# All Research Paper Summaries

## Summary 1

| Field | Details |
|---|---|
| title | Enhancing Biomedical Question Answering with Parameter-Efficient Fine-Tuning and Hierarchical Retrieval Augmented Generation |
| authors | Yichen Gao, Licheng Zong, Yu Li |
| keywords | Biomedical Question-Answering, Retrieval-Augmented Generation, Large Language Model, Parameter-Efficient Fine-Tuning, BioASQ |
| method_model | Corpus PEFT Searching (CPS): A system integrating Parameter-Efficient Fine-Tuning (PEFT) on a corpus with a hierarchical retrieval-based searching method for biomedical QA. |
| goal_problem | The goal is to enhance biomedical question answering accuracy and comprehensiveness by integrating intricate medical knowledge into LLMs. The problem is to effectively manage and retrieve information from a massive corpus of medical literature to provide accurate answers to complex medical queries. |
| components | BM25 index, Llama2-chat-7B model, Parameter-efficient Fine-Tuning (PEFT) using Low-Rank Adaptation (LoRA), Hierarchical Retrieval-Augmented Generation (RAG) pipeline, Ensemble retriever (sparse and dense retrievers), Text Splitter |
| process | <table><tr><td>step</td><td>mechanism</td><td>advantages</td><td>disadvantages</td></tr><tr><td>Phase A: Document Retrieval</td><td>Build BM25 indexes for PubMed Central (PMC) abstracts. Use the question as a keyword to retrieve relevant documents.</td><td>Efficient initial document retrieval.</td><td>May retrieve irrelevant documents if keywords are ambiguous.</td></tr><tr><td>Phase B: Answer Generation</td><td>Fine-tune Llama2-chat-7B</td><td>Improves answer accuracy</td><td>Requires significant computational</td></tr></table> |

| step | mechanism | advantages | disadvantages |
|---|---|---|---|
| Phase A: Document Retrieval | Build BM25 indexes for PubMed Central (PMC) abstracts. Use the question as a keyword to retrieve relevant documents. | Efficient initial document retrieval. | May retrieve irrelevant documents if keywords are ambiguous. |
| Phase B: Answer Generation | Fine-tune Llama2-chat-7B | Improves answer accuracy | Requires significant computational |

| | | | |
|---|---|---|---|
| with PEFT | model using PEFT (LoRA) on BioASQ training set. Use an ensemble retriever to find relevant snippets, then feed question and snippets to the fine-tuned model for answer generation. | by fine-tuning on domain-specific data. | resources for fine-tuning. |
| Phase A+: Hierarchical RAG for Answer Generation | Implement a two-stage hierarchical RAG pipeline. First-level BM25 retrieves relevant documents, which are split into chunks. Second-level ensemble retriever finds the most relevant chunks. These | Combines broad and fine-grained retrieval for better context. | Increased complexity due to multi-stage retrieval. |

| | |
|---|---|
| | chunks and the question are fed to the fine-tuned model for answer generation. |
| variables | |

| Key | Value |
|---|---|
| dependent | ['Answer Accuracy', 'Recall', 'F-Measure', 'MAP', 'GMAP', 'R-2(Rec)', 'R-2(F1)', 'R-SU4(Rec)', 'R-SU4(F1)'] |
| independent | ['Parameter-Efficient Fine-Tuning (PEFT)', 'Retrieval-Augmented Generation (RAG)', 'Retrieval Unit (Chunk vs. Snippet)', 'Retrieval Source (Test Set vs. Training Set)', 'Ensemble Retriever', 'BM25 Retriever', 'Dense Retriever'] |
| mediating | [] |
| moderating | [] |

| | |
|---|---|
| inputs | Biomedical question, PubMed Central (PMC) documents (Phase A), Golden enriched snippets (Phase B), Relevant document chunks (Phase A+) |
| outputs | List of relevant documents (Phase A), Ideal answers to biomedical questions (Phase B, A+) |
| features | Parameter-Efficient Fine-Tuning (PEFT) with LoRA, |

| | Hierarchical Retrieval-Augmented Generation (RAG), Ensemble retriever combining sparse and dense methods, BM25 indexing and retrieval, Llama2-chat-7B model |
|---|---|
| contribution_value | This work demonstrates that PEFT and hierarchical RAG can significantly improve performance in biomedical QA tasks, achieving competitive results on the BioASQ challenge by efficiently fine-tuning a smaller LLM and leveraging a sophisticated retrieval strategy. |
| positive_impacts | Improved accuracy and relevance of answers in biomedical question answering., Efficient utilization of computational resources through PEFT., Enhanced information retrieval by combining sparse and dense methods., Provides a robust framework for future biomedical QA research. |
| negative_impacts | Performance might be limited by the scale and context window of the base LLM., Potential for noise in retrieved information if not carefully managed., Reliance on specific datasets (BioASQ) may limit generalizability without further adaptation. |
| critical_analysis | The study effectively combines PEFT and RAG to enhance biomedical QA, showing strong performance on the BioASQ challenge. The hierarchical retrieval in Phase A+ is a key innovation. However, the paper could benefit from a more in-depth comparison with state-of-the-art dense retrieval methods and a broader discussion on the ethical implications of using LLMs in healthcare. |
| tools_used | Pyserini, LangChain, Llama2-chat-7B, BM25, LoRA, bge-large-en |
| paper_structure | The paper is structured into Introduction, Related Work, Methodology (detailing Phase A, B, and A+), Results and Analyses (including official evaluations and ablation studies), Discussion, Acknowledgments, and References. |
| diagrams_flowcharts | Yes, the paper includes three figures illustrating the system overview for Phase A, Phase B, and Phase A+ respectively. |
| url | None |
| pdfurl | None |
| year | None |