

---

# MedRAG-2

---

Ahmad Sait  
KAUST

ahmed.sait@kaust.edu.sa

Wenxuan Zhu  
KAUST

wenxuan.zhu@kaust.edu.sa

Aznaur Aliev  
KAUST

aznaur.aliev@kaust.edu.sa

## Abstract

Large Language Models (LLMs) show promise in medical question answering (QA), but remain prone to hallucination, poor citation handling, and unreliable output formatting, especially in high-stakes clinical domains. We analyze two recent retrieval-augmented generation (RAG) systems, MedGraphRAG and MedRAG, and identify core limitations including excessive LLM usage, lack of structured output enforcement, and ungrounded or repetitive reasoning. We then introduce **MedRAG-2**, a system designed to address these gaps through five key enhancements: (1) domain-specific corpus curation; (2) prompt redesign enforcing citation grounding, fallback logic, and anti-hallucination constraints; (3) Pydantic-based structured output enforcement; (4) cross-encoder-based document re-ranking; and (5) retrieval-aware search-augmented prompting. Experimental results on five medical QA benchmarks (MMLU, MedQA, MedMCQA, PubMedQA, BioASQ) show consistent performance gains over MedRAG and LLM-only baselines. Qwen2.5-7B with search augmentation achieves up to +28.4% on PubMedQA, and LLaMA 3-8B with our prompt and output structure gains up to +11% on key datasets. Our findings highlight the importance of grounded reasoning, structured output, and semantic retrieval for reliable medical QA.

## 1 Introduction

Recent studies have shown that Large Language Models (LLMs), when used in medical settings such as diagnosis or clinical decision support, often exhibit a critical limitation: they hallucinate. That is, they generate confident but factually incorrect or misleading outputs, especially when faced with queries outside the scope of their training data. This is particularly dangerous in medicine, where incorrect information can have serious consequences.

The issue is exacerbated in clinical domains where data is sparse (rare diseases) and temporally evolving (new diseases are introduced)—conditions under which LLMs struggle due to the lack of grounding in real-time or structured domain knowledge (e.g., from EHRs or clinical databases). As a result, there is an urgent need for methods that enhance the factual accuracy of LLMs in medical tasks by grounding their outputs in trusted, structured sources such as knowledge graphs derived from datasets like MIMIC-III and MIMIC-IV.

While recent work has attempted to address these limitations through Retrieval-Augmented Generation (RAG) systems that integrate biomedical corpora or knowledge graphs, existing approaches remain inefficient, brittle, and prone to hallucinations due to poor prompt design, redundant retrieval, and unstructured outputs. In this work, we propose **MedRAG-2**, a more robust medical RAG pipeline

that improves grounding and accuracy through prompt enforcement, structured output with Pydantic, cross-encoder re-ranking, and real-time search-augmented QA.

Code for MedRAG-2 are publicly available at: <https://github.com/AhmadSait/MedRAG-2>

## 2 Background

To address hallucinations and lack of factual grounding in LLMs, several recent methods have emerged in the medical Retrieval-Augmented Generation (RAG) domain. These approaches integrate structured biomedical knowledge, either as knowledge graphs or as curated corpora and snippet-level evidence retrieved from biomedical literature, clinical textbooks, and decision-support tools - as seen in systems like MedGraphRAG [1] and MedRAG [2] - into the language model pipeline to improve accuracy and trustworthiness.

### 2.1 MedGraphRAG

MedGraphRAG introduces a triple-linked graph structure that connects user data, trusted clinical knowledge, and controlled vocabularies. This structure supports evidence-based medical reasoning by reducing hallucination risks during response generation.

For MedGraphRAG’s implementation, we have found major inefficiencies while reproducing their code:

#### Excessive LLM API Calls

`agentic_chunker.py` repeatedly calls GPT-4 to summarize each chunk (`_get_new_chunk_summary`), to title each chunk, to update the chunk metadata with each new proposition (e.g., “The patient was administered 500mg of acetaminophen.”), and to decide whether a proposition fits an existing chunk (`_find_relevant_chunk`). This results in many redundant LLM calls per document, leading to high token usage and increased API costs.

