generative capabilities with precise retrieval from social media sources. Dakshit [10] explored the use of RAG in higher education, focusing on applications as virtual teaching assistants and teaching aids. Faculty perspectives gathered in the study highlighted the benefits of RAG in supporting teaching processes, such as the generation of study guides, quizzes, and assignment questions, while also assisting students by providing precise answers to academic queries. Faculty members emphasized the importance of integrating broader data sources and advanced functionalities, including the ability to process mathematical content and image-based inputs, to improve the system's effectiveness.

The potential of RAG-powered systems lies in their ability to provide accurate, contextually relevant, and personalized support by combining retrieval mechanisms with generation capabilities [7]. A retrieval component first searches for relevant information from a curated set of academic resources, ensuring the content is accurate and domain-specific. The generation component then synthesizes this information to produce coherent and contextually appropriate responses [11]. This dual approach not only improves the reliability of responses but also enables the system to adapt to individual learning styles and paces, making it a valuable tool for personalized education. These findings align with the goals of Uni-Mate, particularly in creating a chatbot that integrates multiple functions—academic guidance, counseling services, and information retrieval—into a cohesive platform. Drawing from the studies mentioned above, we plan to leverage RAG's strengths to ensure that Uni-Mate not only meets students' informational needs but also provides reliable, context-aware responses to enhance their educational

journey.

### 3. Dataset

To collect the documents for our task, we accessed the course websites of Unitn<sup>6</sup> and Univr<sup>7</sup> to gather the necessary data. Since the main objective of this project is to provide orientation for high school students, we selected undergraduate degrees and single cycle master's degrees. For Unitn, we obtained data from the "Prospective Student" section, which is divided into three parts: "Course Programme," providing an overview of the degree; "Course Content," listing all courses offered over the years along with their respective ECTS credits, and in some cases, detailed course descriptions; and "Application," which contains enrollment information. For Univr, we collected similar information. After selecting a degree, we retrieved the "Overview" section under the "Find out more" option, the study plan from the "Modules" section, and enrollment details from the "How to apply" option. All collected data of the courses was converted into Markdown format with the help of an extension of ChatGPT-4 called Markdown converter<sup>8</sup>. ChatGPT-4 does not always structure the data in the same way, so we manually adjusted the formatting when discrepancies were too large. We also collected data on the libraries of both universities. In this case, the data were gathered manually to ensure a consistent file structure and order. The collected library data included: a general overview, with information on access, location, staff, and available spaces; the services offered by the libraries; and the opening hours.

We used Markdown language for several reasons, including efficiency and flexibility. This format allows for a clear structuring of data through the use of headings, enabling the RAG to subsequently divide the information into well-defined and interconnected sections. This optimization facilitates the retrieval process, making it easier to identify and associate relevant information. Another advantage of Markdown is its ability to include tables, which are clearer and more understandable as responses for users. Finally, the Markdown format is more practical during the dataset creation phase, as it allows for the use of tools like scrapers to quickly extract text from web pages. This process simplifies and accelerates the assembly of necessary information while ensuring greater consistency and quality of the data. In total, we collected data for 29 degrees from Unitn and 41 degrees from Univr, resulting in 70 course documents. Additionally, we collected data from 5 libraries from Unitn and 34 libraries from Univr, resulting in 39 library documents. This yielded a total of 109 documents.

<sup>6</sup><https://www.unitn.it/en/ateneo/1819/programmes-of-study>

<sup>7</sup><https://www.univr.it/en/degree-programmes>

<sup>8</sup><https://chatgpt.com/g/g-lnlmekbGd-markdown-converter>

Additionally, all data were translated into English when the English version of the site did not contain sufficient information compared to its Italian counterpart, as the answers provided by our RAG system were more accurate due to the embedding model introduced during the course. The English version of the embedding model is trained and tested on more data and has access to a larger corpus than the Italian version, which typically results in better training, improved generalization, and richer language representations [12]. To verify this, we consulted the literature and found a paper titled “Retrieval-augmented generation in multilingual settings” [13], which confirms our hypothesis.

