

# Investigating LLM Capabilities on Long Context Comprehension for Medical Question Answering

Anonymous ACL submission

## Abstract

This study is the first to investigate LLM comprehension capabilities over long-context (LC), clinically relevant medical Question Answering (QA) beyond MCQA. Our comprehensive approach considers a range of settings based on content inclusion of varying size and relevance, LLM models of different capabilities and a variety of datasets across task formulations. We reveal insights on model size effects and their limitations, underlying memorization issues and the benefits of reasoning models, while demonstrating the value and challenges of leveraging the full long patient’s context. Importantly, we examine the effect of Retrieval Augmented Generation (RAG) on medical LC comprehension, uncovering best settings in single versus multi-document QA datasets. We shed light into some of the evaluation aspects using a multi-faceted approach uncovering common metric challenges. Our quantitative analysis reveals challenging cases where RAG excels while still showing limitations in cases requiring temporal reasoning.

## 1 Introduction

Large Language Models (LLMs) (Achiam et al., 2023; Dubey et al., 2024; Yang et al., 2024) perform impressively on medical tasks (Borgeaud et al., 2021; Nori et al., 2023), achieving super-human scores on United States Medical Licensing Examination (USMLE)-style exam question-answering (QA) (Chen et al., 2023b; Tang et al., 2023; Pal and Sankarasubbu, 2024; Singhal et al., 2025). Yet evaluation is mostly based on multiple-choice QA (MCQA) (Liévin et al., 2024; Xiong et al., 2024; Chen et al., 2024b; Singhal et al., 2025) which doesn’t reflect performance in complex tasks (Arias-Duart et al., 2025), such as open-ended medical QA (Sandeep Nachane et al., 2024). Furthermore, medical board exams designed to assess professional knowledge and decision-making rely primarily on textbook knowledge which is likely

available during pre-training (Chen et al., 2025a). Assessing LLM capabilities on context from expert-curated electronic health records (EHRs) could provide crucial signals on how well models address complexities in real data, including domain-specific vocabulary, heterogeneous document types, multi-document reasoning, data noise, linguistic diversity, long contexts, and long-range dependencies (Wornow et al., 2023). However, research on open-form QA pertaining to EHRs focuses on creating datasets (Yang et al., 2022) or EHR-specialized models (Fleming et al., 2023), rather than investigating LLM performance in real-world QA.

Moreover, most EHR-based benchmarks consider single-document contexts (Pampari et al., 2018; Yue et al., 2020), while in practice there is no guarantee that the necessary information will be included within a specific patient document. Reasoning over multiple long, temporally-dependent documents further poses questions around long-context (LC) limitations and evidence positioning sensitivity (Liu et al., 2023; Xiao et al., 2023), also observed in medically focused studies (Adams et al., 2024). Work on the effect of long-range dependencies within medical machine reading comprehension (MRC) is limited (Vatsal and Singh, 2024), focusing on shorter texts (~1.5-4K tokens) and span-based QA rather than complex EHRs. Work on LC medical QA involves either synthetic data, is MCQA-based and still relatively short (up to 6k) (Adams et al., 2024), or small in sample size and not human-validated (Fan et al., 2024). We address this gap by investigating LLM capabilities to answer questions given longitudinal patient-centric EHRs, where QA pairs are human-validated.

While recent LLM releases push the frontiers of context size, performance does not scale with increased context (Levy et al., 2024; Modarressi et al., 2025). Retrieval Augmented Generation (RAG) (Lewis et al., 2021) often forms a cost efficient alternative by selecting relevant context

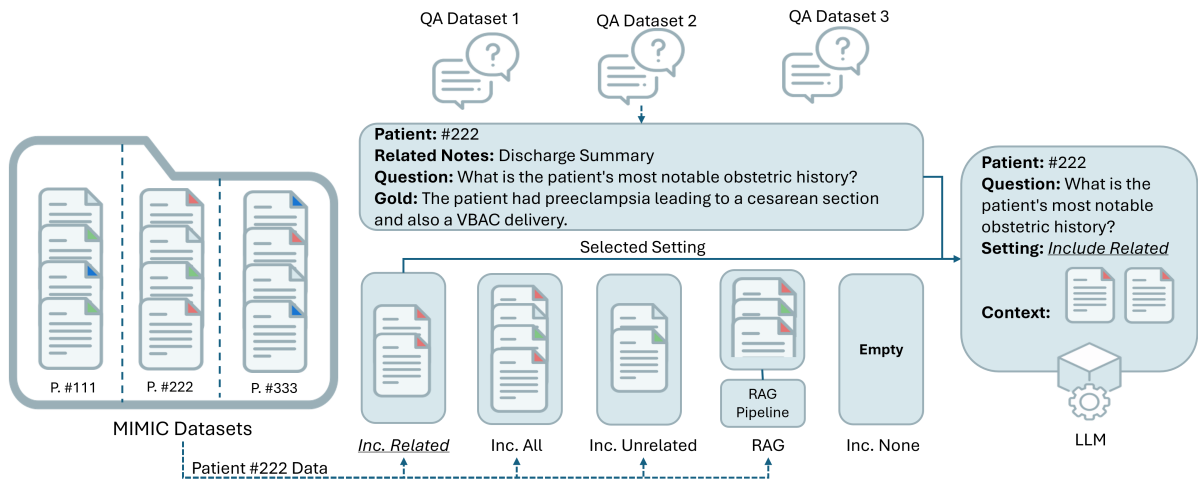


Figure 1: Overview of our Approach

on demand. What is preferable is subject to debate, with each showing benefits for different tasks and settings (Li et al., 2024b): RAG has been reported to enhance LLM performance compared to LC prompt compression for QA (Zhang et al., 2024a), outperforming SOTA LLMs on both 32K and 128K benchmarks (Xu et al., 2024a; Bai et al., 2024). By contrast, others showed that RAG performance peaks at 32K context length for large LLMs (Leng et al., 2024). While work on medical QA reports promising gains when retrieval is carefully designed (Xiong et al., 2024), LC vs RAG superiority in long-form medical QA remains an open question. We make the following contributions:

- We are the first to assess how LLMs of different underlying capabilities and sizes perform on long-context medical QA given longitudinal EHR patient notes, covering three task formulations: 1) MCQA, 2) Extractive and 3) Open-ended generative QA (§3.1) and a variety of LC and RAG settings, as depicted in Fig. 1.
- We perform comprehensive experiments evaluating different note inclusion strategies in the input context, including note memorization (§4.2).
- We carefully design a hybrid RAG pipeline examining its performance over Full-context (FC) across task formulations and context sizes, showing its superiority across tasks (§5.1, §5.2).
- Our experiments and multi-faceted evaluation methodology allow for a direct comparison between long-form reasoning when the answer is part of a single versus multiple documents (§5.1).
- Our quantitative and qualitative analyses uncover some of the LLM and metrics challenges in reasoning over long EHRs, while demonstrating

such cases where RAG excels (§5.2, §5.3).

## 2 Related Work

### 2.1 Leveraging Long Context

While increasing long-context window capabilities of LLMs, e.g. GPT-4o (Hurst et al., 2024) (128K), Claude 3 (Anthropic, 2023) (200K), Gemini 1.5 (Team et al., 2024a) (1 million), boost their long-context performance they still struggle in general purpose open-form QA (beyond short passage retrieval) (Zhang et al., 2024b; Chen et al., 2025b). Despite positional embedding techniques like RoPE (Su et al., 2024b), YaRN (Peng et al., 2023) and Position Interpolation (Chen et al., 2023a) allowing context extrapolation, QA task performance is low even for 32K contexts (Chen et al., 2025b). Furthermore, the performance gap between open-source and longer-context close-source models (Li et al., 2024a) poses questions around the potential avenues for patient-centric long-dependency data where privacy matters.

Assessment of long-context LLM performance (Dong et al., 2023; Bai et al., 2023; Liu et al., 2023; Zhang et al., 2024b; Hsieh et al., 2024), shows that apart from performance drops with increasing token length, LLMs are particularly challenged by long-range dependencies (Li et al., 2024a) that potentially require reasoning. Fan et al. (2024) show notable LLM performance drop in medical QA tasks with contexts up to 200K tokens, with open-source LLMs struggling to produce output given larger contexts. Although long and longitudinal context is prominent in real-world patient and medical data, there is little work investigating LLM capabilities in such settings.

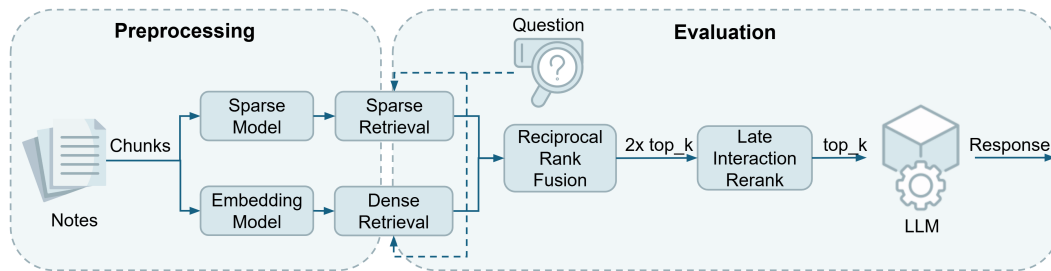


Figure 2: RAG Pipeline

## 2.2 Retrieval-augmented Generation (RAG)

RAG has been leveraged in a prompt-based plug-in manner to address black-box long-context challenges (Yu et al., 2023). Current retrieval strategies are primarily: chunk-based (Izacard et al., 2021), index-based (Liu, 2022) and summarization-based (Sarathi et al., 2024). LLM-based large-scale retrievers for chunk-based approaches like BGE M3 (Chen et al., 2024a) give notable performance improvement. Furthermore, approaches like Self-RAG (Asai et al., 2023) outperform ChatGPT in open-domain QA through on-demand retrieval followed by generation and self-reflection.

Biomedical retrievers showed benefits in domain retrieval, with the MedCPT dense semantic retriever and re-ranker (Jin et al., 2023), achieving top performance on 6 biomedical tasks surpassing LLM-based counterparts. The BMRETRIEVER (Xu et al., 2024b) further surpasses MedCPT’s performance on most downstream tasks. More recently, methods like Med-RAG (Zhao et al., 2025) showed better diagnostics over EHRs, through enhancement by knowledge graph-elicited reasoning, while Self-BioRAG (Jeong et al., 2024) employs MedCPT and generates through self-reflection using a domain-specific LLM. Yet Fan et al. (2025) underscore the medical retrieval challenges despite strong in-domain retrieval capabilities across large corpora, with models struggling with specialized medical content such as EHRs. Myers et al. (2025) further demonstrated variability in retrieval across EHR datasets. Therefore a thorough examination of RAG vs LC over realistic long-context patient-oriented QA is needed.

## 3 Approach

Since strong medical QA performance offers the potential of LLM integration in clinical workflows we aim to uncover LLM limitations given complex patient-oriented scenarios. We focus on **content-comprehension** settings using longitudinal EHRs,

thus emphasizing relevant patient-centric **long-context** across time and documents.

Our comprehensive approach spans across expert annotated datasets covering **three task formulations**: 1) MCQA, 2) Extractive, and 3) Open-ended QA. We assess each dataset across **four settings with varying note relevance** (see §4.2).

### 3.1 Dataset Selection Criteria

The following criteria were used for benchmark selection in our study (full comparison in Table 10):

**Beyond MCQA**: Medical MCQA benchmarks i.e. PubMedQA (Jin et al., 2019), MedQA (Jin et al., 2021), MedMCQA (Pal et al., 2022) evaluate option selection instead of grounded synthesis and risk pretraining leakage from textbooks, which can inflate scores, providing little signal towards noisy clinical settings. Our evaluation goes beyond MCQA to both extractive and open-ended QA.

