



Conversational Agent with Retrieval-Augmented Generation

Katarina Velkov, Nejc Krajšek, Luka Sabotič

Abstract

abstract

Keywords

Conversational agent, Retrieval-Augmented Generation

Advisors: Aleš Žagar

Introduction

Large language models (LLM) rely solely on their pre-trained knowledge to generate responses. These models have proved to be powerful but they are prone to generating hallucinated or inaccurate information. Furthermore their use in specialized fields has proved to be challenging as LLMs are trained mostly on publicly available data and are not updated after new data arises. That is why the need for more advanced models has emerged. Retrieval-Augmented Generation (RAG) enables LLMs to retrieve more relevant information from various external databases and web sources during conversations. This process allows for more accurate, domain-specific and up-to-date responses.

The aim of this project is to build a conversational agent using the RAG technique. We will use an existing large language model, expand it to be able to query linked domain-specific databases. We will evaluate the performance of our model by comparing it to the original LLM to discover improvements our implementation brought and also other RAGs available online to inspect the overall quality of our work.

Related work

The concept of Retrieval-Augmented Generation (RAG) was first introduced in 2020 by [1], who proposed a framework that combines the strengths of large language models (LLMs) with external knowledge retrieval. Since then, RAG has seen significant advancements, particularly in the context of domain-specific applications and conversational agents.

Since we are focusing on RAG of code and code-based answers, we also took a look at [2]. The authors explain how natural language RAG differs from code-based applications. They emphasize chunking as the crucial difference and explain how it should be done to preserve context such as import

statements and dependencies across all chunks of the code file. [3] provide open source code for their chunking algorithm and other RAG components as are used in the Jet Brains Sweep AI coding assistant. Here, the chunking algorithm differs quite a lot from the standard approach for natural language RAG. Programming itself is quite prone to syntax errors, not to mention LLM generated code. This is why [4] provide a RAG pipeline that performs dependency and syntax unit tests on the generated code to ensure correctness of the output, which we see as a good idea for our project.

Another relevant work is the survey titled "Retrieval-Augmented Generation for Large Language Models: A Survey," [5], which provides a comprehensive overview of RAG's evolution. The survey categorizes RAG into three paradigms: Naive RAG, Advanced RAG, and Modular RAG. It highlights the importance of retrieval, generation, and augmentation techniques in enhancing the performance of RAG systems. The survey also discusses the challenges and future directions of RAG, such as improving robustness, handling long contexts, and integrating multimodal data.

Methods

We have developed a retrieval-augmented generation to provide relevant context to a conversational agent, which then produces the answer to the user query. We are focusing on the programming and debugging domain, trying to create a model that will correctly generate, complete and correct the given programming code and answer other programming-related questions.

We used Stack Overflow [6] question-answer pairs as our context database, for now only focusing on python-related content due to the computational complexity of the pipeline. Later, if resources allow, we are going to expand to other

languages as well.

Data preparation

As mentioned, we are using Stack Overflow data, specifically extracted python tagged questions from Google Big Query's public Stack Overflow dataset. We are not only interested in question and answer contents, but also tags and Stack Overflow scores, for pre- and post-retrieval processing.

Chunking

We split each answer-question pair into multiple chunks. For each chunk, we save the text to be embedded and also some metadata that will help at the retrieval step. For a question-answer pair, we embed: the question with the answer stored as metadata, the answer, and all the code blocks we find in answers. We don't use code from the question texts, as it often contains errors and is thus in the question body. We use a very simple code identifier function that claims something is code if it is inside '<code>' blocks, longer than 10 characters and contains at least two lines, as code usually spans across multiple lines.

Embedding

For the embedding, we are using the all-MiniLM-L6-v2 [7], which embeds text into a 384-dimensional dense vector space. Embeddings are then stored in a vector database, provided by the Chroma [8] library. It automatically embeds and stores the chunks. It also provides easy retrieval by only calling the get function with the query, distance function and the number of wanted results.

Retrieval

We retrieve context chunks by embedding the user query and returning the 10 closest results by cosine distance. For the chunks that represent questions, we return the answer instead of the question, as this is the more relevant information for the user; we want to answer their question, not repeat what they asked. Code and answers are just returned.

