

The Impact of Embedding Methods on Faithfulness of Retrieval-Augmented Generation (RAG)

Amandeep Kaur Singh
amandeepkaur@umass.edu
University of Massachusetts
Amherst, USA

Chang Jin
changjin@umass.edu
University of Massachusetts
Amherst, USA

Reagan Keeney
rkeeney@umass.edu
University of Massachusetts
Amherst, USA

ACM Reference Format:

Amandeep Kaur Singh, Chang Jin, and Reagan Keeney. 2024. The Impact of Embedding Methods on Faithfulness of Retrieval-Augmented Generation (RAG). In *Proceedings of Informaition Retrieval (CS646 FA2024)*. ACM, New York, NY, USA, 5 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Abstract

Retrieval Augmented Generation (RAG) is a promising field at the intersection of LLMs and information retrieval which aims to improve the quality and accuracy of LLM responses by providing supplemental context from outside the model's knowledge base. This paper evaluates the faithfulness of the outputs of three RAG pipelines with different embedding models for domain-specific biomedical question-answering on a subset of the RAG-Mini-BioASQ dataset. We observe that the RAG pipeline with the general purpose baseline bge-small-en-v1.5 embedding model achieved the highest faithfulness score on our subset of the Mini RAG BioASQ dataset (97%), outperforming both the fine-tuned MedEmbed-small-v0.1 (93%) and the sentence-transformer all-MiniLM-L6-v2 (92%). These findings highlight the adaptability of the general-purpose models and underscore the importance of embedding selection in optimizing RAG pipelines for faithfulness. Future work should explore larger datasets, a wider variety of models (fine-tuned and general purpose), and additional metrics for evaluation.

2 Problem Statement

Language models are a study topic with a long history, dating as early as 1950. Since then, language models have been applied in dynamic tasks such as information retrieval, speech recognition, and more. Large Language Models (LLMs) refer to language models that are trained on vast volumes of data, making them capable of answering questions in broad contexts and offering impressive natural language generative capabilities [11]. LLMs are designed to generate original output for tasks such as answering questions, translating between languages, and completing sentences. Despite their capabilities, LLMs are not without flaws, due to their text-generative nature and their need to answer a wide range of queries, LLMs are vulnerable to 'hallucinations' [10], where the model incorrectly answers questions with false information instead of reporting

a lack of adequate topic knowledge. Such hallucinations can be dangerous, as they may propagate misinformation and lead to poor decision-making in critical applications such as healthcare, law, and more.

Retrieval-augmented generation (RAG) is a technique that optimizes a large language model's output, utilizing information retrieval to consult documents from a reliable knowledge base outside of its training data sources before generating a response using an LLM. Ensuring these LLM-generated responses are accurate, pertinent, and helpful is important. Current research faces challenges in assuring the faithfulness of the generated responses; that is, that the output stays correct and in line with the acquired data—a significant problem with RAG[5].

The selection of embedding techniques, which translate textual data into numerical vectors used to compute similarities and retrieve pertinent documents, is a key determinant of a Retrieval-Augmented Generation model's faithfulness and performance. Different embedding strategies can lead to drastically different retrieval accuracy, quality, and reliability of the generated responses. For instance, dense embeddings like BERT can capture semantic relationships, improving relevance and potentially maintaining faithfulness [2]; sparse embeddings, on the other hand, yield faster results but may sacrifice contexts, which can lead to less faithful responses [8].

Given these variations in embedding models, optimizing RAG for applications that need high factual consistency and alignment is important. Our study thus attempts to answer the following question: **Which embedding technique best maintains faithfulness in responses produced by RAG models on a domain-specific dataset?** We hope that our findings will support applying RAG for high-stake applications where maintaining faithfulness and building user trust is crucial.

3 Motivation

Retrieval-augmented generation (RAG) has shown significant promise in enhancing the performance of large language models with the integration of retrieved information. The performance of RAG models can vary depending on the choice of embedding models used for information retrieval[4].