## 4. Experiments

The objective of this study was to develop and evaluate a document retrieval system designed to query information from university course descriptions and library details. The system’s performance was assessed based on its accuracy in retrieving relevant and contextually appropriate information. For this purpose, we utilized Groq<sup>9</sup> as the provider for Large Language Models (LLMs). Specifically, two models were employed: Llama3-8b-8192 (8 billion parameters) served as the primary LLM for query processing, while Llama3-70b (70 billion parameters) functioned as the reference (“golden”) model during evaluation.

### 4.1. RAG Pipeline

The experimental workflow starts with a corpus of structured Markdown documents, detailing university courses and library information (as described in Section 3). The documents are loaded manually into the system from the two separated folders for courses and libraries. For each file we then create a LlamaIndex Document object by adding metadata to it, extracting information from the file title. Specifically for the courses we extract the university name, in its shorter form, and the course name, eventually translated in English and dash separated. For the libraries we extract the university name and the name of the library, following the same convention. Because the single documents are considerably long, we decided to split them in smaller chunks to have more meaningful embeddings. Among the different strategies available, our ultimate choice for processing documents relied on a specific node parser: `MarkdownNodeParser`<sup>10</sup>. This is a class provided by LlamaIndex that splits the documents into Nodes following a Markdown splitting

logic, by separating the sources using headings. Moreover, through the use of the `include_prev_next_rel` and `include_metadata` parameters, we keep relationships between the nodes, supporting the retrieval process. Nodes are persisted in a local document store in a Google Drive folder.

Subsequently, these nodes are converted into vector embeddings. As for the model of embedding, we chose the BAAI/bge-m3 model<sup>11</sup> which distinguished itself especially for its multi-granularity and the ability to work with long documents in generating semantic representations of the text. The model is loaded using the `HuggingFaceEmbedding`<sup>12</sup> module of LlamaIndex, which provides a convenient interface for working with Hugging Face models. The embeddings are generated using the GPU acceleration provided by a T4 instance in Google Colab<sup>13</sup>, which significantly speeds up the embedding generation process, and are saved in a cache folder on Google Drive to avoid redundant computations in development.

The retrieval is performed using the BM25 algorithm<sup>14</sup>, a widely used keyword-based retrieval method that employs lexical matching to retrieve relevant document sections. The BM25 algorithm is implemented in LlamaIndex and is used to retrieve the top 15-k nodes based on the similarity with the user query.

A graphical representation of the whole pipeline is shown in Figure 1.

### 4.2. Evaluation

Evaluation of the system’s performance employed a dual approach: automated assessment using the Llama3-70b model and blinded human evaluation, ensuring objectivity. Both methods assessed the quality of the generated answers and, for the automated part, the suitability of the retrieved context.

For the automated evaluation of generated answers, the Llama3-70b model assessed relevance and correctness relative to the user query. It assigned a score on a 1-to-5 scale, which was subsequently normalized to a 0-to-4 scale for direct comparison with human scores. The model also generated a textual justification explaining its assessment, highlighting aspects like completeness or accuracy. Due to API call limitations with standard evaluation frameworks, custom requests were implemented to facilitate this automated assessment process.

Automated context assessment focused on the text passages retrieved by the BM25 algorithm before answer

