

Deep Learning and Retrieval-Augmented Generation Methods for Autonomous Clinical Diagnosis and Decision Support

A Progressive Architectural Analysis of Grounded Clinical Question Answering Systems

Alur Taher Basha^{#1}, Mr. P. Bharath Kumar^{#2}, Dr. D. William Albert^{#3}

^{#1}M.Tech Student, Dept. of CSE, Bheema Institute of Technology and Science, Adoni, India

^{#2}Associate Professor, Dept. of CSE, Bheema Institute of Technology and Science, Adoni, India

^{#3}Professor, Head of Dept. CSE, Bheema Institute of Technology and Science, Adoni, India

#1 taherbasha295@gmail.com , #2 bharathkumar1218@gmail.com , #3 dr.albertdwgtl@gmail.com

Abstract:

Large Language Models (LLMs) have demonstrated remarkable potential in healthcare knowledge processing; however, their susceptibility to hallucination and static knowledge boundaries significantly limits safe clinical deployment. This paper introduces MedRAG, a healthcare conversational assistant built on a multi-tier Retrieval-Augmented Generation (RAG) architecture that grounds every generated response in dynamically retrieved, authoritative medical literature. Three progressively sophisticated architectural configurations are designed, implemented, and systematically evaluated: Foundational RAG employing dense MedCPT retrieval, Optimised RAG combining hybrid BM25 and dense search with cross-encoder re-ranking, and Modular RAG integrating knowledge graph traversal with iterative faithfulness verification. Experiments across five standardised benchmarks—MedQA-USMLE, MedMCQA, PubMedQA, MMLU-Medical, and a private chronic-pain diagnostic corpus—demonstrate statistically significant performance improvements over standalone LLM baselines. The Modular RAG configuration paired with GPT-4-Turbo achieves 91.1% accuracy on MedQA-USMLE, equalling the Med-Gemini state-of-the-art, while the fully open-source variant with LLaMA-3-8B attains 79.7%, directly exceeding GPT-4-Turbo without retrieval (73.4%). Hallucination rates are reduced from 39.1% to 9.2% through RAGAS faithfulness verification. Ablation analysis identifies cross-encoder re-ranking (−6.0 pp) and knowledge graph retrieval (−4.9 pp) as the highest-impact architectural components. These findings demonstrate that retrieval quality is a more decisive performance determinant than generator model scale, enabling resource-constrained healthcare systems to achieve frontier-level clinical AI performance.

Keywords— Retrieval-Augmented Generation; Healthcare AI; Clinical Decision Support; Large Language Models; Hallucination Reduction; Medical Question Answering; Knowledge Graphs; RAGAS; Biomedical NLP; Vector Search.

I. INTRODUCTION

The convergence of large-scale natural language processing and clinical informatics has opened compelling possibilities for intelligent healthcare support. Physicians must synthesise evidence from randomised controlled trials, clinical guidelines, drug databases, and patient histories within the narrow time constraints of clinical encounters. Artificial intelligence, and large language models in particular, offers a means to alleviate this informational burden.

Contemporary LLMs such as GPT-4, Google PaLM-2, and Meta LLaMA-3 achieve human-competitive performance across standardised medical licensing examinations. Despite these achievements, deploying standalone LLMs in live clinical environments raises critical safety concerns. Training corpora carry fixed knowledge cutoff dates, and models exhibit a well-documented tendency to generate factually incorrect yet linguistically fluent statements—a behaviour termed hallucination—that can directly endanger patient outcomes when applied to clinical decisions.

Retrieval-Augmented Generation (RAG) provides a principled architectural response. By coupling generative capacity with a retrieval mechanism that fetches relevant passages from a curated, continuously updated external knowledge base, RAG systems ground every response in verifiable, current evidence. RAG preserves the fluency of generative models while introducing the factual grounding and source attributability required for safe clinical deployment.

This paper presents MedRAG, a healthcare conversational assistant that embeds a multi-stage RAG pipeline at its architectural core. Three progressive RAG configurations are designed, implemented, and comparatively evaluated—Foundational, Optimised, and Modular—against standalone LLM baselines across five medical benchmarks. The principal contributions of this work are: (1) a three-tier progressive RAG architecture tailored for clinical question answering, with systematic evaluation across 12 model-architecture combinations; (2) demonstration that Modular RAG with LLaMA-3-8B (79.7% MedQA) surpasses GPT-

4-Turbo without retrieval (73.4%), democratising clinical AI access; (3) reduction of hallucination rates from 39.1% to 9.2% via RAGAS faithfulness verification and iterative refinement; and (4) ablation analysis identifying cross-encoder re-ranking and medical knowledge graph integration as the two highest-impact architectural components.