As such, while retrieval-augmented generation (RAG) offers a solution to some of the major shortcomings of Large Language Models (LLMs) for question-answering tasks, optimizing this pairing of systems remains inadequately explored. By systematically evaluating and comparing the faithfulness of RAG with different embedding models for a domain-specific dataset, we hope to gain valuable insights into how the choice of embeddings impacts the quality of the generated responses.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CS646 FA2024, December 5, 2024, Amherst, MA

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM

<https://doi.org/XXXXXXX.XXXXXXX>

4 Related Work

Lewis et. al introduced Retrieval-Augmented Generation as a novel approach combining information retrieval (IR) with text generation for knowledge-intensive Natural Language Processing (NLP) tasks in 2020 with their paper *Retrieval-augmented Generation for Knowledge-Intensive NLP Tasks* [10]. The proposed approach includes two types of memories:

- (1) Non-parametric Memory: Document index which can be 'hot-swapped' in order to update knowledge as needed.
- (2) Parametric Memory: Parameters of a generative model such as BART.

A key finding of this paper was that the proposed RAG models perform better and tend to be more 'grounded' i.e., they generate "more specific, diverse and factual language"[10, p. 1] than the then state-of-the-art parametric-only seq2seq baseline chosen for the study (BART). Lewis et. al concluded their findings with a discussion of the broader impact of their work, suggesting applications for medical information retrieval. This is particularly relevant for our project, where we evaluated the impact of the choice of embedding models on the performance of RAG by applying it to the RAG-Mini-BioASQ dataset.

More recently, Şakar and Emekci[15] performed 23,625 grid-search iterations to obtain results comparing several RAG architectures and techniques. The embedding models compared included OpenAI's text-embedding-v3-large model, Beijing Academy of Artificial Intelligence's (BAAI) open-source bge-en-small model, and Cohere's cohere-en-v3 model. As discussed in the results section of [15], the choice of embedding models significantly affects the performance of the RAG systems in terms of hardware utilization (runtime, CPU usage, and memory usage), and performance (evaluated with median similarity scores). While Şakar and Emekci's research included a direct comparison of embedding models, the focus was on efficiency with regard to runtime and hardware requirements, not on the quality of generated outputs. Which embedding model is most optimal seems to vary by use case [4, p.9], leaving ambiguity on the correct conditions for a given model.

Evaluation in and of itself presents a challenge for RAG systems. Yu et.al note that the interplay of retrieval and generation tasks compounded with vague criteria for what constitutes a high-quality response result in uncertainty as to the most effective evaluation metric[18]. By extension, a single evaluation metric is unlikely to encompass all components and criteria in RAG systems.

For this research, we focus on 'faithfulness' as a measure of generated output quality. In broad terms, faithfulness is the level of agreement between the generated text and the contexts of all claims[1]. This makes faithfulness a useful tool to gauge to what extent the LLM hallucinates when generating responses[7]. This can provide information on how retrieved contexts from different embedding models are ingested and incorporated into responses, similar to how modification of a prompt's text can create consistent change in LLM output.

5 Approach

5.1 Overview of RAG Pipeline

Our Retrieval Augmented Generation pipeline begins by embedding passages using one of three models:

- (1) Baseline: BGE-SMALL-EN-v1.5,
- (2) Fine-tuned BERT Model MEDEMBED-SMALL-v0.1,
- (3) or ALL-MINILM-L6-v2

The embeddings are then stored in the Milvus vector database. When a query or question (information need) is passed in, it is first embedded, and then we retrieve documents from the RAG-Mini-BioASQ corpus using Milvus, leveraging vector similarity search. Then, we construct a prompt with the retrieved passages as context, and finally, Mistral AI's Mixtral-8x7B-Instruct-v0.1 Large Language Model generates natural language answers to the prompt based on the retrieved documents. Referenced documentation[19] to build the pipeline.

Detailed overview of the pipeline:

- (1) Embedding the passages in the corpus (documents to retrieve from) using one of the three embedding models.
- (2) Save the embeddings in the Milvus Vector Database: Milvus is a high-performing, open-source vector database optimized for efficient/fast, and scalable similarity searches on high-dimensional data [11, 17, 19]. It enables efficient retrieval of the most relevant documents for each query.
- (3) Embedding the query and retrieving: For our experiments so far, we return the top three closest documents using the inner product as the similarity metric.
- (4) Prompt Construction for LLM: Using the original question and the three retrieved documents as context, the prompt for the LLM model is formulated/structured. This step is crucial for ensuring that the generated answer is grounded in the retrieved evidence.
- (5) Mistral AI LLM: For the final step of the RAG pipeline, we utilize MistralAI's Mixtral-8x7B-Instruct-v0.1 language model[6, 19]. This model was selected due to its capability to handle long prompts effectively and generate coherent, contextually aware responses, making it well-suited for our retrieval-augmented generation (RAG) setup.