#### Double LLM Passes for Proposition Extraction

`data_chunk.py` extracts propositions by first using a LangChain hub chain (`wfh/proposal-indexing`), then piping that output into GPT-4 again for sentence-level extraction, followed by a separate structured parser. This leads to multiple GPT-4 calls per paragraph, compounding token usage and increasing cost.

#### Redundant Graph Construction Calls

`creat_graph.py` builds a Neo4j graph by calling GPT-4 twice through `KnowledgeGraphAgent.run()` (once for freeform output, once for structured), then embedding each node individually via OpenAI and assigning IDs. This creates a new graph for every chunk, with repeated LLM and embedding API calls that significantly increase runtime and cost.

We made several attempts to reproduce MedGraphRAG’s code. Despite numerous bug fixes and lack of open-source corpora, we estimated that it would approximately take 42 days to run their code on a single A100 with 10 workers for API parallelization, processing 48000 MIMIC files. Table 1 shows preliminary results that were obtained from running their system on only 1000 files (less than 5% of the original MIMIC dataset), which took a very time-demanding time of 62 hours. This experiment used a LLaMA-3.3-70B-Instruct-AWQ-INT4 as the knowledge graph agent and LLaMA-3-8B for evaluation.

### 2.2 MedRAG

MedRAG combines multiple domain-specific retrievers (e.g., BM25, MedCPT) and corpora (e.g., PubMed, StatPearls, MedCorp) to enhance medical QA through comprehensive retrieval-augmented generation pipelines.

Despite MedRAG’s robust retrieval setup, our implementation revealed several key limitations that hinder reliability and accuracy.

Table 1: Comparison Between Our MedGraphRAG Attempt and LLM-only Accuracy in %. The bottom row shows the absolute improvement of LLM-only over our approach.

	MMLU	MedQA	MedMCQA	PubMedQA	BioASQ
<b>Our attempt</b>	64.37	56.17	47.62	30.20	51.46
<b>LLM-only Accuracy</b>	<b>71.90</b>	<b>59.94</b>	<b>58.79</b>	<b>58.60</b>	<b>75.89</b>
<b>(Ours - LLM)</b>	-7.53	-3.77	-11.17	-28.40	-24.43

### Ineffective Use of Large General Corpora

Table 7 in MedRAG’s paper shows that Wikipedia, despite its massive size (29.9M snippets), consistently underperforms across nearly all MIRAGE tasks. This confirms that bigger isn’t better if domain relevance is weak. Their choice to include Wikipedia dilutes the relevance of retrieved contexts, adding noise and harming precision.

### Hallucinations

A flaw in MedRAG’s prompt is the absence of grounding constraints. It asks the model to answer using the documents but does not specify how to use or cite them. This leaves room for unsupported claims or blending of sources, which undermines traceability.

Another flaw that was critical when running MedRAG’s code is that the model exhibits a verbose hallucination in which it loops through the same set of document citations multiple times, leading it to infinitely hang without converging to an answer.

### Ambiguous Output

Despite operating in a high-stakes medical domain, MedRAG does not enforce structured output. The prompt merely instructs: “Think step-by-step and generate the output in JSON.” However, prior work has shown that prompt-based instructions alone are insufficient to ensure structured compliance. The StructuredRAG benchmark (Shorten et al., 2024) found that even with explicit prompts, LLMs frequently produce malformed or ambiguous outputs, leading to parse failures and unreliable downstream evaluation [3].

We observed similar behavior. The model often repeated a single document citation or rationale indefinitely, failing to converge or advance through reasoning steps.

#### Example 1: Infinite citation loop

"Document [1] lists the clinical consequences of lesions in the cavernous sinus, including facial paralysis, which may be caused by compression of the facial nerve at the stylomastoid foramen."

(This line was repeated indefinitely; the model failed to proceed to the next question.)

#### Example 2: Redundant reasoning steps

"The patient’s T-score is below -2.5, which indicates osteoporosis. Therefore, pharmacotherapy is indicated."

(This rationale was repeated with minor variations as the number of reasoning steps increased.)

"According to the guidelines... patients be considered for treatment when BMD is >2.5 SD below the mean..."

(The same guideline was reused repeatedly with rewording.)