II. LITERATURE REVIEW

The application of computational text processing to clinical documentation dates to the early 1990s. The transformer architecture [1] fundamentally transformed biomedical NLP: BioBERT [2], pre-trained on 4.5 billion words from PubMed and PMC, established state-of-the-art across nine biomedical tasks; GatorTron [3], trained on 82 billion clinical words, achieved 9.5% accuracy gains on medical question answering. Standalone GPT-4-Turbo achieves 73.4% on MedQA-USMLE [4] but generates unverifiable claims in approximately 39% of atomic statements under RAGAS evaluation, confirming the hallucination problem in clinical contexts.

Med-PaLM 2 [5] reached 86.5% through domain-specific instruction tuning, and Med-Gemini [6] achieved 91.1% via multimodal training; both require closed proprietary infrastructure. Lewis et al. [7] introduced the RAG paradigm, demonstrating that grounding LLM generation in retrieved evidence reduces hallucination while preserving fluency. Zakka et al. [8] applied RAG to clinical medicine through the Almanac system. Xiong et al. [9] benchmarked foundational RAG across medical datasets, reporting 71.6% MedQA accuracy using single-stage dense retrieval.

Hybrid retrieval combining dense and sparse signals consistently outperforms either modality in isolation. Reciprocal Rank Fusion provides parameter-free score aggregation [10]. Cross-encoder re-ranking further refines hybrid candidate lists through joint query-passage encoding. Medical Graph RAG [11] traverses UMLS and SNOMED-CT for multi-hop clinical inference. SELF-RAG [12] introduced adaptive retrieval, achieving accuracy gains on PubMedQA. The RAGAS framework [13] provides automated reference-free assessment of faithfulness, answer relevance, and context relevance. Despite these advances, no prior work has systematically compared all three RAG tiers under a unified experimental protocol with comprehensive ablation.

III. MEDRAG SYSTEM ARCHITECTURE

A. Architectural Overview

MedRAG adopts a layered pipeline architecture comprising five principal stages: (1) Query Processing, (2) Knowledge Retrieval, (3) Context Integration, (4) Response Generation, and (5) Verification and Attribution. Each stage is

parameterised differently across the three RAG configurations, enabling systematic progressive evaluation against a shared experimental baseline.

B. Knowledge Base Construction

The MedRAG knowledge index comprises 4.7 million passages sourced from 2.1 million PubMed abstracts, 890,000 StatPearls clinical summaries, 1.2 million clinical practice guidelines, and 510,000 UMLS concept definitions. All passages are chunked to 512 tokens with 64-token overlap, encoded using MedCPT [15], and indexed using FAISS IVF256. BM25 indices are maintained in parallel using Elasticsearch. Index construction requires 6.2 hours on an A100 GPU cluster.

C. Foundational RAG Configuration

The Foundational configuration implements single-stage MedCPT dense retrieval. An incoming clinical query is encoded into a 768-dimensional embedding vector. FAISS IVF256 retrieves the top-5 most semantically similar passages from the knowledge index in under 5 milliseconds. Retrieved passages are concatenated into a structured prompt and forwarded to the generator LLM. No query reformulation or re-ranking is applied, providing a clean baseline for measuring retrieval augmentation effects.

D. Optimised RAG Configuration

The Optimised configuration introduces hybrid retrieval, cross-encoder re-ranking, and context compression. BM25 and dense retrievers independently return top-20 candidates. Reciprocal Rank Fusion aggregates the two ranked lists: $RRF(d) = \sum 1 / (60 + rank_r(d))$. The merged top-20 passages are re-ranked by a MedBERT cross-encoder that jointly encodes each query-passage pair. The top-5 re-ranked passages undergo positional reordering and extractive context compression. Retrieval Precision@5 improves from 0.61 (Foundational) to 0.79.