<sup>9</sup><https://groq.com/>

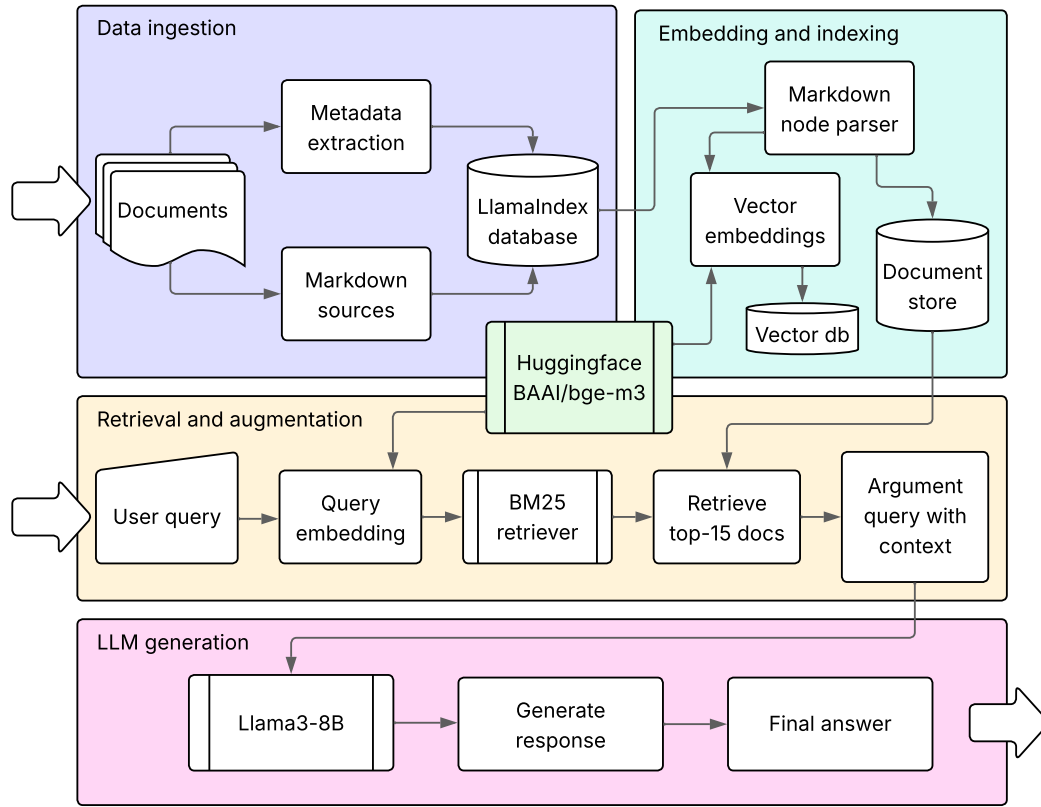
<sup>10</sup>[https://docs.llamaindex.ai/en/v0.10.17/api/llama\\_index.core.node\\_parser.MarkdownNodeParser.html](https://docs.llamaindex.ai/en/v0.10.17/api/llama_index.core.node_parser.MarkdownNodeParser.html)

<sup>11</sup><https://huggingface.co/BAAI/bge-m3>

<sup>12</sup><https://docs.llamaindex.ai/en/stable/examples/embeddings/huggingface/>

<sup>13</sup><https://colab.research.google.com/>

<sup>14</sup>[https://docs.llamaindex.ai/en/stable/examples/retrievers/bm25\\_retriever/](https://docs.llamaindex.ai/en/stable/examples/retrievers/bm25_retriever/)



**Figure 1:** RAG Pipeline Diagram

generation. The Llama3-70b model evaluated the context based on two criteria: (1) the relevance of the retrieved context to the subject matter of the user’s query, and (2) the degree to which the context contained sufficient information to fully answer the query. These assessments contributed to a final context alignment score presented on a 0-to-4 scale.

The prompts used by the Llama3-70b model were adapted from the correctness evaluation<sup>15</sup> and context relevancy evaluation<sup>16</sup> modules available within the LlamaIndex framework. These prompt templates are included as an attachment at the end of this paper for full transparency and reproducibility.

In parallel, two human annotators independently evaluated the final generated answers. They assessed relevance and correctness on a 0-to-4 scale and provided

qualitative notes detailing their reasoning, pointing out strengths or weaknesses such as omissions or inaccuracies. To evaluate the reliability of the annotations, we computed inter-annotator agreement using Krippendorff’s Alpha [14, 15], which is particularly well-suited for ordinal data. The calculation results in a value of 0.90, suggesting strong agreement between annotators. In case of disagreement between the two annotators, a third annotator evaluated the instance to determine which of the two grades was more in line with the guidelines (see C. Guidelines for Human Annotation). Consensus was then reached by majority vote.