These behaviors indicate a lack of structured control and logical progression. Furthermore, the prompt fails to specify what the model should do when no retrieved document contains sufficient information. Without explicit fallback logic, the model resorts to guessing or hallucinating unsupported facts.

Together, these flaws—repetition, hallucination, and absence of structured enforcement—undermine the reliability of MedRAG and pose significant risks in clinical applications.

### 3 Our System

We propose several solutions to solve MedRAG’s limitations:

#### 3.1 Corpus Selection

To address the limitations observed in MedRAG’s use of large, general-purpose corpora like Wikipedia, our system deliberately avoids aggregating broad sources with low domain specificity. Instead, we prioritize targeted medical datasets such as MedCorp, which is a combination of three datasets - PubMed, StatPearls, and medical textbooks - which offer higher semantic relevance and factual density. This improves retrieval precision and reduces the noise introduced by irrelevant or loosely related documents like Wikipedia.

#### 3.2 Prompt Redesign to Prevent Hallucinations and Citation Loops

We observed frequent hallucination patterns in MedRAG, particularly repetitive loops where the model recycled the same set of document citations without reaching a definitive answer. This behavior stems from shortcomings in MedRAG’s prompt, which offers no constraints on evidence citation, factual grounding, or fallback logic.

MedRAG’s prompt suffers from three critical limitations: (1) it provides no rules for how retrieved documents should be cited, leading to repeated or ambiguous references; (2) it lacks a fallback mechanism when no supporting evidence is found, causing the model to guess or stall; and (3) it loosely encourages a “definite answer” without enforcing factual accuracy, allowing for speculation or fabricated content. Our redesigned prompt directly addresses these gaps. It enforces a one-to-one grounding policy, introduces explicit fallback logic, and prohibits hallucinations through the constraint: “Do not invent new facts.” These changes promote factual consistency, prevent citation inflation, and ensure traceable, grounded reasoning.

Figure 1 shows MedRAG’s original prompt, which exhibits these shortcomings in structure, grounding, and citation policy.

To address these issues, we designed a new prompt enforcing three key constraints:

- **One-to-one grounding:** Each reasoning step must cite a distinct document or fallback to prior knowledge.
- **Fallback logic:** If no document supports the answer, the model must explicitly state this and stop citing.
- **Anti-hallucination policy:** The prompt includes a hard constraint: “*Do not invent new facts.*”

These changes ensure factual consistency, prevent citation inflation, and enforce clear attribution throughout the reasoning process. Figure 2 presents our redesigned prompt.

To evaluate its impact, we tested our prompt under both LLaMA 3-8B and Qwen 3-8B configurations. As shown in Table 2, LLaMA gains up to +5% accuracy on MedMCQA and PubMedQA benchmarks. Table 3 shows consistent improvements on Qwen 3-8B, particularly in datasets prone to citation ambiguity.

Table 2: Effect of Our Prompt on LLaMA 3-8B (Accuracy in %). Structured prompt improves performance, particularly in MedMCQA and PubMed QA.

k	Dataset	Enhanced	CoT/MedRAG	MMLU	MedQA	MedMCQA	PubMed QA	BioASQ	Avg
N/A	N/A	No	CoT	<b>66.20</b>	55.02	47.90	42.07	61.50	54.52
16	MedCorp	No	MedRAG	65.01	54.05	47.26	40.20	60.84	53.47
16	MedCorp	Yes	MedRAG	64.74	<b>55.38</b>	<b>51.33</b>	<b>49.60</b>	<b>67.96</b>	<b>57.80</b>
(Enhanced - Baseline)				-0.27	+1.33	+4.07	+9.40	+7.12	+4.33

```

general_medrag_system = '''
You are a helpful medical expert,
and your task is to answer a multi-choice medical question using the relevant documents.
Please first think step-by-step and then choose the answer from the provided options.
Organize your output in a json formatted as Dict{"step_by_step_thinking": Str(explanation),
"answer_choice": Str{A/B/C/...}}.
Your responses will be used for research purposes only, so please have a definite answer.
'''

general_medrag = Template('''
Here are the relevant documents:
{{context}}

Here is the question:
{{question}}

Here are the potential choices:
{{options}}

Please think step-by-step and generate your output in json:
''')

