# RAG for searching for answers in school literature
# Project Report for NLP Course, Winter 2024

**Maja Andrzejczuk**
Warsaw University of Technology
`01161548@pw.edu.pl`

**Piotr Bielecki**
Warsaw University of Technology
`01161399@pw.edu.pl`

**Jakub Kasprzak**
Warsaw University of Technology
`01161454@pw.edu.pl`

**Maciej Orsłowski**
Warsaw University of Technology
`01161526@pw.edu.pl`

**Paweł Gelar**
Warsaw University of Technology
`01161422@pw.edu.pl`

**supervisor: Anna Wróblewska**
Warsaw University of Technology
`anna.wroblewska1@pw.edu.pl`

## Abstract

Retrieval-Augmented Generation (RAG) is becoming an increasingly popular technique of improving Large Language Models (LLMs) accuracy and reliability. This project focusses on developing a RAG model designed to query the content of a school reading effectively (e.g., an assigned book). The model will retrieve relevant information from the text, enabling students to interact with the book's content in a question-answer format. It will be joined with an LLM for easier interaction and possibly even verification of student responses.

## 1 Introduction

Retrieval-Augmented Generation is a crucial branch of Natural Language Processing. Its base use is to locate document passages that best answer queries made by the user. Such models can be combined with Language Models (LMs) to provide a user-friendly interface for document interaction. Combining RAG models with LMs is a prevalent approach with other use cases, such as improving the reliability of language model responses (Li et al., 2024). This is because providing a base for a model's response makes it easier to avoid hallucinations, but it also allows verifying the language model's output. This also allows us to provide up-to-date data to a language model, which on its own relies on training data limited to a certain moment in the past.

Retrieval-Augmented Generation is divided into a few stages (Gao et al., 2024). The first step is to collect the documents and split them into chunks, often referred to as passages. These passages are then indexed and stored in a vector database (Pan et al., 2023). Such databases index the passages to allow for quick searching for them based on some query. This capability is used to search for relevant documents when a user submits a query, for example, a question regarding a book. Indexing is possible because the contents of the document are converted to embeddings, that is, a numerical representation of text. Such embeddings are normally performed by machine learning models. Quick search time allows retrieval of some documents, which then proceed to the next stage, which is referred to as the post-retrieval stage. This step is not part of a traditional RAG pipeline; however, nowadays it is often added to further boost the performance of RAG systems. The documents here are ranked to allow for more efficient injection into a query for the large language model. In this way, less relevant documents can be omitted, allowing LLM to focus on the most important data, instead of overloading it with tons of documents. The reason why re-ranking is done after retrieval is that it is less efficient when

performed on numerous documents. This way, numerous relevant documents can first be searched for using a vector database, and then a more accurate relevance measurement can be performed on the top k documents retrieved from the database by a re-ranker model. Finally, prompt engineering techniques are used to combine the original user query with the contents of the documents into a final query for the LLM, which then provides an answer to the user.

Our project will involve implementing a RAG system for information retrieval from school literature readings. We believe this domain is a perfect use case for such models, and to the best of our knowledge such systems have not been implemented before. We aim at utilising the power of language models to help the students search through their school readings, thus making the studying process significantly more efficient in a world where time is a valuable but limited currency. We also intend to combine our retrieval model with a language model. This will make it easier to use for users, as they will be able to receive the answers to their questions in natural language.

Our goal is to create a system, where students will be able to ask questions related to any given school literature. The system will then point the student to a relevant fragment of the book, which the student can use to answer their question. Additionally, if we integrate this solution with a language model, the model can also use the extracted fragment to provide an answer to the question, which will improve the user experience.

Our solution will be developed in the Python programming language, most likely in the *Google Colab* environment, because it provides access to GPU, which significantly speeds up computations. A detailed concept for our solution is described in section 3.

## 2 SOTA Review

### 2.1 Related works

To deepen our knowledge in the selected area of research, we explored several state-of-the-art works related to Retrieval-Augmented Generation. While we might not eventually incorporate most of them into our project for various reasons, they give us valuable insights into various aspects of the topic. Therefore, we decided to explore these works in hope of finding inspiration for our work,

and perhaps even ideas for future work on our system after the course project is concluded. Let us then explore some of the modern approaches in this field.

Combining Large Language Models with RAG is an effective way to make the models more trustworthy, as their answers can be verified using the parts of the document which were selected to generate the response. This paradigm is called Verifiable Generation and one of the proposed approaches is LLatrieval, or *Large Language Model Verified Retrieval* (Li et al., 2024). This method addresses the issue of RAG models hindering the performance of the whole pipeline due to them being significantly less complex than LLMs, thus rendering the pipeline less effective than standalone LLMs. It introduces a verify-update iteration. In the *verify* step, the LMM provides feedback to the RAG model on the document it selected as a base for generating the answer. This feedback then can be used to avoid relying on low-quality retrievals. If this is the case, the next step, *update*, aims to update and improve the retrieval. Through experiments, the authors of this framework were able to achieve new state-of-the-art results.

Document retrieval is key to obtain an accurate response to queries in RAG systems. One of the works proposes *Unsupervised Passage Re-ranker*, or UPR in short (Sachan et al., 2022), which aims to improve passage (document) retrieval process. Its main feature is the zero-shot generation of questions when re-ranking. It then utilises a pretrained LLM, without any fine-tuning, to estimate the probability of obtaining a question given a passage text. The zero-shot approach makes the evaluation process more robust and forces a model to better explain its generated questions. The authors prove the effectiveness of this approach compared to other state-of-the-art solutions and claim that the gains from using this method will grow with the improvements of available language models. While our work plan does not involve utilising this method in the re-ranking stage, we might consider experimenting with it should the time allow us to do so, or, more likely, as a part of future work on our solution after the project.

A common problem when developing systems based on RAG is their evaluation. This is because it heavily relies on manual human work of domain experts, which is time-consuming and expensive.

To minimise the reliance on this aspect, a RAGAS framework was proposed (Es et al., 2023). The authors stress that the biggest strength of this solution is its ability to evaluate RAG systems even without access to ground-truth answers to evaluation questions, which might be useful in our setting. RAGAS measures context relevance (is the retrieved passage relevant to the query), answer faithfulness (is the answer generated relevant to the retrieved passage, or is it hallucinated), and answer relevance (is the answer relevant to both the query and the passage). RAGAS utilises an LLM under the hood to perform evaluations. While the authors use powerful models in their implementation, the package allows for using more lightweight models.

Another framework for automatic RAG evaluation, which is even newer, is ARES, which stands for *Automated RAG Evaluation System* (Saad-Falcon et al., 2024). It works by leveraging LLM judges, trained on queries, passages, and answers. Queries and answers are generated synthetically by LLMs using provided passages. These judges are then applied to assess query-document-answer triples obtained from any RAG system. ARES evaluates the same three metrics as RAGAS, and three separate LLM judges are used for this. This framework is proven to be more accurate than other modern evaluation approaches including RAGAS, while maintaining a low cost of evaluation, with reduced need for human annotations. However, one of its core limitations are high computational requirements, which might eventually make it hard for us to utilise.