This comprehensive evaluation process utilized a dataset of 71 question-answer pairs, selected from a larger pool generated across all 109 source documents (covering both university courses and libraries). Notably, 10 of these 71 pairs were specifically designed to query information contained within the library documents, ensuring assessment of the system’s performance on that subset of data. Overall, the system demonstrated comparable

<sup>15</sup>[https://github.com/run-llama/llama\\_index/blob/main/llama-index-core/llama\\_index/core/evaluation/correctness.py](https://github.com/run-llama/llama_index/blob/main/llama-index-core/llama_index/core/evaluation/correctness.py)

<sup>16</sup>[https://github.com/run-llama/llama\\_index/blob/main/llama-index-core/llama\\_index/core/evaluation/context\\_relevancy.py](https://github.com/run-llama/llama_index/blob/main/llama-index-core/llama_index/core/evaluation/context_relevancy.py)



performance across both evaluation methodologies. It achieved an average normalized accuracy score of 83.63% (SD = 16.45%) in the AI evaluation and 79.22% (SD = 28.34%) in the human evaluation. This similarity in overall scores suggests reasonably consistent performance, although individual query evaluations could differ between the AI and human assessors, underscoring the value of the dual approach. Notably, the context evaluation score was 76.36% (SD = 20.89%). Some random test pairs results are shown in Tables 1, 2 3, 4 and 5.

Detailed implementation procedures, including data processing scripts, model configurations, and complete evaluation results, are documented in the associated Jupyter notebook.

## 5. Discussion

### 5.1. Advantages

A significant advantage of the implemented system lies in its rapid deployment capability, stemming from the simplified toolchain. The streamlined setup process enabled quick deployment, facilitating efficient testing and development cycles. This ease of use facilitated the integration of various components, reducing the learning curve and making the system accessible even for individuals with limited prior experience.

Another notable benefit was the availability of multiple components, particularly the Markdown parser, which proved invaluable. The parser effectively handled document processing, ensuring accurate interpretation and formatting of content. This feature enhanced the system's overall functionality, enabling seamless handling of structured documents and consequently improving the user experience.

Despite certain challenges, the system achieved relatively high accuracy in its responses. However, document retrieval remains an area for improvement, presenting an opportunity for optimization to further enhance precision and relevance. Nevertheless, the current results demonstrate promising potential, indicating that the fundamental approach is sound and can be further refined with additional efforts.

### 5.2. Limitations

A primary difficulty encountered was the extensive documentation, which contained a wealth of information requiring considerable time for comprehension and analysis. Understanding the optimal implementation and optimization strategies demanded significant effort due to the complexity of the available options, which necessitated careful evaluation.

Another challenge arose from the numerous potential "blocks," such as different retrievers and rerankers, that

could be integrated into the workflow. The wide array of choices required extensive experimentation to determine the most effective combination, leading to increased development time and complexity.

The necessity of a GPU to support computationally demanding embedding models presented another hurdle. While Google Colab offered an accessible environment for initial development, it occasionally failed to provide adequate hardware resources for intensive tasks. This issue was eventually resolved by transitioning to a local PC equipped with a dedicated graphics card, which provided a more stable and powerful development environment.

A particularly limiting factor was the API rate-limiting imposed on the LLM provider. While high-level methods offered precise functionality, they required multiple API calls per query, resulting in significant costs and increased response times. To mitigate this, a delay was implemented between successive API calls, which, although effective in managing costs, considerably slowed down the evaluation process. Furthermore, the inability to modify built-in API functions to define specific rate limits led to challenges such as unnecessary calls and system crashes.

### 5.3. Other Attempts

One of the most complex approaches attempted was the creation of agents capable of responding to specific questions for each document to enhance response accuracy. However, we ultimately discarded this idea due to the excessive response times, which rendered the approach impractical for real-time applications.

Another challenge was to implement a more comprehensive, state-of-the-art evaluation system, such as Ragas. While this approach showed theoretical promise, API limits prevented us to use more sophisticated evaluation systems.

In conclusion, while the project encountered several challenges, the overall results were promising, demonstrating the potential of the approach. Future efforts should focus on optimizing document retrieval, improving workflow efficiency, and addressing hardware and API limitations to further enhance the system's performance and usability.

## 6. Release

The source code of the RAG pipeline and the dataset used are available on the Github repository of the project.<sup>17</sup>

The data downloaded from the websites of University of Trento and University of Verona is available along with the source where the documents are taken. The Python code of the tool is released under the Apache 2.0 license.