E. Modular RAG Configuration

The Modular configuration introduces query decomposition, knowledge graph traversal, and iterative faithfulness verification. Complex clinical queries are decomposed into atomic sub-queries using UMLS terminology normalisation. Each sub-query activates three parallel retrievers: FAISS dense retrieval, BM25 sparse retrieval, and knowledge graph traversal across UMLS, SNOMED-CT, RxNorm, and DrugBank. RRF merges all three ranked lists; cross-encoder re-ranking selects top-5 passages. After generation, RAGAS faithfulness verification checks all atomic claims against retrieved context. Responses with faithfulness below 0.80 trigger a single-pass refinement with expanded retrieval.

F. Medical Knowledge Graph Integration

The knowledge graph encodes 3.5 million biomedical concepts from UMLS across 200 source vocabularies. Multi-hop path retrieval uses Breadth-First Search to identify shortest connecting paths up to four hops between query entities, enabling reasoning chains such as 'rifampicin → induces → CYP3A4 → metabolises → tacrolimus.' Community detection via the Louvain modularity algorithm identifies densely connected clinical concept clusters, enabling retrieval of thematically coherent evidence beyond simple pairwise entity relationships.

IV. EXPERIMENTAL METHODOLOGY

A. Hardware and Software Configuration

All experiments were conducted on two NVIDIA A100 80GB GPUs with 256 GB RAM and 4 TB NVMe RAID-0 storage. LLM inference was served using vLLM with tensor parallelism. Three generator backbones were evaluated: GPT-4-Turbo (API access), Mixtral-8x7B-Instruct (INT8 quantisation), and LLaMA-3-8B-Instruct (INT4 quantisation). Each of three RAG configurations was paired with each generator, and a no-RAG baseline was established for each, yielding 12 total experimental configurations. Each was evaluated across three independent random seeds.

B. Benchmark Datasets

Five datasets were used: (1) MedQA-USMLE—1,273 four-option MCQs from USMLE Steps 1–3; (2) MedMCQA—4,183 questions across 21 medical subjects [19]; (3) PubMedQA—500 biomedical research questions with yes/no/maybe answers [20]; (4) MMLU-Medical—1,089 questions across clinical knowledge, anatomy, and genetics [21]; and (5) the Chronic Pain Diagnostic Dataset (CPDD)—312 private clinical diagnostic cases. MedQA-USMLE accuracy is the primary comparative metric.

C. Evaluation Metrics

Retrieval quality was assessed using Precision@5, Recall@5, and Mean Reciprocal Rank (MRR). Generation quality was measured with BLEU-4, ROUGE-L, and BERTScore F1. RAG-specific quality was evaluated with RAGAS [13] across faithfulness, answer relevance, and context relevance. Clinical meaningfulness was assessed through the S.C.O.R.E. framework [14] administered by 47 blinded clinicians. Scalability was evaluated using Locust load testing at 50, 200, and 500 concurrent simulated users.

V. RESULTS AND DISCUSSION

A. Baseline LLM Performance

Standalone LLM baselines confirm the severity of the hallucination problem in clinical NLP. GPT-4-Turbo achieves 73.4% on MedQA-USMLE with RAGAS Faithfulness of 0.61, indicating approximately 39% of

atomic claims in generated responses are unverifiable from primary sources. Mixtral-8x7B and LLaMA-3-8B record 61.4% and 52.1% accuracy with faithfulness scores of 0.54 and 0.49 respectively. These baseline hallucination rates are clinically unacceptable for any scenario involving independent decision support.

B. Foundational RAG Results

Single-stage dense retrieval produces substantial improvements across all models. GPT-4-Turbo advances from 73.4% to 79.9% (+6.5 pp, $p < 0.001$). LLaMA-3-8B improves from 52.1% to 64.3% (+12.2 pp), with the smallest model recording the largest proportional gain, indicating that retrieval grounding disproportionately benefits knowledge-deficient generators. Mean RAGAS Faithfulness improves from 0.61 to 0.74, a 21% relative reduction in unsupported claims. Context Relevance averages 0.68.

C. Optimised RAG Results

Hybrid retrieval with cross-encoder re-ranking produces a further substantial accuracy increment. GPT-4-Turbo advances from 79.9% to 84.1%. Retrieval Precision@5 increases from 0.61 to 0.79, a 30% relative improvement attributable to cross-encoder re-ranking. Mean Faithfulness reaches 0.83. Adaptive retrieval (SELF-RAG protocol) achieves 72.4% on PubMedQA compared to 65.8% for fixed-retrieval Optimised RAG, confirming that selective retrieval reduces retrieval noise on abstractive research tasks.