There are also attempts at incorporating the ideas of Active Learning into Retrieval-Augmented Generation, most notably a *Forward-Looking Active REtrieval augmented generation* (FLARE) method has been proposed (Jiang et al., 2023). Its goal is to address the issue of hallucinations in language models, by actively retrieving external documents when generating answers to provide more accurate response. It is different to most models, which perform such retrieval only once upon receiving an input, and then use it to generate the whole response. Some more advanced models might perform this action multiple times in fixed intervals, but often such retrievals are performed when they are not necessary. In contrast, FLARE retrieves only when it generates low-confidence tokens, i.e. part of the response

that likely could be incorrect. It then retrieves the documents which then are used to regenerate the response, leading to it being more accurate. The use of active learning principles also prevents unnecessary retrievals.

## 2.2 Datasets

The data we intend to be using comes from a free library of eBooks, named *Project Gutenberg* (https://www.gutenberg.org/). It claims to contain over 70,000 different eBooks, which includes school literature that we are interested in. We will select some examples of school readings, and use them in our project. They are available in the plain text format, but can also be converted to PDF – our solution will support loading data from both formats. There is no API available for this data, which means collecting it will require substantial manual work. However, it is definitely feasible considering the volume of the data that we can realistically work on given the computational resources available to us. We will collect them in a *GitHub* repository, which we can clone in *Colab* to have easy access to them.

Regarding the model evaluation datasets, unfortunately we were not able to locate a free dataset with school literature retrieval based exam questions which could be used for evaluation of our system. We could keep searching, but as an alternative, we could manually design question-answer pairs for our system. With five people in the team, it should be feasible to develop an exam dataset of at least 150 entries.

## 2.3 Models

We intend to use open-source models, and we will focus on pre-trained ones, since we will not have high computational power at our disposal during the project. There are three types of models we were looking for:

- embedding models,

- re-ranker cross-encoder models,

- large language models.

We found several models on HuggingFace platform, which could be useful for our system. Let us now discuss them in more detail. Starting with embedding models, the first model we consider using is called *all-MiniLM-L6-v2*, which has 22.7M parameters. It is a part of *sentence-transformers*

library (Reimers and Gurevych, 2019). This library includes numerous models aiming to generate embeddings for text input data. These models are based on Sentence-BERT, which itself is a modification of BERT (Devlin et al., 2019). The authors of Sentence-BERT claim that cross encoder embedded into BERT is not efficient, and their proposed method is based on siamese networks. This allows it to generate high-quality fixed-size embeddings, and the *all-MiniLM-L6-v2* model, which we intend to use, is a lightweight example of such a model, and will allow us to generate embeddings for passages for efficient storage in a vector database.

A larger alternative for this model would be *bge-large-en-v1.5* (Xiao et al., 2023), which has 335M parameters. This model comes from a large family of models named BGE (BAAI General Embeddings), which originally was introduced in a paper that proposed a method for embeddings for Chinese language. However, this family contains not just models for Chinese language, but also Multilingual, and models which work just for the English language, and the model we chose is exactly one of such.

When it comes to re-ranking models, the ones that we are considering using are *bge-reranker-large* with 560M parameters, and *bge-reranker-base* with 278M parameters (Xiao et al., 2023). They also come from the BGE model family. These are cross-encoder models, which utilise attention mechanism to rank the retrieved documents. Contrary to embedding models, these models take user query and documents and output similarity scores. These scores are not constrained to any range.

Finally, there is an LLM which will generate the final outputs. We are looking for a lightweight model due to limited resources. One that we could potentially use is Cerebras GPT (Dey et al., 2023), which is available for free on HuggingFace. It comes in seven variants, ranging from 111M to 13B parameters. We will most likely use the variant with 2.7B parameters, as in our opinion it provides a good balance between quality and computational requirements. These models were released by Cerebras Systems and are pre-trained on the Pile dataset (Gao et al., 2020). The authors show in their experiments that they are compute-optimal, at least compared to other state-of-the-art models. They claim that models from this family

offer the most quality out of other models trained on the same budget, due to them using the latest advancements in LLMs training and scaling techniques.

Alternatively, we could also try *Llama3.2* (Touvron et al., 2023), which is a lightweight release from a group of *Llama 3* models, and comes in two sizes: one with 1B parameters, and the other with 3B parameters. This gives us flexibility when it comes to choosing the most optimal one. Llama models are released by Meta AI and are trained solely on publicly available data, unlike some of the other LLMs. What is very useful from our perspective is that the authors state Gutenberg as one of the data sources on which their model was trained – this is the same data source that we intend to source the school readings from. The only problem with using this model is that while it is completely free, it requires approval from its owners. While it is usually granted for students, one needs to provide their personal details in order to obtain it. Once the access is granted, an access key from HuggingFace is needed to use it.

A third alternative that we found is a model called *Qwen2.5-7B-Instruct* (Team, 2024). This model is the newest release from the Qwen2 family of models (Yang et al., 2024), which is open-source and available without any restrictions. It is a larger architecture than all the other proposed, however it is quantised to *Int4* data type, which enables its usage in environments with less computational power, whilst maintaining solid performance. Its parameter number varies from 0.5 up to 72 billion, depending on the variant.

All the models that we listed above can be loaded through either *transformers* (Wolf et al., 2020), or *sentence-transformers* (Reimers and Gurevych, 2019) packages. Instructions on how to do that are included on these models' homepages on HuggingFace.

## 2.4 Vector database

As mentioned earlier, vector database is a key component of any RAG systems as it allows for quick similarity searches used when retrieving documents relevant to a query. The vector database that we intend to use is Meta's FAISS (Douze et al., 2024). It is an open-source library. The authors point out the flexibility of their library in that it can use multiple indexing methods, instead of just some preselected one. FAISS does

not generate embeddings on its own, and it requires an external embedding model to generate vectors for the data.

## 3 Work plan

Let us now discuss our plan for this project. Our work will start by performing exploratory data analysis on our data. We found a *LangChain* (Chase, 2022) library, which will not only help us load the PDF documents, but also contains other useful functionalities for RAG tasks, which we may utilise further in the project. Once we have investigated our data, we will start pre-processing it. There are many state-of-the-art Python libraries we could utilise for this task, and beside *LangChain*, the one we find the most promising is *Unstructured* (`https://unstructured.io/`). Both can be easily integrated with each other and provide many functions for pre-processing of PDF documents. We will definitely need to split selected books into chunks, which we could do, for example, by utilising their structure and dividing them by paragraphs, or by just splitting them into fixed-sized chunks. We also need to clean them, for example, by removing redundant whitespace. However, the exact preprocessing steps will need to be decided only once we have explored our data.

Simultaneously, we will split our efforts to also work on assessing the models. We will need to verify that:

- they load correctly in our environment,

- they are feasible to run and use (in terms of computational resources).

We plan to perform these checks early in the project, so that in case anything goes wrong, we have time to find emergency alternatives for these models.

Once the documents are pre-processed, and the models are loaded, we can feed the documents to the embedding model. It will generate vector embeddings, which we will need to store in a vector database, which allows for efficient vector similarity search, and also is integrated into *LangChain* library. At this point, we will be able to take inputs from a user, and retrieve relevant documents from our database. Then, these documents will be passed to a re-ranker cross-encoder model, which will rank the retrieved documents to help select the most relevant one. We will also add a language model to act as an interface for interacting with our system.