<sup>17</sup><https://github.com/Samu01Tech/myVision-universities-RAG>

## 7. Conclusions and Future Work

In this paper, we presented the development of a Retrieval-Augmented Generation (RAG) system designed to provide students with accurate academic guidance, specifically focusing on university course and library information. The system leverages a streamlined toolchain, incorporating a Markdown parser for efficient data handling and the Llama3-8b-8192 LLM for query processing. While the system demonstrates promising results, there are areas for enhancement.

Future work will concentrate on several key improvements. Firstly, we aim to enhance the evaluation framework to provide a more comprehensive assessment of the RAG model's performance, incorporating metrics for contextual relevance, accuracy, and adaptability. Secondly, the integration of reranking mechanisms will be explored to prioritize retrieved results based on relevance and quality. Thirdly, to ensure robust and scalable performance, we plan to test the model with a wider range of LLMs, such as Gemini, Claude and others.

Finally, we plan to extend the current dataset, which remains relatively small, to improve both the retrieval and generation components of the system. This expansion will allow for more robust model training and better generalization across academic contexts. In addition, we will conduct user studies to evaluate the system's effectiveness in real-world scenarios, gathering insights from student interactions to refine and improve the overall user experience.

Beyond these technical refinements, the myVision service will be expanded to serve a broader audience, including bachelor's degree graduates and students interested in specialized master's programs, and to include more universities. We envision the chatbot as a core component of a larger platform that will offer a dedicated user interface, informative podcasts, and direct interaction with student advisors. Ultimately, this work lays the groundwork for a powerful tool to aid students in navigating their academic journeys.

## References

- [1] K. Ukaoha, J. Ndunagu, F. Osang, et al., A guidance and counseling mobile application (gc mobile app) for educational institutions, *NIPES-Journal of Science and Technology Research* 2 (2020).
- [2] H. Majjate, Y. Bellarhmouch, A. Jeghal, A. Yahyaouy, H. Tairi, K. A. Zidani, Ai-powered academic guidance and counseling system based on student profile and interests, *Applied System Innovation* 7 (2023) 6.
- [3] L.-C. Lin, Y.-C. Lai, W.-C. Chang, H.-L. Chiu, T.-Y. Chen, Unicompass: Helping high school students find the right college major, in: *Extended Abstracts of the 2023 CHI Conference on Human Factors in Computing Systems*, 2023, pp. 1–6.
- [4] A. M. A. Al, Prototype rule-based expert system with an object-oriented database for university undergraduate major selection, *International Journal of Applied Information Systems (IJ AIS) Foundation of Computer Science FCS*, New York, USA (2012).
- [5] M. Jawhar, Z. Bitar, J. R. Miller, S. Jawhar, Ai-powered customized university and career guidance, in: *2024 Intermountain Engineering, Technology and Computing (IETC)*, IEEE, 2024, pp. 157–161.
- [6] S. Mirri, C. Prandi, N. Parisini, M. Amico, M. Bracuto, P. Salomoni, User-driven and open innovation as app design tools for high school students, in: *2018 IEEE 29th Annual International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC)*, IEEE, 2018, pp. 6–10.
- [7] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel, et al., Retrieval-augmented generation for knowledge-intensive nlp tasks, *Advances in neural information processing systems* 33 (2020) 9459–9474.
- [8] T. I. Dieing, M. Scheffler, L. Cohausz, Enhancing chatbot-assisted study program orientation, in: *Proceedings of DELFI Workshops 2024*, Gesellschaft für Informatik eV, 2024, pp. 10–18420.
- [9] B. Saha, U. Saha, Enhancing international graduate student experience through ai-driven support systems: A llm and rag-based approach, in: *2024 International Conference on Data Science and Its Applications (ICoDSA)*, IEEE, 2024, pp. 300–304.
- [10] S. Dakshit, Faculty perspectives on the potential of rag in computer science higher education, in: *Proceedings of the 25th Annual Conference on Information Technology Education*, 2024, pp. 19–24.
- [11] H. Modran, I. C. Bogdan, D. Ursuțiu, C. Samoila, P. L. Modran, Llm intelligent agent tutoring in higher education courses using a rag approach, *Preprints* 2024 2024070519 (2024).
- [12] J. Chen, S. Xiao, P. Zhang, K. Luo, D. Lian, Z. Liu, M3-embedding: Multi-linguality, multi-functionality, multi-granularity text embeddings through self-knowledge distillation, in: L.-W. Ku, A. Martins, V. Srikumar (Eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, Association for Computational Linguistics, Bangkok, Thailand, 2024, pp. 2318–2335. URL: <https://aclanthology.org/2024.findings-acl.137/>. doi:10.18653/v1/2024.findings-acl.137.
- [13] N. Chirkova, D. Rau, H. Déjean, T. Formal, S. Clinchant, V. Nikoulina, Retrieval-augmented generation in multilingual settings, *arXiv preprint arXiv:2407.01463* (2024).