**Long-context Documents**: To address LLM capabilities in handling real-world LC we seek benchmarks extending well-beyond 8K, with inclusion of patient-specific content.

**Expert Annotated QA pairs** grounded on multi-document EHRs were our choice, to avoid clinical reliability issues like non-expert curation (Fan et al., 2024) or synthetic data (Adams et al., 2024).

**Patient-centric Longitudinal Content**: Prioritizing clinically-relevant QA with reasoning, we focus on MIMIC-based EHR datasets (Johnson et al., 2016, 2023), offering timestamped notes across document types with consistent format.

### 3.2 RAG Approach

Retrievers are distinguished into *sparse* and *dense*. Sparse retrievers such as BM25 (Robertson et al., 2009) and Splade (Formal et al., 2021) use overlapping terms to match queries with snippets, while dense retrievers encode them into embeddings and match them based on semantic similarity. Recently,

using LLM-based embeddings for dense retrieval has become standard practice.

Research on biomedical MCQA (Wang et al., 2024) and biomedical document retrieval (Luo et al., 2022) has shown that hybrid approaches combining sparse and dense retrieval outperform single component counterparts. Additionally, non-hybrid experiments on EHR subsets of varying granularity (Fan et al., 2025) show no clear winner between sparse and dense retrievers. Finally, re-ranking of retrieved snippets also exhibits superior performance in biomedical applications (Jin et al., 2023; Wang et al., 2024; Sohn et al., 2025). Based on the above, we leverage a hybrid approach followed by re-ranking, presented in Figure 2. Clinical notes are segmented into 512-token chunks, forming the document collections  $D = \{d_1, d_2, \dots, d_n\}$ . The question respectively forms the query  $Q$ .

*Sparse retrieval:* we employ a lexical retriever as a ranking function of each chunk  $d_j$ <sup>1</sup>:

$$S_{\text{sparse}}(Q, d_j) = \sum_{i=1}^n \text{IDF}(q_i) \cdot \frac{f_{q_i, d_j}(k_1+1)}{f_{q_i, d_j+k_1} \left(1 - b + b \cdot \frac{|d_j|}{\text{avgdl}}\right)} \quad (1)$$

*Dense retriever:* we use the same Encoder,  $E$ , for the query and the chunks and capture semantic similarity of pairs using the cosine similarity:

$$S_{\text{dense}}(Q, d_j) = \frac{E(Q) \cdot E(d_j)}{\|E(Q)\| \|E(d_j)\|} \quad (2)$$

*Reciprocal Rank Fusion (RRF):* combines the top- $k$  chunks from each retriever accounting for each document’s rank, using a smoothing constant  $k_{RRF}$ :

$$\text{RRF}(d_j) = \sum_{s \in S} \frac{1}{k_{RRF} + \text{rank}_s(d_j)} \quad (3)$$

*Late Interaction Reranking:* measures how well every  $Q$ -token is semantically supported by at least one document token, to produce the final top- $k$  set:

$$S_{\text{MaxSim}}(Q, d_j) = \sum_{i=1}^{|Q|} \max_{1 \leq r \leq |d_j|} \langle \mathbf{q}_i, \mathbf{d}_{j,r} \rangle \quad (4)$$

*Document Ordering:* we sort retrieved chunks temporally rather than by retrieval score due to its superior long-context performance (Yu et al., 2024).

### 3.3 Evaluation Methodology

Evaluating medical QA systems has traditionally relied on automatic metrics and MCQA benchmarks (Jin et al., 2019, 2021; Pal et al., 2022),

<sup>1</sup>Here terms document and chunk are used interchangeably.

which probe memorized knowledge and underrepresent the complexity of clinical reasoning, and EHR-grounded context. Towards a clinically meaningful assessment, we adopt a multi-dimensional, reference-based scheme that complements lexical overlap with embedding-based semantic similarity and domain-adapted Natural Language Inference (NLI). In line with emerging clinical evaluation practices (e.g., MEDIC (Kanithi et al., 2024)) and the rise of LLM-as-a-judge, we employ a calibrated rubric and we prioritize widely available metrics while cross-checking signals for robustness. Our datasets comprise gold standard QA pairs that are expert verified/generated enabling clinically-centric evaluation:

**METEOR** (Banerjee and Lavie, 2005) to capture *surface-level* and *lexical similarity*.

**BERTScore** (Zhang et al., 2020), computed with Clinical BioBERT embeddings (Alsentzer et al., 2019) to assess *semantic similarity*.

**NLI Scores**, to capture the logical relationship between reference and candidate answers. We use a domain-adapted NLI model (Deka et al., 2023) and then measure *Factual Consistency* using the probability of non-contradiction (Song et al., 2024) and *Factual Precision* using the entailment probability.

**LLM-as-a-judge**, employed with a 5-point scale rubric on three aspects: *Correctness* capturing factual consistency with respect to the gold answer, penalizing contradictions while allowing compatible additional information, *Completeness* capturing recall by assessing the predicted answer’s coverage of information present in the reference, *Faithfulness* capturing precision assessing if the predicted answer only contains information that is supported by the reference.

**Accuracy** reported on the MCQA task.

## 4 Experiments

### 4.1 Datasets

We leverage MIMIC-III (Johnson et al., 2016) and MIMIC-IV (Johnson et al., 2023) comprising de-identified EHRs covering hospital admissions and ICU stays as the underlying context. They support single- and multi-note settings and include: clinical notes (CN), discharge summaries (DS) and radiology reports (RR) per patient over time.

Our four selected QA datasets only use QA pairs curated or validated by clinical experts to ensure medical validity. We provide a brief dataset

Dataset	Context Data	QA Pairs	Patients	Mean/Max Context	Task Types	Answer Location	Reasoning	Question Generation	Answer Annotation
CliniQG4QA RadQA	MIMIC-III	1,287 3,509	36 80	4K / 7K 6K / 14K	Extractive	Clinical Note Radiology Report	× Single-Note	Experts Human-generated	Clinical experts (3) Physicians (2)
EHR-DS-QA EHRNoteQA	MIMIC-IV	478 962	70 962	46K / 131K 9K / 39K	Open-ended MC,Open-ended	Clinical Note/Dis. Summary Dis. Summaries	× Multi-Note	LLMs LLM	Physician-verified (1) Clinician-refined (3)

Table 1: EHR QA Dataset Statistics. Mean/Max Context corresponds to the *Include All* setting.

overview below, also summarized in Table 1:

**CliniQG4QA** (Yue et al., 2020): An extractive span benchmark based on MIMIC-III. While most of its pairs are machine-generated (MG), the 1,287 QA test pairs are either generated or verified by three clinical experts using the MG questions for reference. We only use the expert-annotated set.

**RadQA** (Soni et al., 2022): An extractive dataset comprising 3,074 physician-crafted questions and 6,148 answer spans from the Findings and Impressions sections of RR in MIMIC-III. Many questions correspond to multiple sentences and require different types of reasoning. Annotation is carried out by two human annotators. After filtering out for *Unanswerable* pairs we obtain 3,509 pairs.

**EHR-DS-QA** (Kotschenreuther, 2024): A synthetic QA corpus of 156,599 pairs generated by two LLMs and guided by predefined prompt templates on DS and CN from MIMIC-IV. A subset of 506 pairs were physician-verified of which we retain 478, marked as correct.

**EHRNoteQA** (Kweon et al., 2024): Built on MIMIC-IV, it contains 962 QA pairs authored by GPT-4 and iteratively refined by three clinicians. Questions span a diverse set of topics, each corresponding to multiple DS (~2.3 per patient), while supporting both open-ended and MCQA formats.

## 4.2 Context Formulation

While most of the datasets are based on a single note/report (CliniQG4QA, RadQA, EHR-DS-QA) or up to three DS (EHRNoteQA), there is a large amount of *longitudinal patient supplementary* content that remains unused. Aiming for a realistic clinical setting where there is no a-priori knowledge of the most relevant note, we leverage the MIMIC resources to augment the context with each patient’s longitudinal notes<sup>2</sup>, resulting in long context from the entire patient’s note history (*Include All*). We then investigate how LLMs leverage long-form context to answer patient-specific questions.

We explore four data inclusion scenarios depicted in Fig. 1, allowing comprehensive analysis:

<sup>2</sup>MIMIC preprocessing is described in Appendix C.2.

1. **Exclude All:** Exclusion of all notes - assessing model note memorization.
2. **Exclude Relevant:** Exclusion of relevant notes - assessing usefulness of supplementary material.
3. **Include All:** Inclusion of all patient notes - assessing LC processing of varying importance.
4. **Include Related:** Inclusion of only the most relevant note(s) (DS, CN or RR) as marked in each dataset - assessing model’s ability to leverage relevant context information.

For the *Include All* setting we report context-length-based performance across four token bins: *Short context:* 0–8K, *Medium context:* 8–16K, *Large context:* 16–32K, *Extended context:* 32–128K.

## 4.3 Models and Experimental Setup

**LLMs:** To explore the effect of model size, domain specialization and reasoning capabilities we study four open Qwen2.5 (Team et al., 2024b) single-family models, ensuring a controlled comparison:

- *Qwen2.5-7B-Instruct*
- *Huatuogpt-o1-Qwen-7B* (Chen et al., 2024b): medical LLM based on Qwen2.5-7B-Instruct
- *Qwen2.5-32B-Instruct-128K*
- *QwQ:32B:* reasoning model

**Evaluation:** Through comparisons against Selene-8B (Alexandru et al., 2025) and Prometheus-8x7B v2.0 (Kim et al., 2024), we finally selected QWEN-2.5-32B-INSTRUCT for stability and agreement (see Appendix D.2).

**Retrieval Settings:** Based on biomedical QA retrieval research that demonstrates the superiority of combining sparse and dense retrievers (see §3.2), we designed a hybrid retrieval approach to properly handle the semantic, lexical, and query complexity. The selected model components are enlisted below:

- *Dense Retrieval* using the open embedding model Qwen3-Embedding-8B (Zhang et al., 2025).
- *Sparse Retrieval* using BM25 selected through two sparse retriever ablation (see Table 9).
- *Late Interaction Reranking* with Reason-ModernColBERT, a ColBERT (Khattab and