At this point, we will have the prototype versions of all the key components in our system (preprocessing and RAG models). We will test them on an example book from our dataset. This will serve as the proof-of-concept for our solution, and therefore we plan to arrive at this stage for the midterm checkpoint in December.

Afterwards, we will work on improving individual components and coupling them all together. Our goal will be to create a reusable, customisable pipeline. Internally, it will consist of two parts. The first part will be responsible for loading the documents, generating embeddings, and storing them in a vector database. Since this is a lengthy process, we plan on performing it once for our whole dataset and preserving the vector database for later use. The second part will perform a search for a given query, using the database from earlier. However, the users will have a choice to either utilise our ready-made database, or to create theirs by uploading their own books. Additionally, the whole pipeline will be contained within a single class, which will make it easier to use. Listing 1 shows an early prototype vision of the design of this pipeline (the names of the class and methods are going to be more specific in the final solution).

Listing 1: Prototype design of the pipeline

```
# Default database
pipeline = Pipeline.default()

# Custom books
pipeline = Pipeline(
    books=['./my_book.pdf']
)

# Using the pipeline
pipeline.query(
    'year of the Easter Rising'
)
```

At this stage, with our system working, we will conduct multiple experiments to test the quality of our solution. Let us discuss our plan for the evaluation of our RAG system. Once we have the evaluation dataset as discussed in one of the earlier sections, we could proceed with one of the following:

- Utilise a framework such as ARES (Saad-Falcon et al., 2024). The strength of ARES

is that it requires very little human annotations, however our initial tests proved that it is computationally heavy and thus might not be feasible to run. However, it will be a good idea to at least try this solution in the experiments phase;

- Utilise a lighter alternative to ARES. An example of such is RAGAS (Es et al., 2023), which, while less accurate than ARES, is still a state-of-the-art method which is a bit more lightweight in terms of computational requirements;

- If both frameworks proposed above end up not being feasible to use, and we are unable to find another alternative, we can always fall back to using an LLM judge. This requires prompt engineering, which aim is to make a query for a judge that will output concrete metrics for each sample generated by our system. At this stage, we could either utilise some lightweight LLM model such as Llama 3.2 (Touvron et al., 2023), or a larger model like GPT-3.5 through its web interface (OpenAI, 2024) by manually pasting the prompts and gathering the responses.

Once we make sure that everything works, we will export the whole pipeline from our notebook into a repository on GitHub (the same one where the data will be stored), so that we can easily access it by cloning the repository in our *Colab* environment. This will also declutter the notebook, and make it easier to use by the users. This will mark the end of the project, which we will present on the deadline day in January. A timeline for our work is presented in Table 1 (dates in bold are checkpoints).

| Date | Planned progress |
|---|---|
| 27/11 | Project proposal |
| 04/12 | EDA, verifying models feasibility |
| **11/12** | Proof of concept |
| 18/12 | Reviews |
| 08/01 | Reusable pipeline |
| 15/01 | Experiments |
| **22/01** | Final presentation and report |

Table 1: Detailed work schedule

## 4 Proof of concept

In this section, we describe the work we have done on the proof of concept. We focus here on our methodology, explaining the environment, data transformations, which models we have eventually used, their parameters, etc., in order to provide grounds for reproducibility. We also mention our first findings related to this system.

### 4.1 Environment

Let us begin by discussing the environment in which we developed our proof of concept. We used Google Colab, and the code required 6GB of Google Drive storage space to store all the data and models downloaded by the code. We run our code on the *T4 GPU* runtime. This is critical, as running this code on CPU would take significantly more time, especially the steps involving models inference. The code cells that can take longer than 15 seconds to run are marked in the notebook with a comment, which informs of the average execution time (in the *T4 GPU* runtime). Note that we did not use any paid plan for the runtime. This made the development difficult due to occasional disconnects, meaning we would have only up to around two hours of development each day, and after that we would be forced to wait around 24 hours before we could connect to any GPU runtime again. This forced us to work on the project in short but regular bursts.

Additionally, right at the beginning of the notebook, we set all relevant random number generators seeds. This ensures that the code in the notebook produces the same results every time it is run.

### 4.2 Data loading

For our analysis, we collected a collection of books available through Project Gutenberg, which is a large online library of free e-books. These books are stored in our repository in the `data` folder (link: `https://github.com/andrzejczukm/NLP-RAG`). In total, the dataset includes 35 books, mostly in `.txt` and `.pdf` formats. These texts were carefully selected to include only books known worldwide and recorded for use in our analysis. We conducted a detailed exploratory data analysis (EDA) to better understand the structure and characteristics of the texts to help guide the next stages of our NLP project.

## 4.3 Exploratory Data Analysis

The EDA process involved analyzing the data structure and results of splitting input texts into chunks for further NLP processing. Initially, we encountered problems with chunk splitting, especially for documents in `.txt` format. These problems stemmed from inconsistencies in the way we handled line splits and separators. The following section outlines our approach and subsequent improvements.

### 4.3.1 Initial Analysis

In our initial analysis, we examined the distribution of passages across documents and the variability in passage lengths. These visualizations revealed significant inconsistencies in chunk sizes, particularly for `.txt` files.

The first plot (Figure 1) shows the number of passages generated per document during the initial chunking process. This revealed significant variation, indicating inconsistent chunking logic. Some documents were split into far too many passages due to the aforementioned newline issues. The next plot (Figure 10) shows the distribution



Figure 1: Number of passages per document in the initial chunking process. Notice the uneven distribution across documents.

of passage lengths for each document. While we aimed for a consistent passage length of approximately 2000 characters, the results deviated significantly. On the logarithmic scale, we can see that five documents that exhibited especially unusual distributions, named:

- `kafka_metamorphosis.txt`

- `milton_paradise_lost.txt`

- `homer_iliad.txt`

- `kafka_the_trial.txt`

- `conrad_heart_of_darkness.txt`

The above-mentioned files are identified by the nearly invisible or highly compressed boxes, suggesting irregular or extreme passage lengths compared to other documents.

The histogram in Figure 2 provides a global view of passage lengths across all documents. While the target passage length was 2000 characters, many passages were excessively short.



Figure 2: Histogram of passage lengths before improvements. Note the cluster of very short passages, especially in `.txt` files.

### 4.3.2 Problem Finding and Resolution

During the initial analysis, we encountered significant discrepancies in the way chunks were generated from `.pdf` and `.txt` files sourced from the Gutenberg library. The problem was caused by differences in the formatting of the two file types. While `.pdf` files retain their original structure, `.txt` files from Gutenberg often contain artificial line divisions that disrupt the natural flow of text, making it difficult to properly divide content into meaningful chunks.

To better understand the issue, we compared sample of `.pdf` and `.txt` files sourced from Gutenberg. As seen in Figures 3 and 4, the formatting of the two file types is quite different. Figure 3 shows a `.pdf` file with well-structured content, where paragraphs are properly aligned and separated, with no artificial line breaks.