```

Figure 1: MedRAG’s original prompt

Table 3: Effect of Our Prompt on Qwen3-8B (Accuracy in %). Structured prompt improves performance, particularly in MedMCQA and PubMed QA.

k	Dataset	Enhanced	CoT or MedRAG	MMLU	MedQA	MedMCQA	PubMed QA	BioASQ	Avg
N/A	N/A	No	CoT	<b>82.80</b>	<b>67.46</b>	55.67	36.60	70.44	62.59
16	MedCorp	No	MedRAG	82.34	67.10	55.61	35.40	70.56	62.20
16	MedCorp	Yes	MedRAG	77.09	64.63	<b>56.79</b>	<b>41.08</b>	<b>74.12</b>	<b>62.74</b>
(Enhanced - Baseline)				-5.25	-2.47	+1.18	+5.68	+3.56	+0.54

These results confirm that enforcing grounding, fallback, and anti-hallucination constraints at the prompt level leads to more reliable and accurate responses across multiple medical benchmarks.

### 3.3 Structured Output with Pydantic

Prior to enforcing structured output with Pydantic, the model frequently produced responses that were either syntactically invalid or semantically ambiguous. Common issues included free-form answers lacking a clearly stated final choice, isolated letters without accompanying reasoning, or malformed JSON—such as missing brackets, inconsistent key names, or multi-answer strings like “The correct answer is C” that violated expected enumeration formats. These inconsistencies made downstream parsing unreliable and evaluation scripts brittle. Below is an example of such an output:

**Observed model output:**

```

general_medrag_system = '''
You are a meticulous, evidence-based medical question-answering assistant.
Answer using only the provided documents—or exactly one standard fact if none suffice.

Follow these rules without exception:

1. Grounding
   - Cite each document once and only once.
   - As soon as you've used "Document X," you may never cite it again.

2. Fallback
   - If no document gives the answer, preface with
     "Applying standard knowledge: [fact]"
   - Then stop using documents entirely.

3. Step-By-Step (Max 3 Steps)
   - Use at most three numbered steps.
   - Each step must reference a different document or the one standard fact.
   - Immediately after your final step, stop reasoning.

4. Immediate JSON Output
   - Without any extra text, output only this JSON:
     ```json
     {
       "step_by_step_thinking": "up to 3 numbered steps...",
       "answer_choice": "<one of 'A','B','C','D'>"
     }
     ```
   - Any deviation (extra keys, additional sentences, repeats) → your answer is dropped.

5. No Hallucinations
   - Do not invent new facts.
   - Do not exceed the 3-step limit.
'''

general_medrag = Template('''
Here are the relevant documents:
{{context}}

Question:
{{question}}

Choices:
{{options}}

Please follow the system instructions exactly:
1. Ground each step in "Document X" or in one standard fact.
2. Number your reasoning steps clearly; avoid repetition.
3. Do not speculate or invent new details.
4. Output only this JSON object:

{
  "step_by_step_thinking": "...",
  "answer_choice": "<one of 'A','B','C','D'>"
}
''')

```

Figure 2: Our redesigned prompt with grounding constraints and fallback logic

"The salivary glands are derived from ectomesenchyme. But the sweat glands are not. So if the options include sweat glands, then D is not correct. So the answer would be C. Melanocytes are derived from the neural crest, which is ectomesenchyme. Therefore, the correct answer is C."

Although the answer is present, it is embedded in an unstructured narrative, making automated extraction non-trivial. By integrating Pydantic, we enforced strict adherence to a predefined schema, ensuring that each model output was both syntactically valid and semantically unambiguous. This substantially improved evaluation reliability and simplified result aggregation. Figure 3 illustrates our Pydantic integration for structured response enforcement.

As shown in Table 4, incorporating structured output led to consistent performance improvements across all benchmarks for the LLaMA 3-8B model, with gains of up to 11 percentage points on PubMedQA and 6 points on BioASQ. On average, structured output improved accuracy by 3.8–4.7 points depending on the prompting strategy.