- [14] A. F. Hayes, K. Krippendorff, Answering the call for a standard reliability measure for coding data, *Communication methods and measures* 1 (2007) 77–89.
- [15] K. Krippendorff, *Content analysis: An introduction to its methodology*, Sage publications, 2018.



## A. Correctness Evaluation Prompt

You are an expert evaluation system for a question answering chatbot. You are given the following information:

- a user query, and
- a generated answer

You may also be given a reference answer to use for reference in your evaluation. Your job is to judge the relevance and correctness of the generated answer. Output a single score that represents a holistic evaluation. You must return your response in a line with only the score. Do not return answers in any other format. On a separate line provide your reasoning for the score as well.

Follow these guidelines for scoring:

- Your score has to be between 1 and 5, where 1 is the worst and 5 is the best.
- If the generated answer is not relevant to the user query, you should give a score of 1.
- If the generated answer is relevant but contains mistakes, you should give a score between 2 and 3.
- If the generated answer is relevant and fully correct, you should give a score between 4 and 5.

Example Response:

4.0

The generated answer has the exact same metrics as the reference answer, but it is not as concise.

## B. Context Relevancy Evaluation Prompt

Your task is to evaluate if the retrieved context from the document sources are relevant to the query. The evaluation should be performed in a step-by-step manner by answering the following questions: 1. Does the retrieved context match the subject matter of the user's query? 2. Can the retrieved context be used exclusively to provide a full answer to the user's query? Each question above is worth 2 points, where partial marks are allowed and encouraged. Provide detailed feedback on the response according to the criteria questions previously mentioned. After your feedback provide a final result by strictly following this format: '[RESULT]' followed by the float number representing the total score assigned to the response'

Query: \n {query\_str}

Context: \n {context\_str}

Feedback:

## C. Guidelines for Human Annotation

**0: Wrong Answer** The RAG pipeline generated a factually incorrect or completely irrelevant response.

**1: Misses Crucial Information** The answer provided is generally correct but fails to include essential details, making it incomplete or misleading.

**2: Responded Only in Part** The RAG pipeline addressed only a fraction of the user's query, leaving significant portions unanswered.

**3: Could Add More Information** The answer is good, but additional relevant details or elaboration would significantly improve its quality and completeness.

**4: Good Answer** The RAG pipeline provided a comprehensive, accurate, and highly relevant response that fully addresses the user's query."

## D. Tables

The following tables contain two examples of questions used to test the system. The full list of questions, along with the answers, the human annotation, and the AI evaluation and context scores are available on the Github repository (see Section 6).

- Table 1 contains the text of two questions (IDs 43 and 60).
- Table 2 shows the gold and RAG answers, respectively.
- Table 3 contains the AI evaluations scores.
- Table 4 includes the AI context scores.
- Table 5 shows the human evaluation.

ID	Question
43	What are the admission requirements for the Environmental Engineering program, and what further study opportunities are available after graduation?
60	What are the two different orientations that students can choose within the Business Administration and Law degree program, and where can students undertake internships?