In contrast, Figure 4 illustrates a `.txt` file from Gutenberg, where artificial line breaks are inserted between lines of text, causing fragmentation and making it difficult for the chunking process to divide the text appropriately. Upon identifying this

Figure 3: Sample chunk from a `.pdf` file. The passage structure is consistent and reflects coherent textual units.



Figure 4: Sample chunk from a '.txt' file. Artificial line breaks caused fragmentation of text into incoherent passages.

issue, we applied the `paragraph-grouper = group-broken-paragraphs` setting during the document loading process. This setting addressed the artificial line breaks in the Gutenberg `.txt` files, grouping fragmented text into coherent paragraphs. As a result, the chunking process generated consistent and meaningful passages, allowing for better handling of the text for further analysis.

### 4.3.3 Improved Chunking Approach

After addressing the issue of artificial line breaks in the Gutenberg `.txt` files, we improved the chunking process by ensuring that text was divided into coherent passages with a consistent length. The result of this improvement is reflected in the following visualizations.

In Figure 11, we present boxplots of passage lengths per document after the improvements. The variability observed in the previous boxplots has been significantly reduced. The plots now demonstrate a more consistent distribution of passage lengths across all documents, with fewer extreme outliers. This indicates that the chunking approach is now more stable and effective at dividing the text into roughly equal-sized passages, as intended.

Similarly, Figure 5 shows the histogram of passage lengths after the improvements. The distribution now closely aligns with the intended chunk size of at most 2000 characters. This histogram illustrates that the chunks are more uniform in length, with fewer instances of unusually short or excessively long passages. However, it is evident that the majority of the chunks are slightly smaller than the intended 2000 characters, likely due to the nature of the text and how the splitting process handles natural language.

We ended up splitting the documents using `RecursiveCharacterTextSplitter` with separators set to double new line, full stop, question mark, and single newline characters, in that order of preference. This means that we try to split documents into chunks shorter than maximum size using double new line, if we fail then we split on full stops etc. With this approach, the only problematic book was 'Ulysses', because there were ten paragraphs without any punctuation. It is a rare quirk of this book, and we assume that it is abnormal. Even those should not be a significant problem, as they are an order of magnitude shorter than the maximum context length for the chosen LLM ( 5.5k letters vs  30k tokens). Around 60 (of 50k) passages are just punctuation. We can just ignore them, as the database and re-ranker should not return them, or we may not add them to the database.

These improvements in chunking are vital for ensuring that the data is properly structured for further NLP analysis. With consistent and coherent passages, we are now able to proceed with more accurate text processing and analysis steps.

### 4.4 RAG Model Performance Testing

For the proof of concept, we limited the data to just one book – "Ulysses" by James Joyce (`joyce_ulysses.pdf` in our dataset).

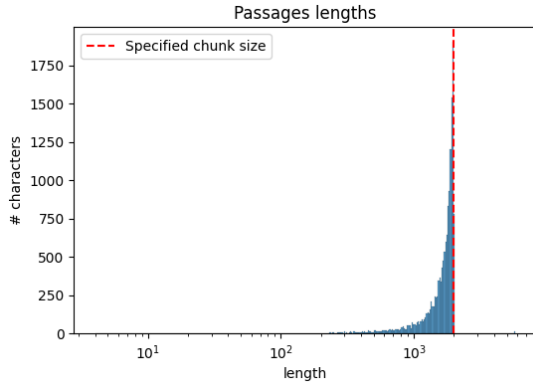We begin by loading the file and splitting the

Figure 5: Histogram of passage lengths after improvements. The distribution aligns closely with the intended chunk size.

book into passages. We use a chunk size of 2000 which we believe should be enough for a question-answer system. We also set a chunk overlap of 200 to allow context to transfer between chunks. We believe that this will help if the information that needs to be retrieved is located near the boundary of some chunk. We split the document by a double newline character. We also utilise two pre-processors from the *Unstructured* library – one for removing excess whitespace (for example, when two or more white spaces occur after one another) and the other for grouping broken paragraphs (paragraphs that are split by newline characters solely for visual purposes). This resulted in 5724 passages. A great advantage of using *LangChain* is that we can load `.pdf` and `.txt` documents exactly the same way.

As discussed in the work plan section (Section 3), we attempted to use all the models that we listed in Section 2.3 in order to test their feasibility in our environment. The first was the embedding model, and the one we eventually selected for our proof of concept is *all-MiniLM-L6-v2*. We tried to use the larger model that we proposed, *bge-large-en-v1.5*, however, due to its higher complexity, it was taking much longer in terms of runtime. However, in this case, we do not exclude the possibility of using it in our final solution. This is because we need the embedding model mostly in the process of creating a vector store, which, for the final solution, will need to be performed only once. Thus, the long computation time might not matter as much, provided that there are no hardware limitations. Once the vector database is created, we can then save it and load it whenever we start

the system, and only use the embedding model on the query, which we assume would not be computationally expensive. However, this will require further testing at later stages of our work. Furthermore, at this stage, where the solution is highly unstable, it is still easier to work with the more lightweight model, thus the current approach uses *all-MiniLM-L6-v2*.

When it comes to the reranker model, again, we tested both proposed architectures: *bge-reranker-base* and *bge-reranker-large*. Fortunately, for this model, the computation times were very quick for both variants, so we were able to use the larger model in our solution.

Finally, there is the large language model, which interacts with the user by producing the final output. Initially, we tested various models from the *CerebrasGPT* family. However, we quickly realised that their performance was underwhelming, no matter the model size. Their answers were usually highly inaccurate, and, in addition, their outputs would often repeat parts of sentences in an infinite loop. Thus, we decided to abandon these models completely and instead try the *Llama3.2* models. We attempted to utilise both 1B and 3B parameter variants, but unfortunately, the environment we used did not have enough virtual memory on the GPU to load the larger model. Thus, we decided to utilise the smaller variant with 1B parameters. However, this model's performance did not satisfy us. While it would occasionally produce correct responses, it was very inconsistent and would often hallucinate or draw incorrect conclusions from the passages. Therefore, we tried the final proposed LLM which was *Qwen2.5-7B-Instruct* quantised to *Int4* data type to enhance performance, which we use in the final proof-of-concept as described further. Choosing instruct-type model greatly increased the quality of generated text, as the model always tried to answer the given questions instead of generating new ones or trying to write more of the book.

Having loaded the data and chosen the models, we constructed a prototype version of the system. Let us now discuss it. In overview, it works as outlined in the work plan section (see section 3). Thus, let us only describe the specific parameter choices for different components. Firstly, we set the number of best-matching documents retrieved from the vector database to 30. We found the value of 30 to be a good balance between the amount

of information retrieved, and the computation time for the re-ranker model. We then use the five best passages according to the re-ranking model to construct the prompt for the LLM. We build it in the following way: It starts with *"Based on these passages:"*, and then the five best passages are concatenated. Then we put *"Answer this question:"*, and we paste the original query. Note that all these components of the query, as well as the passages, are concatenated with a double new line character as a separator. For now, we limit the length of the answer generated to 50 tokens.