We replicated this experiment using Qwen 3-8B to test generalizability. As shown in Table 5, structured output again resulted in higher accuracy, particularly in datasets with longer, more complex answers like MedMCQA and BioASQ.

```

from pydantic import BaseModel
from enum import Enum

class AnswerChoice(str, Enum):
    A = "A"
    B = "B"
    C = "C"
    D = "D"

class MedRAGOutput(BaseModel):
    step_by_step_thinking: str
    answer_choice: AnswerChoice

json_schema = MedRAGOutput.model_json_schema()

def openai_client(messages, **kwargs):
    completion = oai_client.chat.completions.create(
        messages=messages,
        seed=42,
        timeout=150,
        max_tokens=4096,
        extra_body={"guided_json": json_schema},
        **kwargs
    )
    choice = completion.choices[0].message
    return choice.reasoning_content

```

Figure 3: Integration of Pydantic for Reliable Structured Output

Table 4: Effect of Structured Output on LLaMA 3-8B Performance (Accuracy in %)

k	Dataset	Structured	CoT/MedRAG	MMLU	MedQA	MedMCQA	PubMed QA	BioASQ	Avg
N/A	N/A	No	CoT	66.20	<b>55.02</b>	47.90	42.07	61.50	54.52
N/A	N/A	Yes	CoT	66.10	54.91	51.87	51.03	<b>67.50</b>	<b>58.28</b>
16	MedCorp	No	MedRAG	65.01	54.05	47.26	40.20	60.84	53.47
16	MedCorp	Yes	MedRAG	<b>66.82</b>	53.71	<b>52.34</b>	<b>51.13</b>	66.75	58.15
(Structured - Baseline)				+1.81	-0.34	+5.08	+10.93	+5.91	+4.68

Together, these results highlight the robustness and portability of schema-constrained output generation, improving both evaluation stability and overall accuracy in retrieval-augmented medical QA.

### 3.4 Re-ranking

To further improve the precision of document retrieval in our medical RAG pipeline, we incorporate a cross-encoder-based re-ranking stage, following the retrieval pipeline design explored in recent dense and hybrid retrieval systems [5, 6]. While bi-encoder retrievers such as MedCPT and SPECTER provide efficient initial filtering, they rank documents based on independent query and document embeddings. This approach struggles with nuanced semantic dependencies—such as negation, causality, or contradiction—that are prevalent in medical literature. Cross-encoders address this by

Table 5: Effect of Structured Output on Qwen 3-8B Performance (Accuracy in %)

k	Dataset	Structured	CoT/MedRAG	MMLU	MedQA	MedMCQA	PubMed QA	BioASQ	Avg
N/A	N/A	No	CoT	<b>82.80</b>	<b>67.46</b>	55.67	36.60	70.44	62.59
N/A	N/A	Yes	CoT	79.69	64.88	58.10	<b>42.78</b>	<b>74.52</b>	<b>63.99</b>
16	MedCorp	No	MedRAG	82.34	67.10	55.61	35.40	70.56	62.20
16	MedCorp	Yes	MedRAG	80.32	63.27	<b>58.45</b>	42.73	73.47	63.65
(Structured - Baseline)				-2.02	-3.83	+2.84	+7.33	+2.91	+1.45

jointly encoding the (question, document) pair, enabling fine-grained, context-aware relevance scoring.

The pipeline proceeds in three stages:

- The query is encoded using a bi-encoder and compared against a FAISS index of document embeddings. Multiple retrievers (e.g., BM25 and MedCPT) can be fused using Reciprocal Rank Fusion (RRF) to return top-k candidate snippets.
- Each retrieved document is paired with the original query and passed through a cross-encoder model, which outputs a relevance score. The documents are then sorted by this score, yielding a more semantically aligned top-k set.
- The reranked snippets are used for downstream answer generation. Since these documents are better aligned with the intent of the question, they improve the grounding and accuracy of the language model’s responses.

This method allows the model to disambiguate subtle medical nuances. For instance, given the question "What hormone lowers blood glucose?", bi-encoders might surface documents about both insulin and glucagon due to lexical overlap. However, the cross-encoder correctly elevates insulin, the only truly relevant answer, by jointly reasoning over the full statement.

As shown in Table 6, integrating re-ranking into our Qwen 3-8B pipeline led to consistent improvements in key QA benchmarks, demonstrating the benefit of context-sensitive scoring.