Model	Setting	EHRNoteQA				EHR-DS-QA			RadQA			CliniQG4QA		
		Open-ended			MC	Open-ended			Extractive			Extractive		
		LLM	NLI	F1	Acc.	LLM	NLI	F1	LLM	NLI	F1	LLM	NLI	F1
HuatuoGPT-o1 7B	Exclude All	<u>7.85</u>	<b>23.36</b>	<b>68.82</b>	51.08	<b>26.65</b>	<b>33.86</b>	72.21	<u>25.06</u>	<u>27.90</u>	<u>63.78</u>	<u>15.96</u>	<u>23.01</u>	<u>64.44</u>
Qwen2.5 7B		<b>9.73</b>	<u>18.18</u>	<u>68.41</u>	51.46	<u>24.79</u>	<u>33.00</u>	<u>72.35</u>	<b>39.08</b>	<b>29.24</b>	<b>70.73</b>	<b>27.02</b>	<b>26.30</b>	<b>74.06</b>
Qwen2.5 32B		0.71	11.3	66.56	<u>58.34</u>	20.16	20.23	<b>72.67</b>	3.86	2.88	10.35	0.11	0.19	2.16
QwQ 32B		2.68	12.23	66.98	<b>59.24</b>	24.16	23.8	71.70	1.51	6.44	15.29	0.93	5.62	14.77
HuatuoGPT-o1 7B	Exclude Related	<u>26.62</u>	<b>40.26</b>	72.57	58.35	28.49	<b>33.98</b>	71.84	<u>24.39</u>	<u>28.28</u>	<u>63.33</u>	<u>15.70</u>	<u>24.13</u>	<u>64.15</u>
Qwen2.5 7B		25.24	<u>28.64</u>	72.72	<u>59.86</u>	<b>29.9</b>	<u>28.17</u>	<u>73.35</u>	<b>38.86</b>	<b>29.07</b>	<b>70.72</b>	<b>27.01</b>	<b>25.95</b>	<b>74.09</b>
Qwen2.5 32B		25.59	23.67	<u>73.01</u>	58.72	22.16	22.35	73.15	3.77	2.78	10.02	0.13	0.20	2.17
QwQ 32B		<b>27.46</b>	23.01	<b>74.35</b>	<b>60.15</b>	27.19	23.98	<b>73.51</b>	1.14	7.11	17.48	0.99	4.84	15.75
HuatuoGPT-o1 7B	Include All	<u>72.60</u>	45.45	78.21	78.94	<u>59.76</u>	<b>63.66</b>	77.14	63.92	49.41	76.38	67.85	52.45	80.08
Qwen2.5 7B		70.49	38.11	79.95	78.81	<b>62.49</b>	61.00	<u>79.14</u>	62.56	<u>50.13</u>	76.65	68.89	54.13	81.27
Qwen2.5 32B		67.88	<b>55.33</b>	<b>81.50</b>	<b>90.97</b>	57.21	<u>61.96</u>	<b>79.59</b>	67.74	50.10	<b>77.55</b>	<b>80.49</b>	<u>59.22</u>	<u>83.8</u>
QwQ 32B		<b>75.52</b>	<u>47.14</u>	<u>81.02</u>	<u>89.25</u>	57.68	60.54	78.9	<u>67.46</u>	<b>53.26</b>	<u>77.47</u>	<u>78.95</u>	<b>66.00</b>	<b>83.94</b>
HuatuoGPT-o1 7B	Include Related	70.73	39.30	79.71	76.87	<u>63.46</u>	60.65	77.20	64.37	49.19	76.48	70.47	53.19	79.83
Qwen2.5 7B		74.44	39.15	80.78	80.9	<b>64.76</b>	<b>65.14</b>	<u>80.28</u>	63.01	50.25	76.71	69.06	54.87	81.29
Qwen2.5 32B		<u>76.01</u>	<b>61.41</b>	<b>82.48</b>	<b>90.33</b>	59.12	<u>64.86</u>	<b>80.70</b>	<u>67.76</u>	<u>50.78</u>	<u>77.67</u>	<b>79.81</b>	<u>59.21</u>	<u>83.55</u>
QwQ 32B		<b>82.03</b>	<u>47.19</u>	<u>82.36</u>	<u>87.57</u>	61.31	61.68	79.70	<b>67.79</b>	<b>55.48</b>	<b>77.69</b>	<u>79.74</u>	<b>65.86</b>	<b>84.35</b>

Table 2: Results for Full Context in each setting across datasets. **Bold** is best and underlined the second best model in each Setting and Metric. Red highlights the global best across all Settings per Metric. *LLM* corresponds to LLM Correctness, *NLI* to NLI Entailment and *F1* to BioBERT F1

Zaharia, 2020) model finetuned on the ReasonIR dataset (Shao et al., 2025) demonstrating high performance on reasoning-intensive retrieval benchmark (Su et al., 2024a).

We explored two chunk inclusion strategies: i) direct chunk inclusion (RAG), ii) hierarchical parent note inclusion (RAG HIR).

**Experimental Setup:** LLM inference and prompt formulation specifics are in Appendices A and C.3.

## 5 Results and Discussion

### 5.1 Main Results

Performance for the different inclusion settings across datasets and models is shown in Table 2.

**Model Size, Tasks and Inclusion Settings:** Consistent with general findings, larger models perform better across all tasks (Table 2). Qwen2.5:32B and QwQ:32B generally outperformed the 7B parameter models across most metrics and tasks. The performance gap was particularly pronounced in tasks requiring reasoning or information synthesis from multiple sources, namely EHRNoteQA and RadQA. However, smaller models performed better in settings excluding relevant information or completely removing the context. This suggests that they are more prone to memorization. Additionally, unlike the rest of the tasks, MC shows strong evidence of memorization with a performance of 59.24 in the ‘Exclude All’ setting. While MC demonstrates 90%+ performance, open-ended and extractive tasks remain challenging for LLMs.

The ‘Include Related’ setting outperforms ‘Include All’ on most datasets and metrics except for MC. Despite that, the ‘Exclude Related’ setting outperforms the ‘Exclude All’ by a margin across the board showing that the usefulness of the supplementary patient material and suggesting that while the ‘Include All’ setting includes useful additional information, LLMs struggle to process long context and identify the most relevant information.

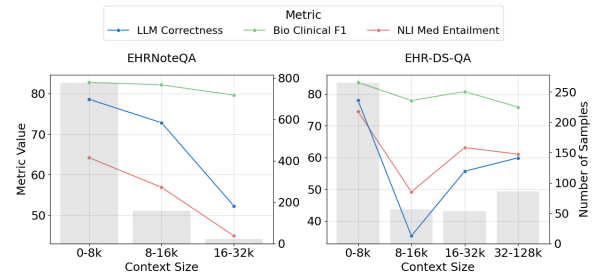


Figure 3: Metric performance over context size on ‘Include All’ for Qwen2.5:32B-Instruct.

**Reasoning and Fine-tuning:** The impact of reasoning-focused models and medical fine-tuning is mixed. The QwQ:32B general reasoning model showed improved results on Table 2 over other models overall on RadQA, EHRNoteQA - the reasoning datasets. This suggests that reasoning models can offer an advantage in tasks that demand complex inference across multiple sources of information, but this benefit may be more apparent in larger model sizes where reasoning capabilities can be effectively leveraged without compromis-

Model	Setting	EHRNoteQA				EHR-DS-QA		
		Open-ended			MC	Open-ended		
		LLM	NLI	F1	Acc.	LLM	NLI	F1
HuatuogPT-o1 7B	Include All	70.68	39.52	77.26	75.36	50.25	59.01	75.30
	RAG 5	59.94	<b>61.34</b>	78.44	<u>75.67</u>	57.57	61.53	<u>78.22</u>
	RAG 10	61.03	48.51	77.99	73.4	<u>58.08</u>	<u>62.45</u>	<b>78.61</b>
	RAG 15	60.84	52.53	77.92	74.97	<b>61.95</b>	<b>64.84</b>	77.68
	RAG HIR 3	70.13	56.91	77.99	70.81	50.06	53.82	76.17
	RAG HIR 5	<b>70.92</b>	<u>57.51</u>	<b>79.52</b>	72.38	48.46	55.21	76.66
	RAG HIR 7	<u>70.8</u>	54.61	<u>79.35</u>	<b>80.46</b>	55.99	57.12	77.35
Qwen2.5 7B	Include All	66.86	28.36	79.23	75.04	54.26	52.78	77.07
	RAG 5	59.54	40.33	78.82	75.67	54.17	56.11	<b>79.82</b>
	RAG 10	64.03	45.02	79.05	73.4	<u>57.11</u>	<b>62.76</b>	79.6
	RAG 15	58.07	<u>52.08</u>	79.20	75.6	<b>59.25</b>	<u>60.32</u>	<u>79.65</u>
	RAG HIR 3	65.54	46.34	<u>80.37</u>	76.93	50.73	53.12	78.65
	RAG HIR 5	<b>72.68</b>	<b>54.69</b>	80.17	<u>79.76</u>	53.59	52.56	77.66
	RAG HIR 7	<u>67.55</u>	46.41	<b>80.69</b>	<b>84.31</b>	52.10	52.39	78.26
Qwen2.5 32B	Include All	62.54	50.90	80.87	89.79	50.27	57.78	78.22
	RAG 5	60.35	46.36	80.28	80.46	<u>58.02</u>	58.91	80.29
	RAG 10	67.00	54.54	80.47	77.87	<b>60.52</b>	60.31	80.24
	RAG 15	<b>73.45</b>	53.94	80.84	81.02	55.68	<u>61.39</u>	<b>80.71</b>
	RAG HIR 3	67.48	<b>57.17</b>	<b>82.41</b>	83.36	55.13	<b>63.32</b>	<u>80.55</u>
	RAG HIR 5	68.86	<u>54.87</u>	81.59	<u>90.74</u>	51.20	60.73	80.12
	RAG HIR 7	<u>72.00</u>	47.33	<u>81.84</u>	<b>91.68</b>	52.82	57.27	80.21

Table 3: Open-ended and MCQA Results for Context of 8K+ tokens including LC and RAG Methods. **Bold** is best and underlined the second best performance per Model and Metric. Red highlights the global best across all models per Metric.

Model	Setting	RadQA			CliniQG4QA		
		LLM	NLI	F1	LLM	NLI	F1
HuatuogPT-o1 7B	Include All	35.60	<b>48.81</b>	75.61	43.73	50.44	78.97
	RAG 5	<b>36.61</b>	29.51	<b>76.61</b>	<b>44.20</b>	<b>54.51</b>	<b>79.54</b>
Qwen2.5 7B	Include All	36.10	<b>48.59</b>	75.92	44.55	53.26	80.04
	RAG 5	<b>37.09</b>	28.39	<b>77.10</b>	<b>47.03</b>	<b>58.36</b>	<b>81.73</b>

Table 4: Extractive QA Results for Context of 4K+ tokens. **Bold** is best performing per Model and Metric and red highlights the global best per Metric.

ing other essential skills. By contrast, *the medically fine-tuned reasoning model, HuatuogPT-o1-Qwen2.5:7B*, did not show any improvement over its standard instruction-tuned base model Qwen2.5:7B-instruct, in line with studies revealing that biomedical LLMs often lead to reduced performance (Dorfner et al., 2024; Dada et al., 2025).

**Context Size:** Fig. 3 shows results for Qwen2.5:32B-Instruct over different context sizes. Generally, performance decreases across metrics with increased context size on open-ended QA particularly when multi-note reasoning is required (EHRNoteQA), showing that even large LLMs are challenged by long context, with similar trends observed across models (see Appendix, Fig. 8). Worth noting that for context range (8-16K) EHR-DS-QA exhibits a dip in performance across all metrics and models. Such dip is likely attributed to data noise, based on manually observed inconsistencies and ambiguous gold answers in some pairs especially in EHR-DS-QA.

correlating with strange dips in performance at mid-range context windows. Fan et al. (2024) and Ma et al. (2025) observed similar dips.

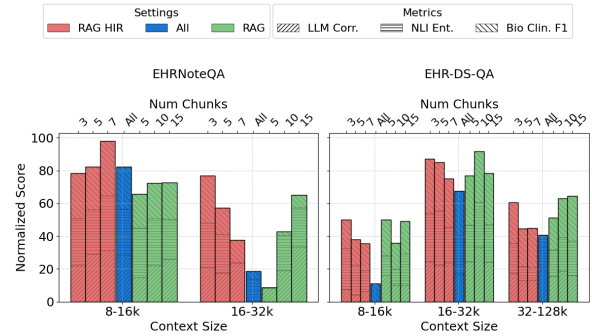


Figure 4: RAG performance across settings with Qwen2.5:32B-Instruct (min-max normalized).