We run a simple test on this prototype system, to verify that it works and has potential. The query that we passed to the system was: *"Who is Leopold Bloom's father?"*. While the response was abruptly cut at the end due to the maximum tokens limit, the model correctly stated that the father of Leopold Bloom is Rudolph Bloom (also referred to as Rudolf Virag). It also provided some additional minor facts about him. Upon the inspection of the passages that were used to generate the query, we saw that the answer was in both the first and the second passage, as ranked by the re-ranking model. The other listed facts were also clearly based on the best-ranked passages. Not only does it prove the re-ranker's quality, but also the quality of the retrieval from the vector database, as it was able to locate these passages from the thousands of others available. The model also did very well when it came to describing Buck Mulligan's appearance.

We tested whether the answers were in fact the result of retrieval augmentation, or if LLM was just able to answer those questions on its own. We did that by asking it the same questions but without the passages. It then answered that the Leopold Bloom's father is named David, which is straight-up wrong. It also describes Buck Mulligan's appearance confidently, but is wrong again.

To conclude our proof of concept, despite some minor problems along the way, the draft of the pipeline we provided proves that this system looks promising.

## 5  Final Solution

Our final pipeline works in a way similar to that described in Section 3. Let us now precisely describe the most important aspects of its API.

The whole logic for the RAG system is contained within one class named `RAG`. Its default

parameters are identical to the ones used in the Proof of Concept, since they proved to be the best then, however unlike then, the system now offers high customisability of these. There are two ways of initialisation; the first one is through a `quickstart` method, which is the fastest and is recommended for most use cases. This method takes two parameters:

- `use_default_db` – whether to initialise the system with the database that contains default books (i.e., books that we have collected; see Sections 2.2 and 4.3). This setting is `True` by default.

- `book_dir` – a path to a directory that contains any extra books that a user wishes to load into the RAG system. This is set to `None` by default, which means no additional books are added.

The other method of initialisation of the system involves calling the constructor directly. Then, the user has more customisation options, which include:

- postprocessors to use when loading the documents;

- chunk size and overlap;

- prompt template used to query the LLM;

- separator for the retrieved passages, used when forming a prompt;

- number of documents retrieved from the vector store, and the number of documents used for the final query.

Note that when initialised this way, the user needs to manually load all models and databases through methods beginning with `load_` (these are called automatically when initialising through `quickstart`).

Once initialised, new books can be added at any point using the `add_books` method, which takes a list of paths to the book files to add. This is naturally optional if the RAG system is initialised with default books but can still be used to enrich the collection.

Then, to use the system, a user can call an `answer` method, which takes a query as the only required parameter. Additionally, the user can specify the maximum length of the system output

through the `max_new_tokens` parameter, as well as apply filtering using some document by passing a path to the `document` parameter.

# 6 Evaluation

The final phase of our project work involved evaluating our system to ensure its quality. The tests we conduct here are much more robust than the simple tests we conducted as a part of the Proof of Concept, which served more as a sanity check and a proof that our solution could work, than actual proper tests.

For our tests, we focused on a few selected books, namely *Romeo and Juliet* by Shakespeare and *Odyssey* by Homer. However, the system would have the database full of the full collection of books, which could make it difficult to extract the necessary data. This was actually brought to our attention as feedback from our Proof of Concept, where the database had only one book loaded. This is why during evaluation, we focused not only on RAG's answers, but we also assessed the passages it used to generate it.

Just like with Proof of Concept, all the code that we run for evaluation was run in a Google Colab environment with T4 GPU runtime. We also always set fixed random number generator seeds to ensure reproducibility.

## 6.1 Testing questions

The first step in our system's evaluation involved generating a dataset of testing questions related to school literature, which our RAG would answer. Unfortunately, to the best of our knowledge, there are no such datasets freely available on the internet. Such questions could be created manually, but this would require much time and should be conducted by domain experts to ensure the quality of such questions.

Therefore, for this task we employed a powerful and popular natural language model GPT-4o through its publicly and freely available interface, ChatGPT (OpenAI, 2024). We used it to generate sample questions, answers, and locations of answers for the specified books. As mentioned earlier, we focused on two books – *Romeo and Juliet* and *Odyssey*. This particular choice was important as we are familiar with these books, meaning it was easier to assess whether the generated questions and answers were of high quality.

We prompted the model to generate the ques-

tions datasets in CSV format. We focused primarily on closed questions. This is important for two reasons. Firstly, open-ended questions do not fit our business case, as the system we develop is meant to help students and teachers extract facts from books, and not answer their exam questions for them. Secondly, evaluating open-ended questions is not a trivial task, and would most likely require work of domain experts. Our final evaluation dataset comprised 48 questions related to *Romeo and Juliet*, and 44 questions about *Odyssey*.

## 6.2 Evaluation results

To comprehensively assess the performance of our system, we used in our testing a range of evaluation metrics. These included:

- **BLEU**: Measures the overlap of n-grams between the generated and reference text, commonly used for machine translation evaluation.

- **ROUGE**: A family of metrics (ROUGE-1, ROUGE-2, ROUGE-L) that evaluate the recall of n-grams and longest common subsequences, widely employed in text summarization tasks.

- **BERT-Based Scores**: Metrics such as BERT Precision, BERT Recall, and BERT F1, which use deep contextual embeddings to evaluate semantic similarity between sentences, offering robust detection of paraphrased or semantically equivalent expressions.

The evaluation results are presented separately for each book to provide a detailed analysis of the system's performance.

### 6.2.1 Romeo and Juliet

We began the analysis by calculating all metrics for each reference-output pair. The distribution of evaluation metrics for *Romeo and Juliet* is illustrated in Figure 6, showcasing the variability across BLEU, ROUGE, and BERT-based scores.

The performance metrics for the Romeo and Juliet dataset indicate that the model performs reasonably well in terms of understanding the content, even though it struggles with exact word matching. The low BLEU score suggests that the model does not generate text that matches the reference exactly, especially in terms of n-gram overlap. This is expected, as n-gram matching is a
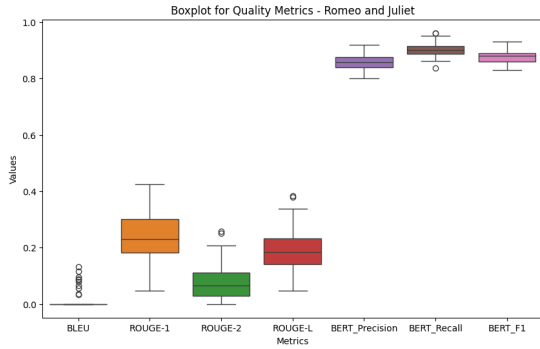
Figure 6: Distribution of Evaluation Metrics for *Romeo and Juliet*



Figure 7: Distribution of BERT-F1 score for *Romeo and Juliet*

| Metric | Mean | Median | Std Dev |
|---|---|---|---|
| BLEU | 0.02 | 3.95e-155 | 0.04 |
| ROUGE-1 | 0.24 | 0.23 | 0.09 |
| ROUGE-2 | 0.08 | 0.06 | 0.07 |
| ROUGE-L | 0.19 | 0.18 | 0.08 |
| BERT Precision | 0.86 | 0.86 | 0.03 |
| BERT Recall | 0.90 | 0.90 | 0.02 |
| BERT F1 | 0.88 | 0.88 | 0.02 |

Table 2: Performance Metrics for *Romeo and Juliet*

strict metric that often doesn't reflect the model's ability to capture the overall meaning or context. However, the ROUGE-1 score, which measures unigram recall, shows better performance, indicating that the system is able to identify many of the keywords from the reference text. The ROUGE-2 score, while lower, still demonstrates some capacity for capturing relationships between consecutive words, though there is room for improvement in fluency and coherence.