The only difference between the three pipelines is the chosen embedding model, all other variables are held constant.

5.2 Baseline

As our baseline, we selected the Beijing Academy of Artificial Intelligence's open-source model BGE-SMALL-EN-v1.5[14], which is the current version of the top performing model in Şakar and Emekci's resource efficiency comparison [15]. This model serves as a baseline for evaluating the performance of other models in our study.

5.3 Models

We selected two additional models based on their rankings and parameter size relative to the baseline, using the METB leaderboard as a reference [12].

5.3.1 Baseline: BAAI/bge-small-en-v1.5 (rank 70 on MTEB). As introduced earlier, this was the Beijing Academy of Artificial Intelligence's open model, a top-performing model in Şakar and Emekci's resource efficiency comparison [15]. This model ranks 70 in the English Retrieval section of the METB and serves as our baseline for evaluating embedding techniques.

5.3.2 Model 2: abhinand/MedEmbed-small-v0.1(rank 58 on MTEB). MedEmbed is a model fine-tuned for medicinal and clinical data, intended for use in medical and clinical contexts. We selected this model to investigate whether its domain-specific design will improve the faithfulness of the responses generated. Its size is comparable to the baseline model but ranks higher in the English Retrieval section of METB.

5.3.3 Model 3: sentence-transformer/all-MiniLM-L6-v2(rank 133 on MTEB). This model is a sentence-transformer model that performs well with semantic search, clustering, and sentence similarity [16]. This model presents an advantage due to its' small size, similar to that of the baseline model. In the English Retrieval section of METB, this model ranks lower than our baseline at 133.

6 Experiments

6.1 Dataset(s)

For the purpose of our study, we used a subset of the BioASQ 11b dataset [13] in the domain of biomedical question answering. Specifically, 100 questions from the pre-generated RAG-Mini-BioASQ dataset[3], which is designed for Retrieval-Augmented Generation use cases.

6.1.1 Key Features of the RAG-Mini-BioASQ dataset:

- (1) Subset of BioASQ dataset: Given the limited resources and timeline of the project, using a smaller dataset is advantageous for ensuring the scope is achievable while still retaining the complexity of the original dataset.
- (2) Structured for RAG-based approaches/experiments: The RAG-Mini-BioASQ dataset already has a correct response to the question, which could be useful for evaluation purposes. The dataset consists of the following:
 1. Text-Corpus: Consists of 40181 biomedical text passages and their IDs.
 2. Question-Answer-Passages: Consists of 4719 questions (representing a user's biomedical information need), a correct answer to those questions, and a list of relevant passage IDs from the Text-Corpus data. For our application, only the questions are utilized.
- (3) Biomedical focus: The chosen dataset is domain-specific; enabling us to perform a more accurate assessment of various embedding models for this specific use case.
- (4) Complexity: The questions are designed to be particularly challenging[9], which may make disparities between the faithfulness of different embeddings more apparent.

In critical domains such as biomedical question-answering, hallucinations from generative language models can raise significant safety concerns. We chose this dataset as it is uniquely important to understand, evaluate, and ensure faithfulness to verify correct reference information in such fields.

6.1.2 Data Considerations: Given that the chosen dataset is already processed and designed for RAG tasks, it does not require extensive additional pre-processing. Additionally, due to computing limitations, we sampled a subset of 100 questions from the dataset and evaluated the RAG with the different embedding models on those queries only.

6.2 Evaluation

Given that our research centers on the quality of the generated text, so do our evaluation metrics. We used the RAG-specific evaluation of faithfulness, as defined in the DeepEval framework [1]. This metric aims to gauge how successfully the LLM incorporates context documents into its generated answer. To measure this, another "evaluator" LLM is utilized to identify claims made within both the answer and context. From there, the number of truthful (according to context passages) claims is divided by the total number of claims within the response.

We employed faithfulness as a metric to monitor the degree to which the LLM integrates niche knowledge from the corpus to produce correct answers. A high faithfulness score would represent both integration and a lack of hallucination, both of which are important to ensure trustworthy models. Our specific tool for measuring faithfulness is from DeepEval's[1] package and employs a 4-bit quantized MISTRAL-7B-v0.3 as the evaluator LLM. The evaluator is provided the RAG-generated output, retrieved-context passages, and the original question for analysis. All parameters in the evaluator (such as passing threshold) are kept at package default.