Table 6: Effect of Re-ranking on Qwen 3-8B with Our Prompt and Structured Output (Accuracy in %). Experiments use MedCorp (k=16) under the MedRAG setting.

Enhanced	Structured	Re-ranking	MMLU	MedQA	MedMCQA	PubMedQA	BioASQ	Avg
Yes	Yes	No	82.25	<b>64.60</b>	<b>59.88</b>	44.05	74.17	64.99
Yes	Yes	Yes	<b>82.81</b>	64.41	59.59	<b>44.43</b>	<b>74.69</b>	<b>65.19</b>
(Re-ranked - No Re-ranked)			+0.56	-0.19	-0.29	+0.38	+0.52	+0.20

These results confirm that even modest reordering of retrieved documents using semantic relevance contributes to improved grounding and accuracy in complex medical QA tasks.

### 3.5 Search-Augmented Medical Question Answering

We propose a system that enhances Large Language Models (LLMs) for medical question answering by combining strategic tool use with search-augmented generation. Rather than relying solely on the model’s internal knowledge, our approach leverages web search to inject up-to-date and contextually relevant medical information into the answering pipeline. Central to our method is the generation of optimized search queries that are more effective than directly using the original question.

In traditional retrieval-augmented medical QA, using the raw question as a search query often leads to suboptimal retrieval. It may return copies of the question itself, overlook critical medical terminology, or miss relevant answer choice concepts. To address this, we prompt the Qwen2.5-7B-Instruct model to reformulate the input question into a search-optimized query. The prompt guides the model to:

- Extract core medical entities and concepts
- Identify clinically relevant relationships
- Eliminate question-specific or redundant phrasing



- Preserve domain-specific terminology
- Generate concise and targeted queries suitable for retrieval

This reformulated query is then used to retrieve external evidence, which is fed back into the LLM for answer generation. Our experimental results show that this query optimization and retrieval-aware prompting pipeline leads to measurable gains in accuracy and reliability on medical QA benchmarks.

We assess the effectiveness of this query optimization pipeline in Section ??, where we show that search-augmented prompting improves answer accuracy across multiple datasets.

## 4 Evaluation

### 4.1 MedRAG vs MedRAG-2

We first evaluate the combined effect of our enhanced prompt and structured output enforcement. Tables ?? and 8 report the accuracy of LLaMA 3-8B and Qwen 3-8B, respectively, across five medical QA datasets. We compare CoT prompting, MedRAG baseline, and our proposed techniques in isolation and in combination.

Table 7: Ablation Study: Accuracy (%) of LLaMA 3-8B and Qwen 3-8B on MedCorp (k=16) using the MedRAG setup. We evaluate the effect of enhanced prompting, structured output, and cross-encoder re-ranking across five medical QA datasets.

LLM	Enhanced	Structured	Re-ranked	MMLU	MedQA	MedMCQA	PubMed QA	BioASQ	Avg
LLaMA 3-8B	No	No	No	65.01	54.05	47.26	40.20	60.84	53.47
LLaMA 3-8B	Yes	Yes	No	68.04	54.52	53.31	<b>52.40</b>	66.99	59.05
Qwen 3-8B	Yes	Yes	No	82.34	<b>67.10</b>	55.61	35.40	70.56	62.20
Qwen 3-8B	Yes	Yes	Yes	82.25	64.60	<b>59.88</b>	44.05	74.17	64.99
Qwen 3-8B	Yes	Yes	Yes	<b>82.81</b>	64.41	59.59	44.43	<b>74.69</b>	<b>65.19</b>

The best performance is achieved when both enhancements are applied, with LLaMA 3-8B reaching 59.05% average accuracy—an improvement of 5.6 points over its MedRAG baseline and 4.5 points over CoT.