D. Modular RAG Results

Configuration	Model	MedQA %	MMLU %	Fait h.	SCORE
No RAG	GPT-4-Turbo	73.4	75.3	0.61	2.94
No RAG	Mixtral-8x7B	61.4	62.1	0.54	2.61
No RAG	LLaMA-3-8B	52.1	54.8	0.49	2.38
Found. RAG	GPT-4-Turbo	79.9	81.2	0.74	3.89
Found. RAG	Mixtral-8x7B	69.5	70.4	0.70	3.61
Found. RAG	LLaMA-3-8B	64.3	65.7	0.67	3.42
Optim. RAG	GPT-4-Turbo	84.1	85.7	0.83	4.31
Optim. RAG	Mixtral-8x7B	74.8	75.9	0.79	4.07
Optim. RAG	LLaMA-3-8B	70.2	71.3	0.76	3.88
Modular RAG	GPT-4-Turbo	91.1	92.3	0.91	4.73
Modular RAG	Mixtral-8x7B	84.4	83.8	0.87	4.49
Modular RAG	LLaMA-3-8B	79.7	78.9	0.84	4.32
Configuration	MedQA (%)	Faithfulness	Δ Acc. (pp)		
Full Modular RAG (baseline)	91.1	0.91	—		
– Query Decomposition	87.4	0.88	–3.7		

- Knowledge Graph Retrieval	86.2	0.87	-4.9
- Cross-Encoder Re-ranking	85.1	0.86	-6.0
- Context Compression	88.7	0.89	-2.4
- Faithfulness Refinement	87.9	0.83	-3.2
- UMLS Normalisation	88.1	0.88	-3.0

The Modular configuration achieves the highest performance across all metrics. GPT-4-Turbo reaches 91.1% on MedQA-USMLE—a 24% relative error reduction from no-RAG and equal to Med-Gemini, which relies on closed multimodal infrastructure. LLaMA-3-8B with Modular RAG reaches 79.7%, directly exceeding standalone GPT-4-Turbo (73.4%). RAGAS Faithfulness reaches 0.91. The faithfulness refinement loop was triggered for 18.4% of initial drafts, elevating mean Faithfulness from 0.83 to 0.91. Clinician S.C.O.R.E. ratings improve from 2.94/5.0 (no-RAG) to 4.73/5.0 (Modular RAG + GPT-4), with 95.7% of evaluating clinicians finding Modular RAG useful for unfamiliar clinical topics.

TABLE I
COMPREHENSIVE PERFORMANCE
COMPARISON ACROSS ALL CONFIGURATIONS

E. Retrieval Quality Analysis

Retrieval quality metrics confirm incremental improvements across all configurations. Foundational RAG achieves Precision@5 of 0.61, Recall@5 of 0.54, and MRR of 0.68 at 42 ms latency. Optimised RAG improves to Precision@5 of 0.79, Recall@5 of 0.71, MRR of 0.82 at 187 ms. Modular RAG achieves Precision@5 of 0.88, Recall@5 of 0.83, MRR of 0.91 at 394 ms. The threefold latency increase from Foundational to Modular RAG remains within clinically acceptable bounds for non-emergency decision support scenarios where response quality outweighs sub-second latency requirements.

F. Ablation Study

The ablation analysis (Table II) reveals cross-encoder re-ranking as the single highest-impact component (-6.0 pp when removed), followed by knowledge graph retrieval (-4.9 pp), query decomposition (-3.7 pp), faithfulness refinement (-3.2 pp), UMLS normalisation (-3.0 pp), and context compression (-2.4 pp). These findings provide evidence-based guidance for resource-constrained deployments: cross-encoder re-ranking should be prioritised as the first architectural enhancement beyond single-stage dense retrieval, as it delivers the largest accuracy gain per unit of additional computational cost.