6.3 Results

Results on a 100-question subset are included below.

Model	Faithfulness
Baseline: bge-small-en-v1.5	97%
MedEmbed-small-v0.1	93%
all-MiniLM-L6-v2	92%

Table 1: Embedding Model Performance

6.4 Analysis of Results

The chosen general-purpose models (baseline model bge-small-en-v1.5, and all-MiniLM-L6-v2) displayed varying levels of adaptability to the task of question answering within the biomedical domain of the RAG-Mini-BioASQ dataset; when evaluated for faithfulness. The bge-small-en-v1.5 achieved the highest faithfulness scores for our experiments (97%), with MedEmbed-small-v0.1 and all-MiniLM-L6-v2 scoring 93% and 92% respectively. Counterintuitively, the fine-tuned domain-specific MedEmbed-small-v0.1 embedding model did not outperform the general use baseline model.

6.5 Limitations and Considerations

The system, domain, and evaluation method described in this paper are highly complex and nuanced. We make note of several points for the reader to be aware of about the results and implications of this work.

(1) Biomedical Domain:

Given that the data is from a biomedical domain, it contains technical terms, abbreviations, et cetera that can make information retrieval and/or question-answering tasks difficult. Additionally, as discussed earlier, credibility/reliability is uniquely essential within this domain as errors in medical advice/diagnoses could have significant consequences.

(2) Evaluating RAG Faithfulness across Embedding Models:

The performance of each RAG pipeline is influenced by not only the embedding model but also by the LLM chosen to generate the final natural language responses. For our experiments, we used MistralAI's Mixtral-8x7B-Instruct-v0.1 language model, which could have impacted the final outputs of the models; thereby affecting the faithfulness scores of all of the models. Additionally, for our experiments, the context passed into the LLM included up to three relevant documents (retrieved from the Milvus database based on the scores calculated using the inner product), this can be improved in various ways; for example, we could increase the number of documents passed in (which could help improve diversity but also potentially increase the noise) and/or set a minimum relevance score threshold for a document to be passed in, et cetera. Further evaluation with alternative LLMs and careful prompt engineering could lead to new results/insights.

(3) Resource Constraints

Due to computational constraints and time limitations, our experiments were restricted to a subset of the dataset (100 questions) and three open-source compact models. Expanding this scope to include other larger, non-open-source models or evaluating with larger datasets would likely provide a more robust evaluation.

(4) Trade-Offs:

Faithfulness is not an all-encompassing metric, and cannot give a holistic view of response quality. Besides faithfulness, it is also important to evaluate other metrics such as diversity, inference time, et cetera for practical industry applications.

6.6 Conclusion

Overall, these results highlight the strengths and potential for application of both general-purpose and fine-tuned models for domain-specific tasks; such as our experiment focusing on biomedical question-answering. From our experiments, we observed that the RAG pipeline with the general purpose baseline bge-small-en-v1.5 embedding model achieved the highest faithfulness score on our subset of the Mini RAG BioASQ dataset. Although the fine-tuned MedEmbed-small-v0.1 performed well (achieving a faithfulness score of 93%), the bge-small-en-v1.5 model's higher score showcases the adaptability and robustness of the BAAI General Embedding models.

However, as discussed in the previous section, in order to gain a more comprehensive understanding of the models' embedding capabilities, future research should consider exploring the following:

- Larger and more diverse datasets within the domain.
- RAG pipelines with alternate LLMs and careful prompt engineering.
- Comparisons with larger and/or more fine-tuned domain-specific models; potentially even non-open-source fine-tuned models to better understand the potential of applying to domain-specific applications.
- Additional metrics to evaluate trade-offs.

There is significant potential to apply Retrieval-Augmented Generation (RAG) models to the task of biomedical question answering. Developing credible and fine-tuned RAG systems with robust citation mechanisms could assist healthcare professionals by providing

preliminary evidence-based insights (which they can verify by either using their domain expertise or reviewing the cited documents). Additionally, it could potentially help enable improved patient education and help accelerate biomedical research by synthesizing large volumes of complex data for medical research professionals. By evaluating the faithfulness and domain adaptation of embedding models, we hope that this research contributes to the long-term goal of creating reliable and efficient question-answering tools for domain-specific fields, such as the biomedical field.

7 Miscellaneous

[A] The code used for this analysis and results are available via **clicking here** or going to <https://github.com/jinc0930/RAG-Faith-Embed>

7.1 Contributions

All authors contributed equally in this project, we list their contributions below.

