# Efficient and Reproducible Biomedical Question Answering using Retrieval Augmented Generation

Linus Stuhlmann*, Michael Alexander Saxer*, Jonathan Fürst

School of Engineering, Zurich University of Applied Sciences, Winterthur, Switzerland

E-mails: linus.stuhlmann@zhaw.ch, michael.saxer@zhaw.ch jonathan.fuerst@zhaw.ch

*Equal contribution.

*Abstract*—**Biomedical question-answering (QA) systems require effective retrieval and generation components to ensure accuracy, efficiency, and scalability. This study systematically examines a Retrieval-Augmented Generation (RAG) system for biomedical QA, evaluating retrieval strategies and response time trade-offs. We first assess state-of-the-art retrieval methods, including BM25, BioBERT, MedCPT, and a hybrid approach, alongside common data stores such as Elasticsearch, MongoDB, and FAISS, on a $\approx$ 10% subset of PubMed (2.4M documents) to measure indexing efficiency, retrieval latency, and retriever performance in the end-to-end RAG system. Based on these insights, we deploy the final RAG system on the full 24M PubMed corpus, comparing different retrievers' impact on overall performance. Evaluations of the retrieval depth show that retrieving 50 documents with BM25 before reranking with MedCPT optimally balances accuracy (0.90), recall (0.90), and response time (1.91s). BM25 retrieval time remains stable (82ms ± 37ms), while MedCPT incurs the main computational cost. These results highlight previously not well-known trade-offs in retrieval depth, efficiency, and scalability for biomedical QA. With open-source code, the system is fully reproducible and extensible.**

*Index Terms*—**Biomedical Information Retrieval, Retrieval-Augmented Generation, Hybrid Retrieval, Large Language Models, PubMed, Information Retrieval Systems.**

## I. INTRODUCTION

Large Language Models (LLMs) have demonstrated strong biomedical question-answering (QA) capabilities [1]. However, LLMs can produce factual inaccuracies, lack specific domain knowledge, and lack verifiability [2]. A major concern is *hallucination*, where LLMs generate factually incorrect responses due to their probabilistic nature. These hallucinations, together with a lack of verifiability, are particularly problematic in healthcare, where misinformation can lead to serious consequences. To mitigate these risks, *Retrieval-Augmented Generation (RAG)* systems leverage external knowledge sources at inference time by selecting relevant documents from a data store to enhance accuracy, transparency, and traceability [3].

Despite the potential of biomedical RAG systems, existing solutions often suffer from limited scalability, poor reproducibility, and suboptimal retrieval performance on large datasets such as PubMed. Existing benchmarks for medical question answering, such as MedExpQA [4] and MIRAGE [5], lack reproducibility and scalable retrieval solutions. Most

retrieval methods rely on either *sparse* bag-of-words vectors such as BM25 [6] or *dense* vectors created through transformer-based models such as BioBERT [7] and MedCPT [8]. However, **hybrid approaches that integrate both techniques remain under-investigated, especially from a system perspective**: the inherent trade-offs between retrieval strategies, their indexing and response times, and the resulting generator accuracy are crucial for practical RAG applications and have been largely unexplored. Hybrid retrieval methods combine the strengths of sparse and dense retrieval: a *probabilistic retriever* (e.g., BM25) efficiently reduces the search space by filtering a large corpus, while a *neural reranker* (e.g., MedCPT's cross-encoder) refines document rankings based on semantic relevance. This two-step approach balances computational efficiency and retrieval precision, mitigating the limitations of stand-alone methods. Although hybrid retrieval has been explored in general NLP tasks [9], *its application in large-scale biomedical QA remains limited, particularly in real-world implementations*. Developing an effective biomedical QA system requires addressing several challenges: (i) **Efficient Retrieval at Scale**: Processing millions of biomedical documents under reasonable *indexing times* while maintaining *low-latency retrieval*; (ii) **Relevance Optimization**: Improving *document ranking* by integrating lexical retrieval with neural reranking; (iii) **Context Integration**: Structuring retrieved documents effectively to generate *factually accurate and verifiable responses*. This work presents a **scalable and reproducible** RAG system for biomedical QA, systematically evaluating hybrid retrieval strategies. The key contributions include:

- **Hybrid Retrieval Approach**: A two-stage retrieval pipeline that integrates BM25 (lexical retrieval) with MedCPT's cross-encoder (semantic reranking), improving recall and precision.
- **Scalability and Performance Analysis**: Comparative evaluation of three common methods and systems, MongoDB, Elasticsearch and FAISS, for large-scale document retrieval efficiency.
- **Reproducibility and Transparency**: Explicit citation of retrieved documents using PubMed IDs to ensure traceability in biomedical QA.