The most encouraging result comes from the BERT F1 score, which reflects the model's semantic understanding. With a mean of 0.88, it shows that the system is proficient at grasping the meaning of the reference text, even if it doesn't match every word exactly. This is particularly important for tasks involving responses to book-related questions, as understanding the context and key themes is often more valuable than word-for-word matching. In this sense, the model is performing well, as it generates responses that are semantically aligned with the original text.

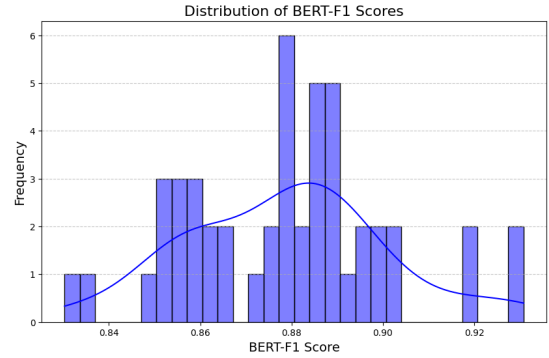We selected 5 pairs that achieved the lowest results according to metric Bert F1, as well as 5 pairs with the highest scores, in order to assess how well the model performs. Below are two examples with the highest and lowest scores:

**High score**

- **How does Benvolio try to help Romeo?**
  **Answer:** Benvolio advises Romeo to forget Rosaline and look at other women to ease his heartache.
  **RAG Output:** Benvolio tries to help Romeo by suggesting he forget about Rosaline and examine other women, and by encouraging him to join in the party to distract him from his love for Rosaline.
  **BERT F1:** 0.9308

- **What is Romeo's reaction to his banishment?**
  **Answer:** Romeo is devastated by his banishment, feeling it is worse than death because it separates him from Juliet.
  **RAG Output:** Romeo reacts with extreme distress and despair to his banishment, viewing it as worse than death due to his deep love for Juliet and the separation it imposes between them. He believes he has been "banish'd from the world" and that his life has ended.
  **BERT F1:** 0.9284

**Low score**

- **What is Romeo's mood at the beginning of the play?**
  **Answer:** Romeo is melancholic and lovesick over Rosaline, who does not return his affections.
  **RAG Output:** Romeo seems to be in a state of confusion and anxiety at the beginning of the play, as he questions his feelings towards

his recent dance partner and worries about the consequences of his love for a girl from his family's enemy. He also expresses a sense of urgency to attend the party despite his misgivings, indicating a mix of excitement and apprehension. (approx. 35 words)
**BERT F1:** 0.8303

- **What is the significance of the poison?**
  **Answer:** The poison represents Romeo's desperate choice, believing Juliet is truly dead.
  **RAG Output:** The poison, particularly arsenic, serves as a powerful tool for murder that leaves no physical trace, making it difficult to detect and prosecute the perpetrator, highlighting the contrast between Eastern and Western approaches to crime and punishment. Its use reflects the sophistication of chemical knowledge and the potential for both malicious intent and hidden motives in the characters' actions.
  **BERT F1:** 0.8521

The high BERT F1 score responses demonstrate that the model can effectively capture the essence of a question, even if the phrasing differs. The answers are consistent in meaning, showing that the model is capable of rewording responses without altering the key message. However, in the low-scoring responses, there is a noticeable divergence in content, often introducing unnecessary details. For example, when explaining the symbolism of the poison, the model overcomplicates the answer by introducing external context, such as the comparison to Eastern and Western crime and punishment systems. This highlights how more abstract or nuanced topics, like symbolism, can be difficult for the model to handle, causing it to add irrelevant or overly complex information instead of staying focused on the main point.

These two examples demonstrate some of the lower-performing pairs, but they reflect an overall accurate understanding of the text. Further analysis of the entire dataset showed that the model performs consistently well, with most pairs achieving higher F1 scores.

As part of the evaluation, an additional test was conducted to verify whether the passages retrieved to answer questions were accurate and relevant or if the model exhibited hallucination — producing responses not grounded in the provided text. A random sample of 10 outputs was manually reviewed, and all of them contained passages that correctly corresponded to the outputs generated by the RAG model. This demonstrates that the model is effectively retrieving relevant passages and generating responses aligned with them. Example:

- **Question:** How does Juliet react when she discovers Romeo's identity?
  **RAG Output:** When Juliet discovers Romeo's identity, she is shocked and conflicted, as he is the son of her family's long-standing enemy, prompting her to exclaim, ""My only love sprung from my only hate!"" and expressing the depth of her newfound love despite their families' feud. She is also dismayed by his banishment, leading her to seek a solution through Friar Lawrence.
  **Fragment of RAG Passage:**

  > NURSE. His name is Romeo, and a Montague,
  > The only son of your great enemy.
  >
  > JULIET. My only love sprung from my only hate!
  > Too early seen unknown, and known too late!
  > Prodigious birth of love it is to me,
  > That I must love a loathed enemy.

Upon reviewing the output, we notice the inclusion of an additional sentence: "She is also dismayed by his banishment, leading her to seek a solution through Friar Lawrence." While this information is relevant to Juliet's broader emotional arc, it wasn't directly prompted by the question regarding her reaction to Romeo's identity. This is a result of the model retrieving passages from multiple chunks, pulling this detail from another part of the passage related to her response to Romeo's banishment.

This observation suggests that providing the model with a larger set of chunks may lead to the inclusion of more context than necessary, potentially resulting in irrelevant information being added to the response. Although the extra detail is accurate and rooted in the text, it indicates that refining the selection of chunks given to the model could help maintain a more focused and relevant response. By limiting the amount of information retrieved, we could reduce the chance of introducing unrelated details, which would help improve

the precision of the model's outputs.

### 6.2.2 Odyssey

The evaluation of the model using Homer's *Odyssey* followed the same scheme as when using *Romeo and Juliet*. The distribution of the metric scores is illustrated in Figure 8. Same as with the previous example, BERT-based metrics have much higher values than the others, BLEU being the lowest by far. The exact results can be seen in Table 3. Though similar, all metrics are slightly lower for *Odyssey* in comparison to *Romeo and Juliet*, meaning the model cannot extract information about it quite as well. One reason could be due to *Odyssey* being significantly longer or *Romeo and Juliet* being a drama, resulting in more dialogues and fewer descriptions, making the chunks more dense in sense of information and action, making it easier for the model to get all the necessary facts.