Table 8: Effect of Enhanced Prompt and Structured Output on Qwen 3-8B (Accuracy in %). Experiments use MedCorp (k=16) under both CoT and MedRAG settings.

k	Dataset	Enhanced	Structured	CoT/MedRAG	MMLU	MedQA	MedMCQA	PubMed QA	BioASQ	Avg
N/A	N/A	No	No	CoT	<b>82.80</b>	<b>67.46</b>	55.67	36.60	70.44	62.59
N/A	N/A	No	Yes	CoT	79.69	64.88	58.10	42.78	<b>74.52</b>	63.99
16	MedCorp	No	No	MedRAG	82.34	67.10	55.61	35.40	70.56	62.20
16	MedCorp	Yes	No	MedRAG	77.09	64.63	56.79	41.08	74.12	62.74
16	MedCorp	No	Yes	MedRAG	80.32	63.27	58.45	42.73	73.47	63.65
16	MedCorp	Yes	Yes	MedRAG	82.25	64.60	<b>59.88</b>	<b>44.05</b>	74.17	<b>64.99</b>

Qwen 3-8B also benefits from prompt and output structure alignment, gaining 2.8 points over its MedRAG baseline and outperforming GPT-4 in average score.

### 4.2 MedRAG vs MedRAG-2 vs Search-Augmented QA

We also compare the use of search-augmented QA via Qwen2.5-7B with our pipeline against the MedCorp-based MedRAG baseline and other open and closed LLMs. Table ?? summarizes accuracy across five medical benchmarks.

The search-augmented Qwen2.5-7B model outperforms all other baselines by a wide margin, especially on PubMedQA and BioASQ, where real-time retrieval fills knowledge gaps that static corpora cannot address.

Table 9: Ablation Study: Accuracy (%) of different LLM configurations on MedCorp (k=16) using the MedRAG setup. We evaluate the effect of enhanced prompting, structured output, cross-encoder re-ranking, and Search-Engine-Augmented QA across five medical QA datasets.

LLM	Enhanced	Structured	Re-ranked	MMLU	MedQA	MedMCQA	PubMedQA	BioASQ	Avg
LLaMA 3-8B	No	No	No	65.01	54.05	47.26	40.20	60.84	53.47
LLaMA 3-8B	Yes	Yes	No	68.04	54.52	53.31	52.40	66.99	59.05
Qwen 3-8B	No	No	No	82.34	<b>67.10</b>	55.61	35.40	70.56	62.20
Qwen 3-8B	Yes	Yes	No	82.25	64.60	59.88	44.05	74.17	64.99
Qwen 3-8B	Yes	Yes	Yes	<b>82.81</b>	64.41	59.59	44.43	74.69	65.19
Qwen2.5-7B	N/A	N/A	N/A	74.00	66.00	<b>63.00</b>	<b>89.00</b>	<b>90.00</b>	<b>76.40</b>

## 5 Conclusion

## References

- [1] Wu, J., Zhu, J., Qi, Y., Chen, J., Xu, M., Menolascina, F. & Grau, V. (2024) Medical Graph RAG: Towards safe medical large language models via graph retrieval-augmented generation. *arXiv preprint arXiv:2408.04187*. <https://arxiv.org/abs/2408.04187>
- [2] Zhao, X., Liu, S., Yang, S.Y. & Miao, C. (2025) MedRAG: Enhancing retrieval-augmented generation with knowledge graph-elicited reasoning for healthcare copilot. *arXiv preprint arXiv:2502.04413*. <https://arxiv.org/abs/2502.04413>
- [3] Guo, Z., Xia, L., Yu, Y., Ao, T. & Huang, C. (2024) LightRAG: Simple and fast retrieval-augmented generation. *arXiv preprint arXiv:2410.05779*. <https://arxiv.org/abs/2410.05779>
- [4] Shorten, C., Pierse, C., Smith, T. B., Cardenas, E., Sharma, A., Trengrove, J., & van Luijt, B. (2024) StructuredRAG: Measuring structured response adherence in retrieval-augmented generation. *arXiv preprint arXiv:2408.11061*. <https://arxiv.org/abs/2408.11061>
- [5] Khattab, O., & Zaharia, M. (2020). ColBERT: Efficient and effective passage search via contextualized late interaction over BERT. Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval, 39–48. <https://arxiv.org/abs/2004.12832>
- [6] Qu, Y., Liu, Y., Yang, J., Chen, W., Tang, D., Duan, N., & Zhou, M. (2021). RocketQA: An optimized training approach to dense passage retrieval for open-domain question answering. Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP), 5835–5847. <https://arxiv.org/abs/2010.08191>