**RAG Performance:** Table 3 shows the two evaluated chunk settings on the open-ended generative datasets<sup>3</sup>. RAG demonstrated clear performance improvement for the single-note generative task of EHR-DS-QA. In this scenario, the chunk-in-context strategy consistently provided the best results. For the more complex, multi-note reasoning task of EHRNoteQA, the hierarchical RAG strategy (RAG HIR) was the best with few exceptions. The results also show similar patterns on MC indicating clear performance improvement for all models on the RAG HIR 7 setting. Our findings are consistent across medium and larger (8K+) context sizes (Fig. 4). Worth noting that using RAG, the medically fine-tuned reasoning model HuatuogPT-o1-Qwen2.5:7B scored the best results on multiple metrics across the open-ended datasets, showing that RAG could benefit even more so the smaller models that struggle in Full-context (FC) settings ('Include All').

In Table 4 we only evaluate 7B sized models using a single RAG setting, due to the task nature and limited context size of the extractive datasets. While the extractive datasets have shorter contexts overall RAG still yields the best performance.

## 5.2 Quantitative Analysis: RAG vs LC

In our analysis we gathered challenging cases of disagreement between metrics (LLM Correctness, NLI Entailment, and Bio-F1) across EHRNoteQA and EHR-DS-QA. Metrics were first z-score standardized and samples were gathered for cases where metrics belonged to the top 60th percentile while

<sup>3</sup>We excluded QwQ:32B due to time and budget constraints

Category	Total	Favored by RAG	Favored by FC
SEE	18	12 (66.7%)	6 (33.3%)
TLR	10	3 (30.0%)	7 (70.0%)
CNS	10	6 (60.0%)	4 (40.0%)
ANI	3	2 (66.7%)	1 (33.3%)
<b>Total</b>	<b>41</b>	<b>23 (56.1%)</b>	<b>18 (43.9%)</b>

Table 5: Breakdown of sampled favored by RAG vs FC across four categories in open-ended QA.

another fell within the bottom 40th percentile. The 41 examples gathered were then manually examined in order to categorize them across QA types and determine which answer between RAG and FC is more reliable. The rubric is provided below:

- **Specific Entity Extraction (SEE):** Requiring retrieval of high cardinality facts such as medications and lab values or surgical procedures.
- **Temporal & Longitudinal Reasoning (TLR):** Requiring tracking of changes over time, chronological event ordering or visit comparisons.
- **Clinical Nuance & Status (CNS):** Questions regarding discharge instructions, patient mental status, or synthesized clinical course summaries.
- **Absence or Negative Information (ANI):** Requiring the identification of absent information.

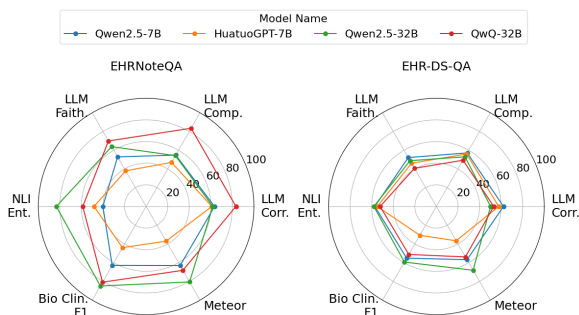


Figure 5: Model performance across metrics for open-ended QA, normalized with min-max per metric.

In Table 5, RAG shows higher performance across the SEE, CNS and ANI categories, while FC leads in the TLR one, demonstrating that while RAG effectively preserves the relevant information for challenging cases, it falls behind when temporal reasoning is required. A further qualitative analysis providing insights on cases where RAG or FC perform better is provided in Appendix G<sup>4</sup>.

### 5.3 Metrics and Insights

**Metric Comparisons:** Fig. 5 provides insights into the relation between models, metrics, and datasets

<sup>4</sup>manually analyzed failure cases in Appendix F.

while focusing on the open-ended generative tasks for the ‘Include All’ and ‘Include Related’ settings combined. *Instruct models lead on the single-note non-reasoning task (EHR-DS-QA), while larger models lead on the multi-note reasoning one (EHRNoteQA).* Instruct models are better based on semantic and NLI metrics, while reasoning ones score better on LLM-as-a-judge metrics<sup>5</sup>.

**Metric Insights:** We performed a qualitative analysis to investigate the reasons behind metrics’ disagreements on our open-ended datasets (see Appendix E.2). The 233 disagreements in EHRNoteQA and 47 in EHR-DS-QA are mainly due to:

- Correct/Complete but Unfaithful answers with addition of unsupported details in both datasets but more pronounced in EHR-DS-QA.
- Correct answer but low surface overlap (lower Meteor/Bert) in both datasets.
- Correct answer but low NLI entailment occurs due to negation/clinical phrasing in both datasets.
- Several short, risk/negation items where NLI flags contradiction despite clinically aligned answers in EHR-DS-QA.

These point to the need for thinking about more reasoning appropriate metrics.

## 6 Conclusion

We studied long-context, patient-centric clinical QA across EHR-grounded datasets, contrasting Full-Context (FC) prompting with Retrieval-Augmented generation (RAG). Our work shows that highly-relevant context often outperforms feeding all notes, likely due to LLMs struggling with long context. Hybrid RAG pipelines with reranking overperform FC, especially for specific factual queries and multi-note synthesis. Larger models generally help especially for reasoning datasets, but performance remains sensitive to context length and task formulation, while medically fine-tuned models fall behind. Future directions include: (i) scaling beyond single-note or single-source assumptions toward richer multi-note, multi-visit reasoning; (ii) advancing temporal and causal reasoning over longitudinal records, with explicit timeline grounding; (iii) strengthening evaluation via clinically faithful, judge-robust rubrics and cross-metric reliability analyses.

<sup>5</sup>We include all-context and long-context correlations between metrics in Appendix E.1

## 584 Limitations

585 Our work focuses on MIMIC-derived datasets,  
586 which represents English-only data in a single U.S.  
587 medical center. As such, performance may not  
588 generalize to other languages, populations, and  
589 healthcare systems. Furthermore, the de-identified  
590 nature of the data limits us from examining the  
591 potential cultural and other biases of LLMs over  
592 long contexts. Importantly, although our work  
593 presents an early step towards benchmarking long-  
594 context LLM performance in longitudinal patient-  
595 centric settings, our findings are intended solely  
596 for research purposes. Results reflect model per-  
597 formance on question answering benchmarks and  
598 should not be interpreted as guarantees of clinical  
599 safety, equitable performance, or readiness for  
600 clinical deployment.

601 Beyond complexities of modeling clinical notes,  
602 real-world records comprise heterogeneous data, in-  
603 cluding data from other sources and in other modal-  
604 ities. Our work poses limitations in the assessment  
605 of long-context by focusing only on the textual  
606 modality, something we aim to address in future  
607 work.

608 Despite the in depth study of literature and our  
609 initial ablations in selecting a competitive RAG  
610 methodology, our work does not exhaustively ex-  
611 amine the full potential of RAG under different  
612 settings, i.e. different retrievers and chunk sizes, in  
613 long-form medical QA. Finally, while we focus on  
614 Qwen-based models due to their competitive per-  
615 formance, we acknowledge that an evaluation of a  
616 wider range of LLMs could offer a more complete  
617 picture of the LLM landscape in long-form medical  
618 QA. Our analysis has also pointed to the limitations  
619 of current evaluation metrics, both in terms of as-  
620 sessing reasoning as well as lack of transparency in  
621 LLM as a judge metrics and how faithfulness can  
622 be truly assessed currently when additional infor-  
623 mation, not present in the benchmark, is added.

## 624 Ethics Statement

625 This work uses MIMIC-III, MIMIC-IV, and four  
626 derivative datasets (EHRNoteQA, EHR-DS-QA,  
627 RadQA, CliniQG4QA) accessed under the Phys-  
628 ioNet Credentialed Health Data License 1.5.0. All  
629 datasets contain de-identified patient records in ac-  
630 cordance with HIPAA Safe Harbor standards. We  
631 exclusively used open-weight LLMs to ensure no  
632 patient data were transmitted to third-party prop-  
633 erty systems.

## References

OpenAI Josh Achiam, Steven Adler, Sandhini Agarwal, 635  
Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale- 636  
man, Diogo Almeida, Janko Altenschmidt, Sam Alt- 637  
man, Shyamal Anadkat, Red Avila, Igor Babuschkin, 638  
Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim 639  
ing Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, 640  
Jake Berdine, Gabriel Bernadett-Shapiro, Christo- 641  
pher Berner, Lenny Bogdonoff, Oleg Boiko, Made 642  
laine Boyd, Anna-Luisa Brakman, Greg Brockman, 643  
Tim Brooks, Miles Brundage, Kevin Button, Trevor 644  
Cai, Rosie Campbell, Andrew Cann, Brittany Carey, 645  
Chelsea Carlson, Rory Carmichael, Brooke Chan, 646  
Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, 647  
Ruby Chen, Jason Chen, Mark Chen, Benjamin 648  
Chess, Chester Cho, Casey Chu, Hyung Won Chung, 649  
Dave Cummings, Jeremiah Currier, Yunxing Dai, 650  
Cory Decareaux, Thomas Degry, Noah Deutsch, 651  
Damien Deville, Arka Dhar, David Dohan, Steve 652  
Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, 653  
Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, 654  
Sim'on Posada Fishman, Juston Forte, Isabella Ful- 655  
ford, Leo Gao, Elie Georges, Christian Gibson, Vik 656  
Goel, Tarun Gogineni, Gabriel Goh, Raphael Gontijo- 657  
Lopes, Jonathan Gordon, Morgan Grafstein, Scott 658  
Gray, Ryan Greene, Joshua Gross, Shixiang Shane 659  
Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, 660  
Yuchen He, Mike Heaton, Johannes Heidecke, Chris 661  
Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, 662  
Brandon Houghton, Kenny Hsu, Shengli Hu, Xin 663  
Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, 664  
Joanne Jang, Angela Jiang, Roger Jiang, Haozhun 665  
Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo 666  
Jun, Tomer Kaftan, Lukasz Kaiser, Ali Kamali, Ing- 667  
mar Kanitscheider, Nitish Shirish Keskar, Tabarak 668  
Khan, Logan Kilpatrick, Jong Wook Kim, Christina 669  
Kim, Yongjik Kim, Hendrik Kirchner, Jamie Ryan 670  
Kiros, Matthew Knight, Daniel Kokotajlo, Lukasz 671  
Kondraciuk, Andrew Kondrich, Aris Konstantini- 672  
dis, Kyle Kopic, Gretchen Krueger, Vishal Kuo, 673  
Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, 674  
Jade Leung, Daniel Levy, Chak Li, Rachel Lim, 675  
Molly Lin, Stephanie Lin, Ma teusz Litwin, Theresa 676  
Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, 677  
Kim Malfacini, Sam Manning, Todor Markov, Yaniv 678  
Markovski, Bianca Martin, Katie Mayer, Andrew 679  
Mayne, Bob McGrew, Scott Mayer McKinney, 680  
Christine McLeavey, Paul McMillan, Jake McNeil, 681  
David Medina, Aalok Mehta, Jacob Menick, Luke 682  
Metz, An drey Mishchenko, Pamela Mishkin, Vinnie 683  
Monaco, Evan Morikawa, Daniel P. Mossing, Tong 684  
Mu, Mira Murati, Oleg Murk, David M'ely, Ashvin 685  
Nair, Reiichiro Nakano, Rajevee Nayak, Arvind Neel- 686  
akantan, Richard Ngo, Hyeonwoo Noh, Ouyang 687  
Long, Cullen O'Keefe, Jakub W. Pachocki, Alex 688  
Paino, Joe Palermo, Ashley Pantuliano, Giambattista 689  
Parascandolo, Joel Parish, Emy Parparita, Alexandre 690  
Passos, Mikhail Pavlov, Andrew Peng, Adam Perel- 691  
man, Filipe de Avila Belbute Peres, Michael Petrov, 692  
Henrique Pondé de Oliveira Pinto, Michael Pokorny, 693  
Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, 694  
Alethea Power, Boris Power, Elizabeth Proehl, Raul 695  
Puri, Alec Radford, Jack W. Rae, Aditya Ramesh, 696