Once again, BERT F1 score value seems very promising, as at 0.86 it means that the model handles semantic reasoning well. The score might be slightly lower than the one for *Romeo and Juliet*, but with *Odyssey* being much longer a book, it still being that high is indicative of the RAG component performing well, even on large texts. Its distribution is presented in Figure 9. The distribution is similar to that for *Romeo and Juliet*, yet clearly shifted slightly to the left, aligning with the marginally lower value of the metric.
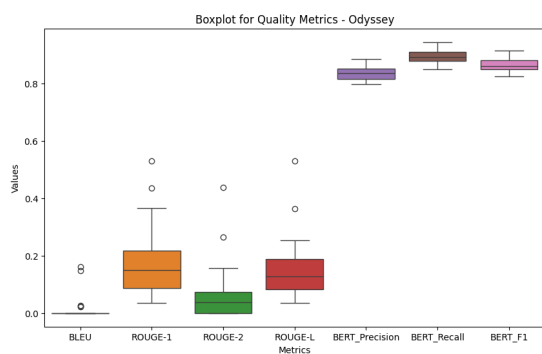


Figure 8: Distribution of Evaluation Metrics for *Odyssey*

Same as with the previous book, a few examples of best and worst in BERT F1 score question-answer pairs were chosen for further evaluation. Below, two examples of each are given.

**High score**

- **How long does Odysseus stay on Calypso's**

| Metric | Mean | Median | Std Dev |
|---|---|---|---|
| BLEU | 0.009 | 1.81e-155 | 0.03 |
| ROUGE-1 | 0.17 | 0.15 | 0.1 |
| ROUGE-2 | 0.06 | 0.04 | 0.08 |
| ROUGE-L | 0.15 | 0.13 | 0.09 |
| BERT Precision | 0.84 | 0.84 | 0.02 |
| BERT Recall | 0.89 | 0.89 | 0.02 |
| BERT F1 | 0.86 | 0.86 | 0.02 |

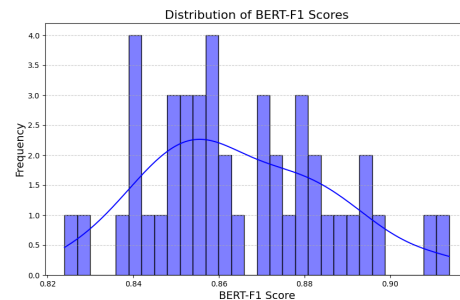Table 3: Performance Metrics for *Odyssey*



Figure 9: Distribution of BERT-F1 score for *Odyssey*

island?
**Answer:** Odysseus stays on Calypso's island for seven years.
**RAG Output:** According to Odysseus' account, he stayed on Calypso's island for seven years before she sent him on his way of her own free will.
**BERT F1:** 0.9102

- **What effect does the lotus plant have on Odysseus's men?**
**Answer:** The lotus plant makes them forget their desire to return home.
**RAG Output:** The lotus plant causes Odysseus's men to become indifferent to home and desire to stay and consume the plant without wanting to return, despite weeping bitterly. This shows the powerful and intoxicating effect the lotus has on those who consume it.
**BERT F1:** 0.8930

**Low score**

- **What is the fate of the Phaeacian ship that brings Odysseus home?**
**Answer:** Poseidon turns the ship into stone as punishment.

**RAG Output:** The Phaeacian ship that brings Odysseus home is rooted in the sea just as it is getting into port by the earth-encircling Neptune, who then drives it down with the flat of his hand so as to root it in the ground.
**BERT F1:** 0.8365

- **What is the omen that confirms Odysseus's return to Ithaca?**
  **Answer:** An eagle carrying a dead goose appears as an omen.
  **RAG Output:** Based on the passage, the omen confirming Odysseus's return to Ithaca is the sight of a hawk, Apollo's messenger, carrying a dove, which Theoclymenus interprets as a sign that Telemachus will remain powerful, and his house will be the most royal in Ithaca. However, this seems to be an interpretation of the omen for Telemachus rather than a direct confirmation of Odysseus's return. The actual omen for Odysseus's return is not explicitly stated in the given passages. The passage focuses more on the gathering of the gods and Telemachus's departure. The mention of the hawk and dove seems to be a separate event not directly tied to Odysseus's return. If an omen for Odysseus's return is needed, it would likely be inferred from the overall context of the passage, which suggests that the gods are favouring Odysseus's journey home despite
  **BERT F1:** 0.8398

Important to remember, the prompt passed to the RAG system consisted of "Use 1–2 sentences if possible.". In the second example, among the low scoring answers, the message was too long to fit in the defined length limit. For all four answers, it can be observed, that the system clearly provides too much redundant information on top of the actual answer to the asked question. It can be seen especially in the second lower rated answer, which provides the correct information at the very beginning, only to digress later on.

Lastly, a random sample of 10 outputs was again chosen and reviewed manually to see whether the passages found by RAG were correct. An example of such passage:

- **Question:** How does Juliet react when she discovers Romeo's identity?
  **RAG Output:** According to the passage, Odysseus chooses to navigate past Scylla, sacrificing six of his men to avoid the worse fate of being sucked into Charybdis' whirlpool, as the goddess Circe advised him. He drives his ship past Scylla at full speed and roars out to her, hoping to deter her from attacking again.
  **Fragments of RAG Passage:**

  > You will find the other rock lie lower, but they are so close together that there is not more than a bow-shot between them. [A large fig tree in full leaf grows upon it], and under it lies the sucking whirlpool of Charybdis. Three times in the day does she vomit forth her waters, and three times she sucks them down again; see that you be not there when she is sucking, for if you are, Neptune himself could not save you; you must hug the Scylla side and drive ship by as fast as you can, for you had better lose six men than your whole crew.
  > You dare devil,' replied the goddess, 'you are always wanting to fight somebody or something; you will not let yourself be beaten even by the immortals. For Scylla is not mortal; moreover she is savage, extreme, rude, cruel and invincible. There is no help for it; your best chance will be to get by her as fast as ever you can, for if you dawdle about her rock while you are putting on your armour, she may catch you with a second cast of her six heads, and snap up another half dozen of your men; so drive your ship past her at full speed, and roar out lustily to Crataiis who is Scylla's dam, bad luck to her; she will then stop her from making a second raid upon you.

The passage does not specify the number of people sacrificed by Odysseus, it mentions Circe's advice as to how he should navigate between the monsters, which Odysseus later does. It is not the

perfect choice, but shows the RAG identifies key phrases and words correctly and provides relevant passages.

It is important to note, that some passages contained information from another book in our dataset – Homer's *Iliad*. It is understandable, since some characters, like Odysseus, appear in both books. It might indicate an issue with RAG focusing on names and proper names rather than the actual action taking place in the fragments and its relevance to the questions.

# References

[Chase2022] Harrison Chase. 2022. Langchain, oct.

[Devlin et al.2019] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding.

[Dey et al.2023] Nolan Dey, Gurpreet Gosal, Zhiming, Chen, Hemant Khachane, William Marshall, Ribhu Pathria, Marvin Tom, and Joel Hestness. 2023. Cerebras-gpt: Open compute-optimal language models trained on the cerebras wafer-scale cluster.

[Douze et al.2024] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library.

[Es et al.2023] Shahul Es, Jithin James, Luis Espinosa-Anke, and Steven Schockaert. 2023. Ragas: Automated evaluation of retrieval augmented generation.

[Gao et al.2020] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. The pile: An 800gb dataset of diverse text for language modeling.

[Gao et al.2024] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024. Retrieval-augmented generation for large language models: A survey.