This work advances scalable and reproducible biomedical QA systems, enhancing their real-world applicability in clinical and research environments. All our code is open-sourced.[1]

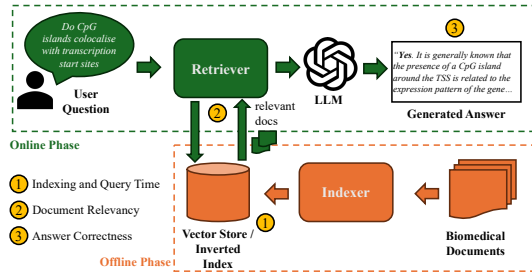## II. BIOMEDICAL QUESTION-ANSWERING WITH RAG



Fig. 1: Biomedical Question Answering with Retrieval-Augmented Generation (RAG). Offline phase: biomedical documents are processed, indexed, and stored in a vector store. Online phase: users ask questions, for which a retriever retrieves relevant documents that are appended with the question and fed to an LLM. Based on the question and context, the LLM generates an answer the PubMed IDs it used.

Retrieval-Augmented Generation enhances the capabilities of LLMs by incorporating external data and grounding responses in verifiable and up-to-date information. This ensures that outputs incorporate relevant biomedical knowledge from sources such as PubMed [10], improving accuracy and transparency. Figure 1 illustrates our biomedical QA system based on RAG, which consists of two main phases: an offline phase for indexing biomedical literature and an online phase for retrieving relevant documents and generating responses.

In the offline phase, biomedical documents are preprocessed and indexed into a vector store (dense vectors) and/or an inverted index (sparse vectors) for efficient retrieval. The choices directly affect retrieval speed and system scalability.

During the online phase, users submit biomedical queries, which the retriever processes to fetch relevant documents. These retrieved documents, along with the query, are provided as context to an LLM, ensuring responses remain grounded in authoritative biomedical sources. The used sources are cited (PMIDs) and provided as references to users. In this setting, system performance should be evaluated based on three key aspects, as highlighted in Figure 1: (1) indexing and query time, which measures retrieval efficiency; (2) relevance of retrieved documents, ensuring the most informative sources are selected; and (3) answer correctness, verifying that the LLM-generated response aligns with biomedical evidence. Evaluating and optimizing these factors improves the reliability and transparency of biomedical QA systems.

## III. EXPERIMENTAL EVALUATION

We evaluate the efficiency and effectiveness of different retrieval and text-generation methods for biomedical question-answering (QA). Our experiments focus on selecting optimal

components for document retrieval, text generation, and overall system performance following the three key aspects from Section II and using a common biomedical QA benchmark.

### A. Experimental Setup

**Datasets.** We evaluate our biomedical RAG system on the BIOASQ [11] QA benchmark, which builds on the PubMed database [5]. Specifically, we use a 10% randomly sampled subset of 2.4M biomedical papers for the component analysis, while for the final system, we use the entire dataset of **24M**. Each entry includes a PubMed ID (PMID), title, and abstract, with an average abstract length of 296 tokens. We evaluate our system using the Task-B dataset, which contains expert-annotated questions paired with supporting PubMed IDs (PMIDs). To ensure answerability within our PubMed subset, we first exclude factoid and list questions, which often require full-text access, making evaluation less precise. Second, we retain only questions with at least one PMID in our dataset to ensure they can be answered using our subset.

**Indexing and Retrieval Systems.** We compare three storage and query systems: Elasticsearch, FAISS, and MongoDB.

MongoDB[2], a NoSQL document database, supports full-text search with TF-IDF-like scoring in its self-hosted version. While BM25 ranking is available in MongoDB Atlas Search, it is a cloud-only service and was not used.

Elasticsearch[3], built on Apache Lucene, uses BM25 ranking and inverted indexing for efficient text-based retrieval.

FAISS (Facebook AI Similarity Search) optimizes dense vector similarity search, commonly used in NLP and recommendation systems [12]. We deployed FAISS using a Flask-based server with a FlatL2 index for exhaustive search.

Metrics: We evaluate *indexing speed* and *response time* on 2.4M PubMed papers to determine the best trade-off between efficiency and retrieval performance.

**Retrieval Methods.** Based on recall and precision, we evaluate four retrieval methods—BM25, BioBERT, MedCPT, and a hybrid approach (BM25 + MedCPT).