697	Cameron Raymond, Francis Real, Kendra Rimbach,	Liu, Aohan Zeng, Lei Hou, et al. 2023. Longbench:	755
698	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder,	A bilingual, multitask benchmark for long context	756
699	Mario D. Saltarelli, Ted Sanders, Shibani Santurkar,	understanding. <i>arXiv preprint arXiv:2308.14508</i> .	757
700	Girish Sastry, Heather Schmidt, David Schnurr, John		
701	Schulman, Daniel Selsam, Kyla Sheppard, Toki	Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xi-	758
702	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	aozhi Wang, Xin Lv, Shulin Cao, Jiazheng Xu, Lei	759
703	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,	Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024.	760
704	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	Longbench v2: Towards deeper understanding and	761
705	Sokolowsky, Yang Song, Natalie Staudacher, Fe-	reasoning on realistic long-context multitasks. <i>ArXiv</i> ,	762
706	lipe Petroski Such, Natalie Summers, Ilya Sutskever,	abs/2412.15204.	763
707	Jie Tang, Nikolas A. Tezak, Madeleine Thompson,		
708	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	Satanjeev Banerjee and Alon Lavie. 2005. {METEOR}:	764
709	Preston Tuggle, Nick Turley, Jerry Tworek, Juan	An Automatic Metric for {MT} Evaluation with Im-	765
710	Felipe Cer'on Uribe, Andrea Vallone, Arun Vi-	proved Correlation with Human Judgments. In <i>Pro-</i>	766
711	jayvergiya, Chelsea Voss, Carroll L. Wainwright,	<i>ceedings of the {ACL} Workshop on Intrinsic and</i>	767
712	Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan	<i>Extrinsic Evaluation Measures for Machine Transla-</i>	768
713	Ward, Jason Wei, CJ Weinmann, Akila Welihinda,	<i>tion and/or Summarization</i> , pages 65–72, Ann Arbor,	769
714	Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wi-	Michigan. Association for Computational Linguis-	770
715	ethoff, Dave Willner, Clemens Winter, Samuel Wol-	tics.	771
716	rich, Hannah Wong, Lauren Workman, Sherwin Wu,		
717	Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo,	Asma Ben Abacha and Dina Demner-Fushman. 2019. A	772
718	Kevin Yu, Qim ing Yuan, Wojciech Zaremba, Rowan	question-entailment approach to question answering.	773
719	Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao,	<i>BMC Bioinform.</i> , 20(1):511:1–511:23.	774
720	Tianhao Zheng, Juntang Zhuang, William Zhuk, and		
721	Barret Zoph. 2023. <a href="#">Gpt-4 technical report</a> .	Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann,	775
		Trevor Cai, Eliza Rutherford, Katie Millican, George	776
722	Lisa C. Adams, Felix Busch, Tianyu Han, Jean-Baptiste	van den Driessche, Jean-Baptiste Lespiau, Bogdan	777
723	Excoffier, Matthieu Ortala, Alexander Loser, Hugo J.	Damoc, Aidan Clark, Diego de Las Casas, Aurelia	778
724	W. L. Aerts, Jakob Nikolas Kather, Daniel Truhn, and	Guy, Jacob Menick, Roman Ring, T. W. Hennigan,	779
725	Keno Kyrill Bressme. 2024. <a href="#">Longhealth: A question</a>	Saffron Huang, Lorenzo Maggiore, Chris Jones, Al-	780
726	<a href="#">answering benchmark with long clinical documents</a> .	bin Cassirer, Andy Brock, Michela Paganini, Geof-	781
727	<i>ArXiv</i> , abs/2401.14490.	frey Irving, Oriol Vinyals, Simon Osindero, Karen	782
		Simonyan, Jack W. Rae, Erich Elsen, and L. Sifre.	783
728	Andrei Alexandru, Antonia Calvi, Henry Broomfield,	2021. <a href="#">Improving language models by retrieving from</a>	784
729	Jackson Golden, Kyle Dai, Mathias Leys, Maurice	<a href="#">trillions of tokens</a> . In <i>International Conference on</i>	785
730	Burger, Max Bartolo, Roman Engeler, Sashank Pisu-	<i>Machine Learning</i> .	786
731	pati, et al. 2025. <a href="#">Atla selene mini: A general purpose</a>		
732	<a href="#">evaluation model</a> . <i>arXiv preprint arXiv:2501.17195</i> .	Hanjie Chen, Zhouxiang Fang, Yash Singla, and Mark	787
		Dredze. 2025a. Benchmarking large language mod-	788
733	Emily Alsentzer, John Murphy, William Boag, Wei-	els on answering and explaining challenging medical	789
734	Hung Weng, Di Jindi, Tristan Naumann, and	questions. In <i>Proceedings of the 2025 Conference</i>	790
735	Matthew McDermott. 2019. Publicly Available Clin-	<i>of the Nations of the Americas Chapter of the Asso-</i>	791
736	ical BERT Embeddings. In <i>Proceedings of the 2nd</i>	<i>ciation for Computational Linguistics: Human Lan-</i>	792
737	<i>Clinical Natural Language Processing Workshop</i> ,	<i>guage Technologies (Volume 1: Long Papers)</i> , pages	793
738	pages 72–78.	3563–3599.	794
739	Anthropic. 2023. Model card and evaluations	Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu	795
740	for claude models. <a href="https://www-files.anthropic.com/production/images/Model-Card-Claude-2.pdf">https://www-files.anthropic.</a>	Lian, and Zheng Liu. 2024a. <a href="#">Bge m3-embedding:</a>	796
741	<a href="https://www-files.anthropic.com/production/images/Model-Card-Claude-2.pdf">com/production/images/Model-Card-Claude-2.pdf</a> .	<a href="#">Multi-lingual, multi-functionality, multi-granularity</a>	797
		<a href="#">text embeddings through self-knowledge distillation</a> .	798
742	Anna Arias-Duart, Pablo Agustin Martin-Torres, Daniel	In <i>Annual Meeting of the Association for Computa-</i>	799
743	Hinjos, Pablo Bernabeu-Perez, Lucia Urcelay Ganza-	<i>tional Linguistics</i> .	800
744	bal, Marta Gonzalez Mallo, Ashwin Kumar Gururaj-		
745	an, Enrique Lopez-Cuena, Sergio Alvarez-Napagao,	Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang,	801
746	and Dario Garcia-Gasulla. 2025. Automatic evalua-	Wanlong Liu, Rongsheng Wang, Jianye Hou, and	802
747	tion of healthcare llms beyond question-answering.	Benyou Wang. 2024b. <a href="#">Huatuogpt-o1, towards med-</a>	803
748	<i>arXiv preprint arXiv:2502.06666</i> .	<a href="#">ical complex reasoning with llms</a> . <i>arXiv preprint</i>	804
		<i>arXiv:2412.18925</i> .	805
749	Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and	Pei Chen, Hongye Jin, Cheng-Che Lee, Rulin Shao,	806
750	Hannaneh Hajishirzi. 2023. <a href="#">Self-rag: Learning to</a>	Jingfeng Yang, Mingyu Zhao, Zhaoyu Zhang, Qin Lu,	807
751	<a href="#">retrieve, generate, and critique through self-reflection</a> .	Kaiwen Men, Ning Xie, et al. 2025b. Longleader: A	808
752	<i>ArXiv</i> , abs/2310.11511.	comprehensive leaderboard for large language mod-	809
		els in long-context scenarios. In <i>Proceedings of</i>	810
753	Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu,	<i>the 2025 Conference of the Nations of the Ameri-</i>	811
754	Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao	<i>cas Chapter of the Association for Computational</i>	812