TABLE II
ABLATION STUDY: COMPONENT CONTRIBUTION TO MODULAR RAG PERFORMANCE

G. State-of-the-Art Comparison

System	MedQA Acc. (%)	Architecture	Open Source
Med-PaLM 2 [5]	86.5	Specialised LLM (closed)	No
Med-Gemini [6]	91.1	Multimodal LLM (closed)	No
Med42-v2 [16]	80.4	Fine-tuned LLaMA-3	Yes
Gazal-R1 [17]	87.1	RL-trained mid-size LLM	Yes
MEDRAG (Xiong) [9]	71.6	Foundational RAG	Yes
MedRAG-Modular (Ours)	91.1	Modular RAG + GPT-4-Turbo	Pipeline
MedRAG-OSS (Ours)	84.4	Modular RAG + Mixtral	Full

Table III positions MedRAG relative to contemporary medical AI systems. MedRAG-Modular achieves accuracy parity with Med-Gemini (91.1%) on MedQA-USMLE while remaining deployable as a fully transparent, open-source pipeline. MedRAG-OSS with Mixtral achieves 84.4%, outperforming Med42-v2 (80.4%) without requiring supervised fine-tuning on clinical datasets. This establishes sophisticated retrieval orchestration as a competitive alternative to domain-specific model adaptation, with significant advantages in interpretability and regulatory compliance.

TABLE III
COMPARISON WITH STATE-OF-THE-ART MEDICAL AI SYSTEMS

H. Key Findings

Four principal findings emerge. First, retrieval quality is the dominant performance determinant: the 11.2 pp gain from Foundational to Modular RAG exceeds the gain from scaling the generator from 8B to frontier scale within the same RAG tier. Second, sophisticated retrieval orchestration compensates for limited parametric knowledge: LLaMA-3-8B with Modular RAG exceeds standalone GPT-4-Turbo by 6.3 pp. Third, improving Precision@5 from 0.61 to 0.88 reduces hallucination from 26.3% to 9.2%—a 65% relative reduction from a 44% improvement in retrieval precision. Fourth, cross-encoder re-ranking is the single most impactful component and should be the priority enhancement in resource-constrained deployment scenarios.

VI. CHALLENGES AND LIMITATIONS

Despite strong empirical performance, several limitations merit acknowledgment. Modular RAG end-to-end latency of 5–9 seconds precludes application in time-critical emergency scenarios. Foundational RAG (42 ms retrieval, ~2 seconds end-to-end) is more appropriate for real-time clinical integration. Knowledge base currency requires active maintenance infrastructure: continuous automated ingestion with duplicate detection across pre-print and peer-reviewed publication stages represents a non-trivial engineering challenge that must be addressed before production deployment.

Demographic equity in retrieval performance remains an ongoing concern. Evidence systems may surface literature less representative of underrepresented clinical populations if the underlying knowledge base reflects historical publication biases. Active curation of diverse guideline sources and subgroup study populations is necessary. Formal regulatory compliance requires prospective clinical validation: RAGAS faithfulness metrics provide automated proxies for grounding quality, but FDA SaMD approval requires prospective trials with primary outcome measures of diagnostic accuracy and treatment appropriateness, which have not yet been completed for MedRAG.

VII. FUTURE SCOPE

Multimodal RAG integration represents the highest-priority near-term extension. Clinical decision-making routinely involves radiology images, pathology slides, ECG traces, and dermatological photography. Cross-modal embedding spaces built on BiomedCLIP and MedTrinity-25M foundations, combined with image-query retrieval mechanisms, would enable comprehensive multimodal evidence grounding unavailable to current text-only systems.

Federated privacy-preserving RAG architectures would enable each institution to maintain a local knowledge index from de-identified patient records without centralising sensitive data. A federated query protocol aggregating rankings from participating institution indexes via RRF, combined with homomorphic embedding encryption and differential privacy mechanisms, would substantially improve

regulatory compliance and institutional adoption. Personalised medicine integration through patient-adaptive retrieval, incorporating pharmacogenomic profiles and comorbidity patterns via privacy-preserving anonymisation gateways, and multilingual query processing supporting WHO official languages are additional planned development priorities.

VIII. CONCLUSION

This paper presented MedRAG, a healthcare conversational assistant grounded in a multi-tier Retrieval-Augmented Generation architecture. Three progressive configurations were designed, implemented, and evaluated across five medical benchmarks in 12 model-architecture combinations. The Modular RAG configuration achieves 91.1% accuracy on MedQA-USMLE with RAGAS Faithfulness of 0.91 and hallucination rate of 9.2%—a 24% relative error reduction from standalone LLM baseline and parity with closed-source Med-Gemini.