BM25 [13] is a ranking algorithm that improves upon TF-IDF [14] by incorporating term frequency, document length normalization, and inverse document frequency. It ranks documents based on query relevance using a probabilistic scoring function. We implemented BM25 in Elasticsearch, with stopword removal for improved efficiency. BioBERT [7] is a domain-specific adaptation of BERT [15], pre-trained on PubMed abstracts and PMC articles to enhance biomedical text understanding. We use BioBERT to encode PubMed abstracts into semantic vectors via FAISS, computing document-query similarity with squared Euclidean distance. MedCPT [8] is a contrastive learning-based retriever trained on 255M PubMed query-article interactions. It consists of a query encoder, document encoder, and a cross-encoder reranker. The cross-encoder refines retrieval results by reranking top candidates based on query-document contextual interactions. We use

---

[1]https://github.com/slinusc/medical_RAG_system

System Prompt: You are a scientific medical assistant designed to synthesize responses from specific medical documents. Only use the information provided in the documents to answer questions. The first documents should be the most relevant. Do not use any other information except for the documents provided. When answering questions, always format your response as a JSON object with fields for 'response', 'used_PMIDs'. Cite all PMIDs your response is based on in the 'used_PMIDs' field. Please think step-by-step before answering questions and provide the most accurate response possible. Provide your answer to the question in the 'response' field.

User Prompt: Answer the following question: ...

Context Prompt: Here are the documents:

```
"doc1": {
    "PMID": {...},
    "title": {...},
    "content": {...}
    "relevance_score": {...}
}, ...
```

Fig. 2: Prompting approach for biomedical QA.

MedCPT to encode 2.4M abstracts. We filter results based on positive relevance scores. The Hybrid Retriever integrates BM25 and MedCPT for enhanced retrieval performance. BM25 first ranks a broad set of documents in Elasticsearch, after which MedCPT's cross-encoder reranks the top-$k$ results. This combination leverages BM25's efficiency and MedCPT's semantic understanding to improve recall and precision.

Metrics: *We assess how well each method retrieves relevant documents*. Since *recall* is critical for ensuring comprehensive retrieval, we prioritize methods that maximize relevant document retrieval while maintaining high *precision*.

**Text Generation.** For text generation, we experiment with different prompting strategies for OpenAI's GPT-3.5-turbo (API version May 2024, temperature=0), ensuring that generated responses are accurate and contextually relevant. Given the biomedical domain's strict accuracy requirements, we focus on structured prompts that enhance factual consistency. We experimented with multiple prompting approaches, following best practices in medical NLP [5], [16]. Due to resource constraints, we evaluated GPT-3.5, with limited testing of GPT-4. Observations showed no significant differences in output quality. As illustrated in Figure 2, our final prompt consists of three components: (1) a system prompt with task-specific instructions, (2) a user query, and (3) retrieved documents with PubMed IDs (PMIDs), titles and content.

Metrics: For text generation, we evaluate *answer correctness* in terms of *accuracy*, *recall*, *precision*, and *F1 score*.

### B. Indexing and Query Time

Table I summarizes the performance of Elasticsearch, FAISS, and MongoDB. Elasticsearch excels in full-text retrieval but is less efficient for semantic vector search, which FAISS optimizes for. However, due to their complex data management and indexing mechanisms, MongoDB and Elasticsearch exhibit the slowest indexing speeds.

MongoDB, while providing a flexible NoSQL document storage solution, uses TF-IDF-based text ranking in its self-

hosted version, which leads to significantly slower query response times compared to Elasticsearch and FAISS. The self-hosted MongoDB lacks efficient semantic retrieval, limiting its effectiveness in large-scale biomedical QA.

Based on these results, we selected Elasticsearch for full-text retrieval and FAISS for semantic vector search. Despite its slower indexing speed, Elasticsearch provides a robust text-based search framework, while FAISS offers superior response times for vector-based queries.

TABLE I: Performance metrics for different search methods.

| Method | Type | Index | Response Time | Indexing Speed |
|---|---|---|---|---|
| MongoDB | Sparse | TF-IDF | 26.4 s ± 1.72 s | 10.41 min |
| Elasticsearch | Sparse | BM25 | 82 ms ± 37 ms | 156 min |
| Elasticsearch | Dense | KNN | 24.6 s ± 1.23 s | 171 min |
| FAISS | Dense | L2 Distance | 657 ms ± 127 ms | 41 min |

### C. Document Relevancy

Table II summarizes the retrievers' performance. The Hybrid Retriever achieved the highest recall (0.567), balancing efficiency and accuracy. BM25 exhibited strong precision but lower recall. MedCPT improved semantic retrieval but underperformed in recall, while BioBERT had the weakest results due to a lack of fine-tuning for question-answering tasks. Note that a low recall score does not necessarily indicate incorrect retrieval; rather, it means that the retrieved documents may not be included in the BioASQ-curated set.