Amandeep: Contributed to coding the exploratory data analysis and basic data loading/setting-up. Contributed to general RAG pipeline setup with baseline embeddings, MedEmbed embeddings, and preliminary testing with dummy data. Heavily contributed to late-stage writing and editing, including abstract, literature review, sections 5.1, 6.1, 6.4, 6.5, and 6.6.

Chang: Contributed heavily in coding, especially "Generate Answers" and part of the "Faithfulness" sections. Tested and run embedding with a subset of real data. Contributed to writing, specifically the introduction and part of the problem statement. Contributed to general editing of writing and formatting.

Reagan: Contributed to coding in the evaluation section. Ran eval for all data. Contributed heavily to sections 2(along with Chang),3, 5.2, and 5.3. Responsible for selecting dataset, evaluation model, and identifying MTEB leaderboard for embedding model selection. Responsible for all LaTeX citations, formatting, and style.

References

- [1] [SW] Confident AI, DeepEval version 1.4.7, Nov. 21, 2024. URL: <https://docs.confident-ai.com/>.
- [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Jill Burstein, Christy Doran, and Thamar Solorio, (Eds.) Association for Computational Linguistics, Minneapolis, Minnesota, (June 2019), 4171–4186. doi: 10.18653/v1/N19-1423.
- [3] enelpol. 2024. Rag-mini-bioasq. (2024). <https://huggingface.co/datasets/enelpol/rag-mini-bioasq>.
- [4] Yunfan Gao et al. 2023. Retrieval-augmented generation for large language models: a survey. *arXiv preprint arXiv:2312.10997*.
- [5] Or Honovich et al. 2022. True: re-evaluating factual consistency evaluation. (2022). <https://arxiv.org/abs/2204.04991> arXiv: 2204.04991 [cs.CL].
- [6] Albert Q Jiang et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- [7] Xiaonan Jing, Srinivas Billa, and Danny Godbout. 2024. On a scale from 1 to 5: quantifying hallucination in faithfulness evaluation. *arXiv preprint arXiv:2410.12222*.
- [8] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, (Eds.) Association for Computational Linguistics, Online, (Nov. 2020), 6769–6781. doi: 10.18653/v1/2020.emnlp-main.550.

- [9] Anastasia Krithara, Anastasios Nentidis, Konstantinos Bougiatiotis, and Georgios Paliouras. 2023. Bioasq-qa: a manually curated corpus for biomedical question answering. *Scientific Data*, 10, 1, 170.
- [10] Patrick Lewis et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems (NIPS '20)* Article 793. Curran Associates Inc., Vancouver, BC, Canada, 16 pages. isbn: 9781713829546.
- [11] Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: a survey. *arXiv preprint arXiv:2402.06196*.
- [12] Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2023. Mteb: massive text embedding benchmark. (2023). <https://arxiv.org/abs/2210.07316> [cs.CL].
- [13] n.d. 2023. Bioasq training 11b. (2023). <http://participants-area.bioasq.org/datasets/>.
- [14] Beijing Academy of Artificial Intelligence. [n. d.] Bge-small-en-v1.5. (). <https://huggingface.co/BAAI/bge-small-en-v1.5>.
- [15] Tolga Şakar and Hakan Emekci. 2024. Maximizing rag efficiency: a comparative analysis of rag methods. *Natural Language Processing*, 1–25.
- [16] Rahultiwari Tiwari. 2024. Unlocking the power of sentence embeddings with all-minilm-l6-v2. (2024). <https://medium.com/@rahultiwari065/unlocking-the-power-of-sentence-embeddings-with-all-minilm-l6-v2-7d6589a5f0aa>.
- [17] Jianguo Wang et al. 2021. Milvus: a purpose-built vector data management system. In *Proceedings of the 2021 International Conference on Management of Data*, 2614–2627.
- [18] Hao Yu, Aoran Gan, Kai Zhang, Shiwei Tong, Qi Liu, and Zhaofeng Liu. 2024. Evaluation of retrieval-augmented generation: a survey. (2024). <https://arxiv.org/abs/2405.07437> arXiv: 2405.07437 [cs.CL].
- [19] Cheney Zhang. 2024. Build rag with milvus. (2024). https://github.com/milvus-io/bootcamp/blob/master/bootcamp/tutorials/quickstart/build_RAG_with_milvus.ipynb.