813	<i>Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 8734–8750.	
814		
815	Shouyuan Chen, Sherman Wong, Liangjian Chen, and	
816	Yuandong Tian. 2023a. Extending context window	
817	of large language models via positional interpolation.	
818	<i>arXiv preprint arXiv:2306.15595</i> .	
819	Zeming Chen, Alejandro Hernández Cano, Angelika	
820	Romanou, Antoine Bonnet, Kyle Matoba, Francesco	
821	Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf,	
822	Amirkeivan Mohtashami, et al. 2023b. Meditron-	
823	70b: Scaling medical pretraining for large language	
824	models. <i>arXiv preprint arXiv:2311.16079</i> .	
825	Amin Dada, Osman Alperen Koraş, Marie Bauer, Jean-	
826	Philippe Corbeil, Amanda Butler Contreras, Con-	
827	stantin Marc Seibold, Kaleb E Smith, Jens Kleesiek,	
828	et al. 2025. Does biomedical training lead to better	
829	medical performance? In <i>Proceedings of the Fourth</i>	
830	<i>Workshop on Generation, Evaluation and Metrics</i>	
831	<i>(GEM<sup>2</sup>)</i> , pages 46–59.	
832	Pritam Deka, Anna Jurek-Loughrey, et al. 2023. Mul-	
833	tiiple evidence combination for fact-checking of	
834	health-related information. In <i>The 22nd Workshop</i>	
835	<i>on Biomedical Natural Language Processing and</i>	
836	<i>BioNLP Shared Tasks</i> , pages 237–247.	
837	Zican Dong, Tianyi Tang, Junyi Li, Wayne Xin Zhao,	
838	and Ji-Rong Wen. 2023. Bamboo: A comprehen-	
839	sive benchmark for evaluating long text modeling	
840	capacities of large language models. <i>arXiv preprint</i>	
841	<i>arXiv:2309.13345</i> .	
842	Felix J Dorfner, Amin Dada, Felix Busch, Marcus R	
843	Makowski, Tianyu Han, Daniel Truhn, Jens Kleesiek,	
844	Madhumita Sushil, Jacqueline Lammert, Lisa C	
845	Adams, et al. 2024. Biomedical large languages mod-	
846	els seem not to be superior to generalist models on un-	
847	seen medical data. <i>arXiv preprint arXiv:2408.13833</i> .	
848	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,	
849	Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,	
850	Akhil Mathur, Alan Schelten, Amy Yang, Angela	
851	Fan, Anirudh Goyal, Anthony S. Hartshorn, Aobo	
852	Yang, Archi Mitra, Archie Sravankumar, Artem Ko-	
853	renev, Arthur Hinsvark, Arun Rao, Aston Zhang,	
854	Aur’elien Rodriguez, Austen Gregerson, Ava Spataru,	
855	Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie	
856	Chern, Charlotte Caucheteux, Chaya Nayak, Chloe	
857	Bi, Chris Marra, Chris McConnell, Christian Keller,	
858	Christophe Touret, Chunyang Wu, Corinne Wong,	
859	Cristian Cantón Ferrer, Cyrus Nikolaidis, Damien Al-	
860	lonsius, Daniel Song, Danielle Pintz, Danny Livshits,	
861	David Esiobu, Dhruv Choudhary, Dhruv Mahajan,	
862	Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,	
863	Egor Lakomkin, Ehab A. AlBadawy, Elina Lobanova,	
864	Emily Dinan, Eric Michael Smith, Filip Radenovic,	
865	Frank Zhang, Gabriele Synnaeve, Gabrielle Lee,	
866	Georgia Lewis Anderson, Graeme Nail, Grégoire	
867	Mialon, Guanglong Pang, Guillem Cucurell, Hai-	
868	ley Nguyen, Hannah Korevaar, Hu Xu, Hugo Tou-	
869	vron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M.	
870	Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet,	
	Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park,	871
	Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,	872
	Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy	873
	Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie	874
	Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo	875
	Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe,	876
	Ju-Qing Jia, Kalyan Vasuden Alwala, K. Upasani,	877
	Kate Plawiak, Keqian Li, Ken-591 neth Heafield,	878
	Kevin R. Stone, Khalid El-Arini, Krithika Iyer, Kshi-	879
	tiz Malik, Kuen Iey Chiu, Kunal Bhalla, Lauren	880
	Rantala-Yearly, Laurens van der Maaten, Lawrence	881
	Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish	882
	Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat,	883
	Luke de Oliveira, Madeline Muzzi, Mahesh Pa-	884
	supuleti, Mannat Singh, Manohar Paluri, Marcin Kar-	885
	das, Mathew Oldham, Mathieu Rita, Maya Pavlova,	886
	Melissa Hall Melanie Kambadur, Mike Lewis, Min	887
	Si, Mitesh Kumar Singh, Mona Hassan, Naman	888
	Goyal, Narjes Torabi, Niko Iay Bashlykov, Nikolay	889
	Bogoychev, Niladri S. Chatterji, Olivier Duchenne,	890
	Onur cCelebi, Patrick Alrassy, Pengchuan Zhang,	891
	Pengwei Li, Petar Vasić, Peter Weng, Prajjwal Bhar-	892
	gava, Pratik Dubal, Praveen Krishnan, Punit Singh	893
	Koura, Puxin Xu, Qing He, Qingxiao Dong, Ra-	894
	gavan Srinivasan, Raj Ganapathy, Ramon Calderer,	895
	Ricardo Silveira Cabral, Robert Stojnic, Roberta	896
	Raileanu, Rohit Girdhar, Rohit Patel, Ro main	897
	Sauvestre, Ron nie Polidoro, Roshan Sumbaly, Ross	898
	Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar	899
	Hosseini, Sa hana Chennabasappa, Sanjay Singh,	900
	Sean Bell, Seohyun Sonia Kim, Sergey Edunov,	901
	Shaoliang Nie, Sharan Narang, Sharath Chandra	902
	Raparthi, Sheng Shen, Shengye Wan, Shruti Bho-	903
	sale, Shun Zhang, Simon Vandenhende, Soumya Ba-	904
	tra, Spencer Whitman, Sten Sootla, Stephane Collot,	905
	Suchin Gururangan, Sydney Borodinsky, Tamar Her-	906
	man, Tara Fowler, Tarek Sheasha, Thomas Georgiou,	907
	Thomas Scialom, Tobias Speckbacher, Todor Mi-	908
	haylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami,	909
	Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez,	910
	Vincent Gonguet, Vir ginie Do, Vish Vogeti, Vladan	911
	Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin	912
	Fu, Whit ney Meers, Xavier Martinet, Xiaodong	913
	Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao	914
	Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur,	915
	Yasmine Babaei, Yiqian Wen, Yiwen Song, Yuchen	916
	Zhang, Yue Li, Yuning Mao, Zacharie Delpierre	917
	Coudert, Zhengxu Yan, Zhengxing Chen, Zoe Pa-	918
	pakipos, Aaditya K. Singh, Aaron Grattafiori, Abha	919
	Jain, Adam Kelsey, Adam Shajnfeld, Adi Gangidi,	920
	Adolfo Victoria, Ahuva Goldstand, Ajay Menon,	921
	Ajay Sharma, Alex Boesenberg, Alex Vaughan,	922
	Alexei Baevski, Allie Feinstein, Amanda Kallet,	923
	Amit Sangani, Anam Yunus, Andrei Lupu, Andres	924
	Alvarado, Andrew Caples, Andrew Gu, Andrew Ho,	925
	Andrew Poulton, Andrew Ryan, Ankit Ramchandani,	926
	Annie Franco, Aparajita Saraf, Arkabandhu Chowd-	927
	hury, Ashley Gabriel, Ashwin Barambe, Assaf	928
	Eisenman, Azadeh Yazdan, Beau James, Ben Mau-	929
	rer, Benjamin Leonhardi, Po-Yao (Bernie) Huang,	930
	Beth Loyd, Beto de Paola, Bhargavi Paranjape, Bing	931
	Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti,	932
	Brandon Spence, Brani Stojkovic, Brian Gamido,	933
	Britt Montalvo, Carl Parker, Carly Burton, Catalina	934

935	Mejia, Changhan Wang, Changkyu Kim, Chao Zhou,	Sy Choudhury, Sydney Goldman, Tal Remez, Tamar	999
936	Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tin-	Glaser, Tamara Best, Thilo Kohler, Thomas Robin-	1000
937	dal, Christoph Feichtenhofer, Damon Civin, Dana	son, Tianhe Li, Tianjun Zhang, Tim Matthews, Timo-	1001
938	Beaty, Daniel Kreymer, Shang-Wen Li, Danny Wy-	thy Chou, Tzook Shaked, Varun Vontimitta, Victoria	1002
939	att, David Adkins, David Xu, Davide Testuggine,	Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish	1003
940	Delia David, Devi Parikh, Diana Liskovich, Didem	Kumar, Vishal Mangla, Vlad Ionescu, Vlad Andrei	1004
941	Foss, Dingkan Wang, Duc Le, Dustin Holland, Ed-	Poenaru, Vlad T. Mihailescu, Vladimir Ivanov, Wei	1005
942	ward Dowling, Eissa Jamil, Elaine Montgomery,	Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz,	1006
943	Eleonora Presani, Emily Hahn, Emily Wood, Erik	Will Constable, Xia Tang, Xiaofang Wang, Xiao-	1007
944	Brinkman, Esteban Arcaute, Evan Dunbar, Evan	jian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo	1008
945	Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat	Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li,	1009
946	Ozgenel, Francesco Caggioni, Francisco Guzm'an,	Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam,	1010
947	Frank J. Kanayet, Frank Seide, Gabriela Medina	Yu Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach	1011
948	Florez, Gabriella Schwarz, Gada Badeer, Georgia	Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen,	1012
949	Swee, Gil Halpern, Govind Thattai, Grant Herman,	Zhenyu Yang, and Zhiwei Zhao. 2024. <i>The llama 3</i>	1013
950	Grigory G. Sizov, Guangyi Zhang, Guna Lakshmi-	<i>herd of models</i> . <i>ArXiv</i> , abs/2407.21783.	1014
951	narayanan, Hamid Shojanazeri, Han Zou, Hannah		
952	Wang, Han Zha, Haroun Habeeb, Harrison Rudolph,	Yongqi Fan, Hongli Sun, Kui Xue, Xiaofan Zhang,	1015
953	Helen Suk, Henry Apegren, Hunter Goldman, Igor	Shaoting Zhang, and Tong Ruan. 2024. Medodyssey:	1016
954	Molybog, Igor Tufanov, Irina-Elena Veliche, Itai	A medical domain benchmark for long context	1017
955	Gat, Jake Weissman, James Geboski, James Kohli,	evaluation up to 200k tokens. <i>arXiv preprint</i>	1018
956	Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff	<i>arXiv:2406.15019</i> .	1019
957	Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizen-		
958	stein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi	Yongqi Fan, Nan Wang, Kui Xue, Jingping Liu, and	1020
959	Yang, Joe Cummings, Jon Carvill, Jon Shepard,	Tong Ruan. 2025. Medeureka: A medical domain	1021
960	Jonathan McPhie, Jonathan Torres, Josh Ginsburg,	benchmark for multi-granularity and multi-data-type	1022
961	Junjie Wang, Kaixing(Kai) Wu, U KamHou, Karan	embedding-based retrieval. In <i>Findings of the Asso-</i>	1023
962	Saxena, Karthik Prasad, Kartikay Khandelwal, Katay-	<i>ciation for Computational Linguistics: NAACL 2025</i> ,	1024
963	oun Zand, Kathy Matosich, Kaushik Veeraragha-	pages 2825–2851.	1025
964	van, Kelly Michelena, Keqian Li, Kun Huang, Kun-		
965	al Chawla, Kushal Lakhota, Kyle Huang, Lailin	S. Fleming, Alejandro Lozano, William J. Haberkorn,	1026
966	Chen, Lakshya Garg, A Lavender, Leandro Silva,	Jenelle A. Jindal, Eduardo Pontes Reis, Rahul Thapa,	1027
967	Lee Bell, Lei Zhang, Liangpeng Guo, Licheng	Louis Blankemeier, Julian Z. Genkins, Ethan H.	1028
968	Yu, Liron Moshkovich, Luca Wehrstedt, Madian	Steinberg, Ashwin Nayak, Birju S. Patel, Chia-Chun	1029
969	Khabsa, Manav Avalani, Manish Bhatt, Maria Tsim-	Chiang, Alison Callahan, Zepeng Huo, Sergios Ga-	1030
970	poukelli, Martynas Mankus, Matan Hasson, Matthew	tidis, Scott J. Adams, Oluseyi Fayanju, Shreya Shah,	1031
971	Lennie, Matthias Reso, Maxim Groshev, Maxim	Thomas Savage, Ethan Goh, Akshay S. Chaudhari,	1032
972	Naumov, Maya Lathi, Meghan Keneally, Michael L.	Nima Aghaeepour, Christopher D. Sharp, Michael A.	1033
973	Seltzer, Michal Valko, Michelle Restrepo, Mihir	Pfeffer, Percy Liang, Jonathan H. Chen, Keith E.	1034
974	Patel, Mik Vyatskov, Mikayel Samvelyan, Mike	Morse, Emma Brunskill, Jason Alan Fries, and	1035
975	Clark, Mike Macey, Mike Wang, Miquel Jubert Her-	Nigam H. Shah. 2023. <i>Medalign: A clinician-</i>	1036
976	moso, Mo Metanat, Mohammad Rastegari, Mun-	<i>generated dataset for instruction following with elec-</i>	1037
977	ish Bansal, Nandhini Santhanam, Natascha Parks,	<i>tronic medical records</i> . In <i>AAAI Conference on Arti-</i>	1038
978	Natasha White, Navyata Bawa, Nayan Singhal, Nick	<i>ficial Intelligence</i> .	1039
979	Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev,		
980	Ning Dong, Ning Zhang, Norman Cheng, Oleg	Thibault Formal, Benjamin Piwowarski, and Stéphane	1040
981	Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem	Clinchant. 2021. Splade: Sparse lexical and expan-	1041
982	Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pa-	sion model for first stage ranking. In <i>Proceedings</i>	1042
983	van Balaji, Pe dro Rittner, Philip Bontrager, Pierre	<i>of the 44th International ACM SIGIR Conference on</i>	1043
984	Roux, Piotr Dollár, Polina Zvyagina, Prashant Ratan-	<i>Research and Development in Information Retrieval</i> ,	1044
985	chandani, Pritish Yuvraj, Qian Liang, Rachad Alao,	pages 2288–2292.	1045
986	Rachel Rodriguez, Rafi Ayub, Raghotham Murthy,		
987	Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah	Cheng-Ping Hsieh, Simeng Sun, Samuel Krizan, Shan-	1046
988	Hogan, Robin Battey, Rocky Wang, Rohan Mah-	tanu Acharya, Dima Rekesh, Fei Jia, Yang Zhang,	1047
989	eswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu,	and Boris Ginsburg. 2024. Ruler: What’s the real	1048
990	Samyak Datta, Sara Chugh, Sara Hunt, Sargun	context size of your long-context language models?	1049
991	Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma,	<i>arXiv preprint arXiv:2404.06654</i> .	1050
992	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-		
993	say, Sheng Feng, Shenghao Lin, Shengxin Cindy	Aaron Hurst, Adam Lerer, Adam P Goucher, Adam	1051
994	Zha, Shiva Shankar, Shuqiang Zhang, Sinong Wang,	Perelman, Aditya Ramesh, Aidan Clark, AJ Os-	1052
995	Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala,	trow, Akila Welihinda, Alan Hayes, Alec Radford,	1053
996	Stephanie Max, Stephen Chen, Steve Kehoe, Steve	et al. 2024. Gpt-4o system card. <i>arXiv preprint</i>	1054
997	Satterfield, Sudarshan Govindaprasad, Sumit Gupta,	<i>arXiv:2410.21276</i> .	1055
998	Sung-Bae Cho, Sunny Virk, Suraj Subramanian,		
		Gautier Izacard, Mathilde Caron, Lucas Hosseini, Se-	1056
		bastian Riedel, Piotr Bojanowski, Armand Joulin,	1057