TABLE II: Performance comparison of different retrievers.

| Retriever | Vector Type | Recall | Precision |
|---|---|---|---|
| Hybrid Retriever | Hybrid | 0.567 | 0.319 |
| BM25 | Sparse | 0.537 | 0.322 |
| MedCPT | Dense | 0.273 | 0.205 |
| BioBERT | Dense | 0.07 | 0.07 |

### D. Answer Correctness of the RAG System (End-to-End)

For **BM25**, the query is processed using term-based retrieval, ranking documents based on query term occurrence. The top $k$ ranked documents are embedded into the LLM context for response generation. For **MedCPT**, the query is encoded into a vector and compared against document embeddings for similarity search. The retrieved documents are reranked by a cross-encoder, and only those with positive relevance scores are used for response generation. For **Hybrid Retrieval**, BM25 first retrieves $k$ candidate documents, which are then reranked by MedCPT's cross-encoder. Only relevant documents are passed to the LLM. Our results show that the *hybrid retriever* achieves the *best answer correctness* on all metrics (Table III).

TABLE III: Performance metrics of the end-to-end RAG system using different retrievers.

| RAG with Retriever | Accuracy | Recall | Precision | F1 Score |
|---|---|---|---|---|
| GPT-3.5 / Hybrid Retriever | 0.86 | 0.86 | 0.89 | 0.86 |
| GPT-3.5 / MedCPT | 0.83 | 0.83 | 0.86 | 0.84 |
| GPT-3.5 / BM25 | 0.72 | 0.72 | 0.83 | 0.74 |
| GPT-3.5 / BioBERT | 0.63 | 0.63 | 0.85 | 0.67 |

## IV. EVALUATION OF THE FINAL SYSTEM

After selecting the most efficient and effective components, we evaluate the final hybrid RAG system on the full 24M-document PubMed corpus for *retrieval effectiveness*, *response time*, and *answer correctness*.

## A. Effect of Retrieval Depth on Performance

To evaluate the impact of retrieval depth on performance, we experimented with different configurations of BM25 retrieval, varying the number of initially retrieved documents while keeping the reranking step fixed at the top 10 (Table IV).

TABLE IV: Comparison for different retrieval depths (BM25), with reranking applied to the top 10 documents.

| Docs | Accuracy | Recall | Precision | F1 Score | Retrieval Time (s) | Total Time (s) |
|------|----------|--------|-----------|----------|--------------------|----------------|
| 20 | 0.89 | 0.88 | 0.89 | 0.88 | 0.39 ± 0.07 | 1.52 ± 0.42 |
| 50 | 0.90 | 0.90 | 0.89 | 0.90 | 0.82 ± 0.13 | 1.91 ± 0.36 |
| 100 | 0.87 | 0.87 | 0.88 | 0.87 | 1.54 ± 0.16 | 2.62 ± 0.44 |

## B. Analysis of Retrieval Depth Trade-offs

Elasticsearch BM25 retrieval has an average response time of $82 \pm 37$ms, which remains constant across all retrieval depths since it ranks all documents regardless of how many are later passed to reranking. The primary factor affecting response time is the cross-encoder reranking step using Med-CPT, which processes a subset of the retrieved documents and incurs additional computational overhead.

Increasing the number of retrieved documents leads to marginal accuracy improvements but significantly increases the rerank time. Retrieving 50 documents before reranking yields the best accuracy (0.90) and F1 score (0.90) while keeping response time manageable at 1.91 seconds. However, retrieving 100 documents leads to a drop in accuracy (0.87) and an increase in total response time to 2.62 seconds, suggesting diminishing returns beyond 50 documents.

The text generation phase relies on the OpenAI API, which introduces additional latency. The mean response time for generation is 1.07 seconds, with a standard deviation of 0.41 seconds. Since the generation time remains stable across configurations, the overall system latency is primarily determined by the retrieval depth and reranking time.

These results demonstrate that increasing the number of retrieved documents beyond a certain threshold does not necessarily improve system performance. Instead, balancing retrieval depth with reranking efficiency is critical for real-world biomedical question-answering applications.

## V. CONCLUSION AND FUTURE DIRECTIONS

Biomedical question-answering (QA) systems require both efficient retrieval and generation components for accuracy and scalability. This study examines a Retrieval-Augmented Generation (RAG) system for biomedical QA, evaluating retrieval strategies and response time trade-offs.