**Table 1**  
Sample Questions

ID	Gold Answer	RAG Answer
43	Admission requires an **Italian high school diploma or equivalent foreign qualification, and passing a TOLC test**. The program has limited enrolment and seats allocated to Italian, EU, and non-EU citizens residing in Italy. Graduates have direct access to the **Master's program in Environmental and Land Engineering** at the University of Trento, and can also access other Master's programs in engineering that meet admission requirements. Transfers from other degree programs or universities are also possible, with applications due by February 3, 2025.	Admission requirements for the Environmental Engineering program: No specific admission requirements mentioned. Further study opportunities after graduation: * Direct access to the Master's program in Environmental and Land Engineering at the University of Trento. * Access to other Master's programs in engineering disciplines that meet admission requirements. * Graduates can continue their education with a Master's Degree in Civil Engineering or related fields such as Energy Engineering or Environmental Meteorology. * Alternative options include enrolling in specialized programs at other universities to further enhance expertise in industrial engineering disciplines.
60	Students can choose between two orientations: **one towards private professions and the other towards public professions**. Those choosing the private professions orientation will complete their mandatory internship mainly within **professional firms or auditing companies**. Students choosing the public professions orientation can complete their internship in **public entities and companies**.	The two different orientations that students can choose within the Business Administration and Law degree program are: 1. Orientation towards private professions, mainly within professional firms or auditing companies. 2. Orientation towards public professions, it can take place in public entities and companies. Students can undertake internships in either of these orientations.

**Table 2**  
Gold Answer vs. RAG Answer

ID	AI Evaluation Score (1-5)
43	3.5 The generated answer is relevant to the user query, as it addresses both the admission requirements and further study opportunities. However, it lacks specific details, such as the Italian high school diploma or equivalent foreign qualification, and the TOLC test, which are mentioned in the reference answer. Additionally, the generated answer provides some correct information, like direct access to the Master's program in Environmental and Land Engineering, but also includes some vague or incorrect information, such as "no specific admission requirements mentioned" and "alternative options include enrolling in specialized programs at other universities".
60	4.5 The generated answer is highly relevant and correct, as it accurately identifies the two orientations within the Business Administration and Law degree program and provides the correct information about where students can undertake internships. The generated answer is very similar to the reference answer, with only minor differences in wording and structure. The generated answer is clear and concise, making it easy to understand for the user.

**Table 3**  
AI Evaluation Scores

ID	AI Context Score (0-4)
43	1. Does the retrieved context match the subject matter of the user's query? The retrieved context partially matches the subject matter of the user's query. The query asks about the admission requirements for the Environmental Engineering program and further study opportunities after graduation. The context provides information about various programs, including Environmental Engineering, but it does not directly address the admission requirements. However, it does provide information about further study opportunities, which is relevant to the query. Score: 1.5 2. Can the retrieved context be used exclusively to provide a full answer to the user's query? The retrieved context cannot be used exclusively to provide a full answer to the user's query. While it provides some relevant information about further study opportunities, it does not address the admission requirements for the Environmental Engineering program. Additional information would be needed to fully answer the query. Score: 1 Final Result: [RESULT] 2.5
60	1. Does the retrieved context match the subject matter of the user's query? The retrieved context is highly relevant to the subject matter of the user's query. The context provides detailed information about the Business Administration and Law degree program, including its objectives, curriculum, and career opportunities. The context specifically mentions the two orientations within the program, which aligns with the user's query. However, the context is quite lengthy and covers a wide range of topics, making it challenging to extract the exact information related to the query. Score: 1.8/2.0 2. Can the retrieved context be used exclusively to provide a full answer to the user's query? The retrieved context provides a comprehensive overview of the Business Administration and Law degree program, including the two orientations mentioned in the query. However, the context does not directly answer the question about where students can undertake internships. Although the context mentions internships and provides some information about the internship experiences, it does not explicitly state where students can undertake them. Score: 1.5/2.0 Final Result: [RESULT] 3.3/4.0

**Table 4**  
AI Context Scores

ID	Human Evaluation Score (0-4)	Human Evaluation Notes
43	2/4	The RAG answer provides an accurate and detailed overview of postgraduate study opportunities, including direct access to the relevant Master's program and other engineering-related fields, which aligns well with the Gold answer. However, it entirely omits the admission requirements, including the essential TOLC test and diploma criteria, as well as the program's limited enrolment structure. This missing information is critical to the question, resulting in a response that is only partially complete.
60	4/4	The RAG answer accurately identifies the two orientations—private professions and public professions—and correctly associates each with the corresponding internship opportunities. The phrasing is slightly different but conveys the same meaning as the Gold answer. The response is complete, accurate, and fully aligned with the reference.

**Table 5**  
Human Evaluation