1058	and Edouard Grave. 2021. <a href="#">Unsupervised dense information retrieval with contrastive learning</a> . <i>Trans. Mach. Learn. Res.</i> , 2022.	1112
1059		1113
1060		1114
1061	Minbyul Jeong, Jiwoong Sohn, Mujeen Sung, and Jae-woo Kang. 2024. Improving medical reasoning through retrieval and self-reflection with retrieval-augmented large language models. <i>Bioinformatics</i> , 40(Supplement_1):i119–i129.	1115
1062		1116
1063		1117
1064		1118
1065		1119
1066	Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. <i>Applied Sciences</i> , 11(14):6421.	1120
1067		1121
1068		1122
1069		1123
1070		1124
1071	Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. <i>arXiv preprint arXiv:1909.06146</i> .	1125
1072		1126
1073		1127
1074		1128
1075	Qiao Jin, Won Kim, Qingyu Chen, Donald C Comeau, Lana Yeganova, W John Wilbur, and Zhiyong Lu. 2023. Medcpt: Contrastive pre-trained transformers with large-scale pubmed search logs for zero-shot biomedical information retrieval. <i>Bioinformatics</i> , 39(11):btad651.	1129
1076		1130
1077		1131
1078		1132
1079		1133
1080		1134
1081	Alistair Johnson, Lucia Bulgarelli, Tom Pollard, Steve Horng, Leo A. Celi, and Roger Mark. 2023. <a href="#">MIMIC-IV-Note: deidentified free-text clinical notes (version 2.2)</a> . <i>PhysioNet</i> .	1135
1082		1136
1083		1137
1084		1138
1085	Alistair E. W. Johnson, Tom J. Pollard, Lu Shen, Liwei H. Lehman, Melania Feng, Marzyeh Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. 2016. <a href="#">Mimic-iii, a freely accessible critical care database</a> . <i>Scientific Data</i> , 3:160035.	1139
1086		1140
1087		1141
1088		1142
1089		1143
1090	Praveen K Kanithi, Clément Christophe, Marco AF Pimentel, Tathagata Raha, Nada Saadi, Hamza Javed, Svetlana Maslenskova, Nasir Hayat, Ronnie Rajan, and Shadab Khan. 2024. Medic: Towards a comprehensive framework for evaluating llms in clinical applications. <i>arXiv preprint arXiv:2409.07314</i> .	1144
1091		1145
1092		1146
1093		1147
1094		1148
1095		1149
1096	Omar Khattab and Matei Zaharia. 2020. Colbert: Efficient and effective passage search via contextualized late interaction over bert. In <i>Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval</i> , pages 39–48.	1150
1097		1151
1098		1152
1099		1153
1100		1154
1101		1155
1102	Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024. Prometheus 2: An open source language model specialized in evaluating other language models. <i>arXiv preprint arXiv:2405.01535</i> .	1156
1103		1157
1104		1158
1105		1159
1106		1160
1107		1161
1108	Konstantin Kotschenreuther. 2024. Ehr-ds-qa: A synthetic qa dataset derived from medical discharge summaries for enhanced medical information retrieval systems. <i>PhysioNet</i> .	1162
1109		1163
1110		1164
1111		1165
		1166
	Sunjun Kweon, Jiyoun Kim, Heeyoung Kwak, Dongchul Cha, Hangyul Yoon, Kwanghyun Kim, Jee-won Yang, Seunghyun Won, and Edward Choi. 2024. <a href="#">Ehrnoteqa: An llm benchmark for real-world clinical practice using discharge summaries</a> . In <i>Neural Information Processing Systems</i> .	1112
		1113
		1114
		1115
		1116
		1117
	Eric Lehman, Vladislav Lialin, Katelyn Edelwina Legaspi, Anne Janelle Sy, Patricia Therese Pile, Nicole Rose Alberto, Richard Raymund Ragasa, Corinna Victoria Puyat, Marianne Katharina Taliño, Isabelle Rose Alberto, Pia Gabrielle Alfonso, Dana Moukheiber, Byron Wallace, Anna Rumshisky, Jennifer Liang, Preethi Raghavan, Leo Anthony Celi, and Peter Szolovits. 2022. <a href="#">Learning to ask like a physician</a> . In <i>Proceedings of the 4th Clinical Natural Language Processing Workshop</i> , pages 74–86, Seattle, WA. Association for Computational Linguistics.	1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126
		1127
		1128
	Quinn Leng, Jacob Portes, Sam Havens, Matei Zaharia, and Michael Carbin. 2024. Long context rag performance of large language models. <i>arXiv preprint arXiv:2411.03538</i> .	1129
		1130
		1131
		1132
	Hui Yi Leong, Yifan Gao, and Shuai Ji. 2024. A gen ai framework for medical note generation. In <i>2024 6th international conference on artificial intelligence and computer applications (ICAICA)</i> , pages 423–429. IEEE.	1133
		1134
		1135
		1136
		1137
	Mosh Levy, Alon Jacoby, and Yoav Goldberg. 2024. Same task, more tokens: the impact of input length on the reasoning performance of large language models. <i>arXiv preprint arXiv:2402.14848</i> .	1138
		1139
		1140
		1141
	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. <a href="#">Retrieval-augmented generation for knowledge-intensive nlp tasks</a> .	1142
		1143
		1144
		1145
		1146
		1147
	Jiaqi Li, Mengmeng Wang, Zilong Zheng, and Muhan Zhang. 2024a. Loogle: Can long-context language models understand long contexts? In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 16304–16333.	1148
		1149
		1150
		1151
		1152
		1153
	Xinze Li, Yixin Cao, Yubo Ma, and Aixin Sun. 2024b. Long context vs. rag for llms: An evaluation and revisits. <i>arXiv preprint arXiv:2501.01880</i> .	1154
		1155
		1156
	Valentin Liévin, Christoffer Egeberg Hother, Andreas Geert Motzfeldt, and Ole Winther. 2024. Can large language models reason about medical questions? <i>Patterns</i> , 5(3).	1157
		1158
		1159
		1160
	Jerry Liu. 2022. Llamaindex. <i>CoRR</i> .	1161
		1162
	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts. <i>arXiv preprint arXiv:2307.03172</i> .	1163
		1164
		1165
		1166

1167	Man Luo, Arindam Mitra, Tejas Gokhale, and Chitta Baral. 2022. Improving biomedical information retrieval with neural retrievers. In <i>proceedings of the AAAI conference on artificial intelligence</i> , volume 36, pages 11038–11046.	1224
1168		1225
1169		1226
1170		1227
1171		
1172	Zizhan Ma, Wenxuan Wang, Guo Yu, Yiu-Fai Cheung, Meidan Ding, Jie Liu, Wenting Chen, and Linlin Shen. 2025. Beyond the leaderboard: Rethinking medical benchmarks for large language models. <i>arXiv preprint arXiv:2508.04325</i> .	1228
1173		1229
1174		1230
1175		1231
1176		
1177	Ali Modarressi, Hanieh Deilamsalehy, Franck Dernoncourt, Trung Bui, Ryan A Rossi, Seunghyun Yoon, and Hinrich Schütze. 2025. Nolima: Long-context evaluation beyond literal matching. <i>arXiv preprint arXiv:2502.05167</i> .	1232
1178		1233
1179		1234
1180		1235
1181		1236
1182		1237
1183	Skatje Myers, Timothy A Miller, Yanjun Gao, Matthew M Churpek, Anoop Mayampurath, Dmitriy Dligach, and Majid Afshar. 2025. Lessons learned on information retrieval in electronic health records: a comparison of embedding models and pooling strategies. <i>Journal of the American Medical Informatics Association</i> , 32(2):357–364.	1238
1184		1239
1185		1240
1186		1241
1187		
1188		
1189	Anastasios Nentidis, Georgios Katsimpras, Anastasia Krithara, Martin Krallinger, Miguel Rodríguez-Ortega, Eduard Rodríguez-López, Natalia Loukachevitch, Andrey Sakhovskiy, Elena Tutubalina, Dimitris Dimitriadis, et al. 2025. Overview of bioasq 2025: The thirteenth bioasq challenge on large-scale biomedical semantic indexing and question answering. In <i>International Conference of the Cross-Language Evaluation Forum for European Languages</i> , pages 173–198. Springer.	1242
1190		1243
1191		1244
1192		1245
1193		1246
1194		1247
1195		
1196		
1197		
1198		
1199	Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023. Capabilities of gpt-4 on medical challenge problems. <i>arXiv preprint arXiv:2303.13375</i> .	1248
1200		1249
1201		1250
1202		1251
1203		1252
1204		1253
1205	Ankit Pal and Malaikannan Sankarasubbu. 2024. <b>Gemini goes to med school: Exploring the capabilities of multimodal large language models on medical challenge problems &amp; hallucinations</b> . In <i>Clinical Natural Language Processing Workshop</i> .	1254
1206		1255
1207		1256
1208		1257
1209		1258
1210		1259
1211		1260
1212		1261
1213		1262
1214		
1215		
1216		
1217		
1218	Anusri Pampari, Preethi Raghavan, Jennifer J. Liang, and Jian Peng. 2018. <b>emrqa: A large corpus for question answering on electronic medical records</b> . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	1263
1219		1264
1220		1265
1221		1266
1222		1267
1223		
1224		
1225		
1226		
1227		
1228		
1229		
1230		
1231		
1232		
1233		
1234		
1235		
1236		
1237		
1238		
1239		
1240		
1241		
1242		
1243		
1244		
1245		
1246		
1247		
1248		
1249		
1250		
1251		
1252		
1253		
1254		
1255		
1256		
1257		
1258		
1259		
1260		
1261		
1262		
1263		
1264		
1265		
1266		
1267		
1268		
1269		
1270		
1271		
1272		
1273		
1274		
1275		
1276		
1277		
1278		