We assess retrieval methods, including BM25, BioBERT, MedCPT, and a hybrid approach, alongside data stores such as Elasticsearch, MongoDB, and FAISS. Despite strong performance, some limitations remain. The reliance on OpenAI's GPT-3.5 for text generation poses reproducibility challenges due to model updates and API latency. Additionally, retriever and database system evaluations remain limited, requiring broader comparisons.

Future work should explore additional retrievers and evaluate alternative databases for indexing efficiency. Efforts should also focus on retrieval optimization, integrating open-source LLMs, and enabling real-time biomedical applications. Our work highlights trade-offs in retrieval depth, efficiency, and scalability. The system is fully reproducible and extensible, supporting future advancements in retrieval and model integration for research and clinical applications.

## REFERENCES

[1] K. Singhal, T. Tu, J. Gottweis, R. Sayres, E. Wulczyn, L. Hou, K. Clark, S. Pfohl, H. Cole-Lewis, D. Neal, M. Schaekermann, A. Wang, M. Amin, S. Lachgar, P. Mansfield, S. Prakash, B. Green, E. Dominowska, B. A. y Arcas, N. Tomasev, Y. Liu, R. Wong, C. Semturs, S. S. Mahdavi, J. Barral, D. Webster, G. S. Corrado, Y. Matias, S. Azizi, A. Karthikesalingam, and V. Natarajan, "Towards expert-level medical question answering with large language models," 2023. [Online]. Available: https://arxiv.org/abs/2305.09617

[2] A. Asai, Z. Zhong, D. Chen, P. W. Koh, L. Zettlemoyer, H. Hajishirzi, and W.-t. Yih, "Reliable, adaptable, and attributable language models with retrieval," *arXiv preprint arXiv:2403.03187*, 2024.

[3] S. M. T. I. Tonmoy, S. M. M. Zaman, V. Jain, A. Rani, V. Rawte, A. Chadha, and A. Das, "A comprehensive survey of hallucination mitigation techniques in large language models," 2024.

[4] I. Alonso, M. Oronoz, and R. Agerri, "Medexpqa: Multilingual benchmarking of large language models for medical question answering," 2024.

[5] G. Xiong, Q. Jin, Z. Lu, and A. Zhang, "Benchmarking retrieval-augmented generation for medicine," *arXiv preprint arXiv:2402.13178*, 2024.

[6] S. Robertson, H. Zaragoza *et al.*, "The probabilistic relevance framework: Bm25 and beyond," *Foundations and Trends® in Information Retrieval*, vol. 3, no. 4, pp. 333–389, 2009.

[7] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, "Biobert: a pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, no. 4, p. 1234–1240, Sep. 2019. [Online]. Available: http://dx.doi.org/10.1093/bioinformatics/btz682

[8] Q. Jin, W. Kim, Q. Chen, D. C. Comeau, L. Yeganova, W. J. Wilbur, and Z. Lu, "Medcpt: Contrastive pre-trained transformers with large-scale pubmed search logs for zero-shot biomedical information retrieval," *Bioinformatics*, vol. 39, no. 11, p. btad651, Nov 2023.

[9] X. Ma, Z. Yang, and H. Zhang, "A hybrid first-stage retrieval model for biomedical literature," in *Proceedings of the CEUR Workshop on Biomedical Information Retrieval*, 2020. [Online]. Available: https://ceur-ws.org/Vol-2696/paper_92.pdf

[10] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel *et al.*, "Retrieval-augmented generation for knowledge-intensive nlp tasks," *Advances in Neural Information Processing Systems*, vol. 33, pp. 9459–9474, 2020.

[11] A. Krithara, A. Nentidis, K. Bougiatiotis, and G. Paliouras, "Bioasq-qa: A manually curated corpus for biomedical question answering," *Scientific Data*, vol. 10, no. 1, p. 170, 2023.

[12] M. Douze, A. Guzhva, C. Deng, J. Johnson, G. Szilvasy, P.-E. Mazaré, M. Lomeli, L. Hosseini, and H. Jégou, "The faiss library," 2024.

[13] S. E. Robertson, S. Walker, S. Jones, M. Hancock-Beaulieu, and M. Gatford, "Okapi at trec-3," *NIST Special Publication*, vol. 500-225, pp. 109–126, 1994.

[14] K. Spärck Jones, "A statistical interpretation of term specificity and its application in retrieval," *Journal of Documentation*, vol. 28, no. 1, pp. 11–21, 1972.

[15] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," 2019.

[16] B. Meskó, "Prompt engineering as an important emerging skill for medical professionals: Tutorial," *J Med Internet Res*, vol. 25, p. e50638, 2023.