Reranking

After retrieval, the results are reranked using the cross encoder model bge-reranker-base [9]

Augmentation

When we have both the user question and the relevant contexts, we build a prompt that will be passed to the LLM. To do this, we use the following function:

```
1 def build_prompt(self, query, results):
2     contexts = self.context_from_results(results)
3     return f'''
4     Answer the following code related question
5     using the context provided inside triple
6     quotes in it is useful.
7     In the answer provide an example of code that
8     is related to the question.
9     If you do not know the answer, say that you do
10    not know. Do not try to invent the solution.
```

```
8
9     Question: {query}
10
11
12    '''"".join(f"Context {i}: {context}{chr(10)}{
13    chr(10)}" for i, context in enumerate(contexts
14    ))'''
15
16    Answer:
17    '''
```

This is a basic prompt with the most important instructions. Further prompt engineering is a subject of later work, when more important parts are finished.

Generation

We chose the StableLM 2 Zephyr 1.6B [?] as our LLM. This was because it is a poorly performing model for programming and mathematics tasks and we hope for significant differences in the results between using the basic model query and our RAG injection.

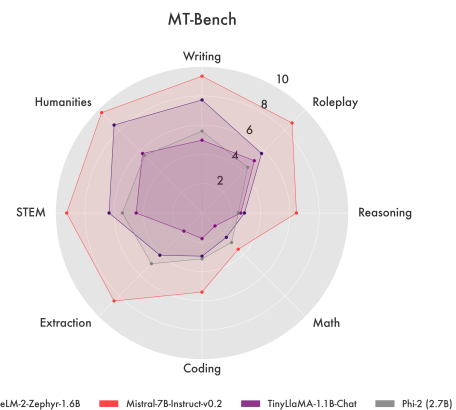


Figure 1. Performance of the selected model compared to other similar-sized models.

Results

For evaluation of our Retrieval-Augmented Generation (RAG) model for assisting with Python coding queries, we made a human evaluation and quantitative evaluation using ROUGE and BERTScore evaluation metrics. We created three sets of queries, first of 14 simple questions, second of 13 semi-advanced instructions and third of 7 advanced tasks. The first set is meant to test model's ability to answer and provide code snippets on simple coding questions. The second set is meant to test the model's ability to provide more complex answers and implement specific algorithms, such as implementing a simple calculator and check for palindromes. The third set is meant to test advanced tasks, such as implementing a machine learning model.

The generated answers using the RAG model were processed as to remove "Incorrect answer" parts, remove mark-down (eg. "python) and the whitespace was normalized. In

the evaluation step, the answers were compared to the ground truth answers, which were processed in the same way. The results of the evaluation are shown in the Table ??.

The human evaluation results were obtained by 3 evaluators (authors), who evaluated the answers on a scale from 1 to 5, regarding 7 predetermined categories. The results are shown in the Table ??.

Discussion

Acknowledgments

References

- [1] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc., 2020.
- [2] Tal Sheffer. Rag for a codebase with 10k repos. <https://www.qodo.ai/blog/rag-for-large-scale-code-repos/>, 2024. Accessed: 01-05-2025.
- [3] JetBrains. Sweepai. <https://github.com/sweepai/sweep>, 2023. Accessed: 01-05-2025.
- [4] Inc. LangChain. Code generation with rag and self-correction. https://langchain-ai.github.io/langgraph/tutorials/code_assistant/langgraph_code_assistant/, 2025. Accessed: 01-05-2025.
- [5] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A survey, 2024.
- [6] Stack Exchange Inc. Stack overflow. <https://console.cloud.google.com/marketplace/product/stack-exchange/stack-overflow?inv=1&inv=AbwUPw>, 2025. Accessed: 28-04-2025.
- [7] Sentence Transformers. all-minilm-l6-v2.
- [8] Chroma. Chroma - the open-source embedding database. <https://github.com/chroma-core/chroma/tree/main/>, 2025. Accessed: 01-05-2025.
- [9] BAAI. Bge reranker base.