[Jiang et al.2023] Zhengbao Jiang, Frank Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7969–7992, Singapore, December. Association for Computational Linguistics.

[Li et al.2024] Xiaonan Li, Changtai Zhu, Linyang Li, Zhangyue Yin, Tianxiang Sun, and Xipeng Qiu.

2024. LLatrieval: LLM-verified retrieval for verifiable generation. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5453–5471, Mexico City, Mexico, June. Association for Computational Linguistics.

[OpenAI2024] OpenAI. 2024. Chatgpt (november 2024 version). https://chat.openai.com/. Accessed: 2024-11-27.

[Pan et al.2023] James Jie Pan, Jianguo Wang, and Guoliang Li. 2023. Survey of vector database management systems.

[Reimers and Gurevych2019] Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11.

[Saad-Falcon et al.2024] Jon Saad-Falcon, Omar Khattab, Christopher Potts, and Matei Zaharia. 2024. ARES: An automated evaluation framework for retrieval-augmented generation systems. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 338–354, Mexico City, Mexico, June. Association for Computational Linguistics.

[Sachan et al.2022] Devendra Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2022. Improving passage retrieval with zero-shot question generation. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3781–3797, Abu Dhabi, United Arab Emirates, December. Association for Computational Linguistics.

[Team2024] Qwen Team. 2024. Qwen2.5: A party of foundation models, September.

[Touvron et al.2023] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models.

[Wolf et al.2020] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Perric Cistac, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods*

*in Natural Language Processing: System Demonstrations*, pages 38–45, Online, oct. Association for Computational Linguistics. If you use this software, please cite it using these metadata.

[Xiao et al.2023] Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. 2023. C-pack: Packaged resources to advance general chinese embedding.

[Yang et al.2024] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

## A   Appendix: Division of work

The project requirements demand to provide workload for each team member at any stage of the project. Therefore, in order to comply with that, we list the contributions of each team member so far below:

- **Maja Andrzejczuk:** Abstract, Introduction, Collecting the data, PoC, Reviews, System evaluation;

- **Piotr Bielecki:** Datasets & Work plan sections, Collecting the data, PoC, Reviews, System evaluation;

- **Paweł Gelar:** Models & Vector database sections, PoC, Collecting the data, Reviews, Final Pipeline;

- **Jakub Kasprzak:** Finding datasets & models, Work plan section, Collecting the data, EDA, Reviews, System evaluation;

- **Maciej Orsłowski:** Related works, Collecting the data, EDA, Reviews, Exporting database.

We state that in terms of workload and time assessment, each member has contributed equally to the project.

## B   Appendix: Wide plots

In this appendix we present the figures that we refer to in the text, which were too wide to nicely place in sections where they are mentioned.
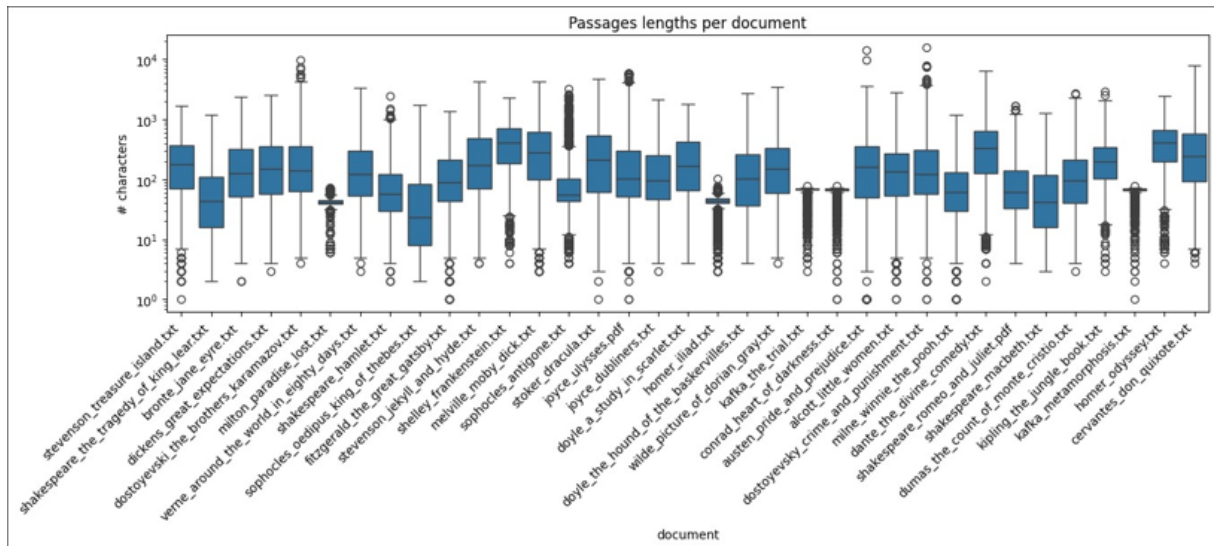
Figure 10: Boxplots of passage lengths per document before improvements. Outliers highlight extreme variability in some documents.
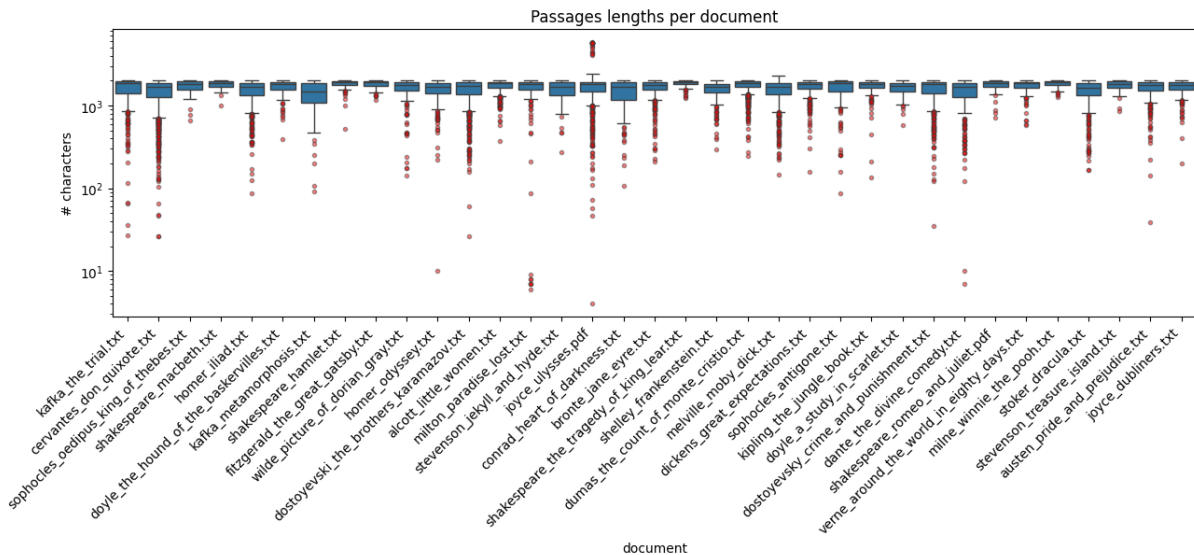


Figure 11: Boxplots of passage lengths per document after improvements. Variability across documents has been significantly reduced.