The primary scientific contribution is the empirical demonstration that retrieval quality is a more decisive performance determinant than generator model scale. LLaMA-3-8B with Modular RAG outperforms standalone GPT-4-Turbo, establishing that sophisticated retrieval orchestration can bridge much of the performance gap between resource-constrained and frontier systems at substantially reduced computational cost. Ablation analysis identifies cross-encoder re-ranking and knowledge graph integration as the highest-impact components, providing evidence-based guidance for clinical deployment prioritisation. MedRAG represents a technically mature, clinically meaningful step toward AI-augmented healthcare that is grounded in evidence, transparent in reasoning, equitable in accessibility, and aligned with emerging regulatory requirements for Software as a Medical Device.

REFERENCES

- [1] A. Vaswani et al., "Attention Is All You Need," in *Advances in Neural Information Processing Systems* (NeurIPS), vol. 30, 2017.
- [2] J. Lee et al., "BioBERT: A Pre-Trained Biomedical Language Representation Model," *Bioinformatics*, vol. 36, pp. 1234–1240, 2020.

- [3] X. Yang et al., "GatorTron: A Large Clinical Language Model," arXiv:2203.03540, 2022.
- [4] D. Jin et al., "What Disease does this Patient Have? A Large-scale Open Domain QA Dataset from Medical Exams," *Applied Sciences*, vol. 11, no. 14, 2021.
- [5] K. Singhal et al., "Towards Expert-Level Medical Question Answering with LLMs," arXiv:2305.09617, 2023.
- [6] K. Saab et al., "Capabilities of Gemini Models in Medicine," arXiv:2404.18416, 2024.
- [7] P. Lewis et al., "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," *NeurIPS*, vol. 33, pp. 9459–9474, 2020.
- [8] C. Zakka et al., "Almanac: Retrieval-Augmented Language Models for Clinical Medicine," *NEJM AI*, vol. 1, p. AIoa2300068, 2024.
- [9] G. Xiong et al., "Benchmarking Retrieval-Augmented Generation for Medicine," arXiv:2402.13178, 2024.
- [10] G. V. Cormack, C. L. A. Clarke, and S. Buettcher, "Reciprocal Rank Fusion Outperforms Condorcet and Individual Rank Learning Methods," in *SIGIR*, 2009.
- [11] J. Wu, J. Zhu, and Y. Qi, "Medical Graph RAG: Towards Safe Medical LLM via Graph RAG," arXiv:2408.04187, 2024.
- [12] A. Asai et al., "Self-RAG: Learning to Retrieve, Generate, and Critique," arXiv:2310.11511, 2023.
- [13] S. Es et al., "RAGAS: Automated Evaluation of Retrieval Augmented Generation," in *Proc. EACL*, pp. 150–158, 2024.
- [14] T. F. Tan et al., "A Proposed S.C.O.R.E. Evaluation Framework for LLMs," arXiv:2407.07666, 2024.
- [15] Q. Jin et al., "MedCPT: Contrastive Pre-Trained Transformers with PubMed Search Logs," *Bioinformatics*, 2023.
- [16] C. Christophe et al., "Med42-v2: A Suite of Clinical LLMs," arXiv:2408.06142, 2024.
- [17] A. M. Adly et al., "Gazal-R1: Achieving State-of-the-Art Medical Reasoning," arXiv:2506.21594, 2025.
- [18] Y. Gao et al., "Retrieval-Augmented Generation for Large Language Models: A Survey," arXiv:2312.10997, 2023.
- [19] A. Pal et al., "MEDMCQA: A Large-Scale Multi-Subject Multi-Choice Dataset," *PMLR ML for Health*, pp. 248–260, 2022.
- [20] Q. Jin et al., "PubMedQA: A Dataset for Biomedical Research QA," in *Proc. EMNLP*, pp. 2567–2577, 2019.
- [21] D. Hendrycks et al., "Measuring Massive Multitask Language Understanding," arXiv:2009.03300, 2020.
- [22] Z. Jiang et al., "Active Retrieval Augmented Generation (FLARE)," arXiv:2305.06983, 2023.
- [23] P. Xia et al., "MMed-RAG: Versatile Multimodal RAG System for Medical VLMs," arXiv:2410.13085, 2024.
- [24] M. Abo El-Enen et al., "A Survey on RAG Models for Healthcare Applications," *Neural Computing and Applications*, vol. 37, pp. 28191–28267, 2025.
- [25] D. Wang and S. Zhang, "Large Language Models in Medical and Healthcare Fields," *Artificial Intelligence Review*, vol. 57, p. 299, 2024.