1279	Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024b. Roformer: Enhanced transformer with rotary position embedding. <i>Neurocomputing</i> , 568:127063.	1333
1280		1334
1281		1335
1282		1336
1283	Simon Šuster and Walter Daelemans. 2018. Clicr: a dataset of clinical case reports for machine reading comprehension. <i>arXiv preprint arXiv:1803.09720</i> .	1337
1284		1338
1285		1339
1286	Xiangru Tang, Anni Zou, Zhuosheng Zhang, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark B. Gerstein. 2023. <a href="#">Medagents: Large language models as collaborators for zero-shot medical reasoning</a> . <i>ArXiv</i> , abs/2311.10537.	1340
1287		1341
1288		1342
1289		1343
1290		1344
1291	Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024a. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	1345
1292		1346
1293		1347
1294		1348
1295		1349
1296		1350
1297	Qwen Team et al. 2024b. Qwen2 technical report. <i>arXiv preprint arXiv:2407.10671</i> , 2:3.	1351
1298		1352
1299		1353
1300	Shubham Vatsal and Ayush Singh. 2024. <a href="#">Can gpt re-define medical understanding? evaluating gpt on biomedical machine reading comprehension</a> . <i>ArXiv</i> , abs/2405.18682.	1354
1301		1355
1302		1356
1303	Yubo Wang, Xueguang Ma, and Wenhui Chen. 2024. Augmenting black-box llms with medical textbooks for biomedical question answering. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 1754–1770.	1357
1304		1358
1305		1359
1306		1360
1307		1361
1308	Michael Wornow, Yizhe Xu, Rahul Thapa, Birju Patel, Ethan Steinberg, Scott Fleming, Michael A Pfeffer, Jason Fries, and Nigam H Shah. 2023. The shaky foundations of large language models and foundation models for electronic health records. <i>npj digital medicine</i> , 6(1):135.	1362
1309		1363
1310		1364
1311		1365
1312		1366
1313		1367
1314	Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2023. Efficient streaming language models with attention sinks. <i>arXiv preprint arXiv:2309.17453</i> .	1368
1315		1369
1316		1370
1317		1371
1318	Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. 2024. Benchmarking retrieval-augmented generation for medicine. In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pages 6233–6251.	1372
1319		1373
1320		1374
1321		1375
1322		1376
1323	Peng Xu, Wei Ping, Xianchao Wu, Chejian Xu, Zihan Liu, Mohammad Shoeybi, and Bryan Catanzaro. 2024a. Chatqa 2: Bridging the gap to proprietary llms in long context and rag capabilities. <i>arXiv preprint arXiv:2407.14482</i> .	1377
1324		1378
1325		1379
1326		1380
1327		1381
1328	Ran Xu, Wenqi Shi, Yue Yu, Yuchen Zhuang, Yanqiao Zhu, May D Wang, Joyce C Ho, Chao Zhang, and Carl Yang. 2024b. Bmretriever: Tuning large language models as better biomedical text retrievers. <i>arXiv preprint arXiv:2404.18443</i> .	1382
1329		1383
1330		1384
1331		1385
1332		1386
	Qwen An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxin Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yi-Chao Zhang, Yunyang Wan, Yuqi Liu, Zeyu Cui, Zhenru Zhang, Zihan Qiu, Shanghaoran Quan, and Zekun Wang. 2024. <a href="#">Qwen2.5 technical report</a> . <i>ArXiv</i> , abs/2412.15115.	1387
		1388
	Xi Yang, Aokun Chen, Nima M. Pournejatian, Hoo-Chang Shin, Kaleb E. Smith, Christopher Parisien, Colin B. Compas, Cheryl Martin, Anthony B Costa, Mona G. Flores, Ying Zhang, Tanja Magoc, Christopher A. Harle, Gloria P. Lipori, Duane A. Mitchell, William R. Hogan, Elizabeth A. Shenkman, Jiang Bian, and Yonghui Wu. 2022. <a href="#">A large language model for electronic health records</a> . <i>NPJ Digital Medicine</i> , 5.	1389
		1390
	Tan Yu, Anbang Xu, and Rama Akkiraju. 2024. In defense of rag in the era of long-context language models, 2024. <a href="https://arxiv.org/abs/2409.01666">URL https://arxiv.org/abs/2409.01666</a> .	1391
		1392
	W. Yu, Zhihan Zhang, Zhenwen Liang, Meng Jiang, and Ashish Sabharwal. 2023. <a href="#">Improving language models via plug-and-play retrieval feedback</a> . <i>ArXiv</i> , abs/2305.14002.	1393
		1394
	Xiang Yue, Xinliang Frederick Zhang, Ziyu Yao, Simon M. Lin, and Huan Sun. 2020. <a href="#">Cliniqg4qa: Generating diverse questions for domain adaptation of clinical question answering</a> . <i>2021 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)</i> , pages 580–587.	1395
		1396
	Lei Zhang, Yunshui Li, Ziqiang Liu, Jiayi Yang, Junhao Liu, Longze Chen, Run Luo, and Min Yang. 2024a. Marathon: A race through the realm of long context with large language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5201–5217.	1397
		1398
	Tianyi Zhang, Varsha Kishore, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. <a href="#">Bertscore: Evaluating text generation with bert</a> . In <i>International Conference on Learning Representations</i> .	1399
		1400
	Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Hao, Xu Han, Zhen Thai, Shuo Wang, Zhiyuan Liu, et al. 2024b. bench: Extending long context evaluation beyond 100k tokens. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15262–15277.	1401
		1402
	Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, et al. 2025. <a href="#">Qwen3</a>	1403

embedding: Advancing text embedding and reranking through foundation models. *arXiv preprint arXiv:2506.05176*.

Xuejiao Zhao, Siyan Liu, Su-Yin Yang, and Chunyan Miao. 2025. Medrag: Enhancing retrieval-augmented generation with knowledge graph-elicited reasoning for healthcare copilot. In *Proceedings of the ACM on Web Conference 2025*, pages 4442–4457.

Ming Zhu, Aman Ahuja, Da-Cheng Juan, Wei Wei, and Chandan K Reddy. 2020. Question answering with long multiple-span answers. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3840–3849.

## A Inference Setup

- **Inference Engine:** Nvidia TensorRT.
- **Quantization:** FP8.
- **Hardware:** 7B parameter models: Deployed on single Nvidia H100 GPUs with 40GB VRAM. 32B parameter models with extended context: Deployed on dual Nvidia H100 GPUs, each equipped with 80GB VRAM. A total of 260 GPU hours was spent.
- **Libraries:** Qdrant, DSPy, Litellm, pandas, Evaluate, Fastembed, PyLate

## B Hyperparameters

Hyperparameter	Value	
LLM	temperature	inst: 0, reas: 1
	freq_penalty	0
	pres_penalty	0
	think tokens	QwQ: 20k HuatuogPT: 8k
RAG	top_k	2 × num chunks
	k <sub>RRF</sub>	60
	k <sub>1</sub>	TF saturation
	b	length norm.
	avgdl	avg. doc length

Table 6: Hyperparameter Configuration

## C Datasets

### C.1 Context Data Distributions

Token-length distributions per dataset are shown in Figures 6 and 7.

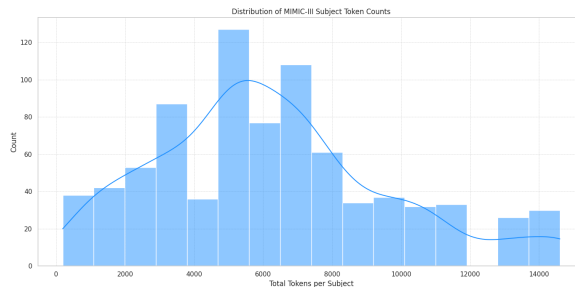


Figure 6: MIMIC-III Context Distributions

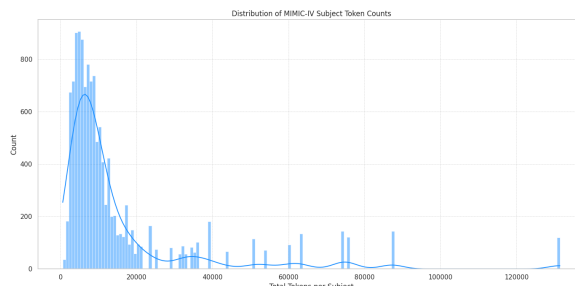


Figure 7: MIMIC-IV Context Distributions

### C.2 Context Data Preparation

We performed minimal cleaning and standardization to support retrieval and prompting, including merging notes by patient and stay, normalizing timestamps, adding note-type metadata, computing token counts for context segmentation, and filtering to human-verified subsets where applicable.

We describe the cleaning and normalization steps applied prior to indexing and prompting:

- Cleaning and de-duplication of clinical notes; removal of template boilerplate where applicable.
- Merging notes by patient and hospital stay; preserving note-type and section headers for downstream retrieval.
- Datetime standardization and chronological ordering; normalization where timezones or partial timestamps occur.
- Metadata augmentation (e.g., note type, encounter identifiers) to support Include DS vs Include All settings.
- Tokenization and token-count computation at note- and patient-level for context segmentation.

These details enable reproducibility of sampling and retrieval segmentation.

Dataset	LLM-as-a-judge	QA Pairs
EHR-DS-QA	Selene-8B	113
	Prometheus-8x7B-v2.0	98
	Qwen2.5-32B-Instruct	71
EHRNoteQA	Selene-8B	180
	Prometheus-8x7B-v2.0	220
	Qwen2.5-32B-Instruct	42

Table 7: Disagreement instances where NLI Med Contradiction is high ( $> 0.7$ ) and LLM Correctness is also high ( $> 0.7$ ).

Dataset	LLM-as-a-judge 1	LLM-as-a-judge 2	QA Pairs
EHR-DS-QA	Prometheus	Qwen2.5-32B-Instruct	520
	Prometheus	Selene-8B	363
	Qwen2.5-32B-Instruct	Selene-8B	664
EHRNoteQA	Prometheus	Qwen2.5-32B-Instruct	690
	Prometheus	Selene-8B	236
	Qwen2.5-32B-Instruct	Selene-8B	519

Table 8: Judge-pair disagreements where the absolute difference in LLM Correctness  $\geq 0.5$ .

Metric	Splade	BM25
LLM Correctness	67.30	<b>69.34</b>
LLM Completeness	65.41	<b>67.14</b>
LLM Faithfulness	<b>47.64</b>	46.23
NLI Med Entailment	52.68	<b>54.35</b>
NLI Med Contradiction	16.98	<b>16.94</b>
Bio BERTScore F1	<b>79.75</b>	79.67
METEOR	<b>43.93</b>	43.63

Table 9: Comparison of Different Sparse Retrievers

### C.3 Prompt Formulations

We used a simple one shot prompt structure and tailored it explicitly for each task formulation:

- **Extractive:** Request the model to answer by extracting the most relevant answer from the context.
- **Multiple-choice:** Standard multiple choice prompt.
- **Open-ended:** Focusing on open ended question answering favoring short, single sentence answers.

Below is the sample prompt for extractive tasks:

```
System message:
Your input fields are:
1. "medical_record" (list[str]): List of patient notes (chronological order).
2. "question" (str): A question about the patient's record.
Your output fields are:
1. "answer" (str): Short single-sentence answer to the question.

All interactions will be structured in the following way, with the appropriate values filled in.

[[ ## medical_record ## ]]
[medical_record]

[[ ## question ## ]]
[question]

[[ ## answer ## ]]
[answer]

[[ ## completed ## ]]

In adhering to this structure, your objective is:
Given a patient's medical record and a question, answer the question correctly and shortly.
```

## C.4 Dataset Comparisons

Aiming to stress-test LLM capabilities under real patient-centric scenarios, our dataset selection process was based on an extensive grid of datasets. Table 10 summarizes the QA datasets we analyzed in our selection process, including the EHR candidates. The chosen EHR-based, human-verified datasets provide diverse but comparable settings across generative, extractive, and MC formulations. While some exceed 8K tokens, supporting our long-context evaluation, they all provide strong augmentation potential towards a multi-note longer evaluation setting that allows content extension up to 128K for each patient (§4.2), enabling a realistic comparison of full-context and RAG pipelines. We provide citations of each considered dataset from Table 10 in Table 11.

## D Model Comparisons

### D.1 Retrievers

On Table 9 we provide a sparse retriever comparison on the EHR-DS-QA dataset using the Qwen2.5-7B-Instruct-1M LLM and Qwen3-Embedding-8B as the dense retriever. BM25 has better performance across LLM Correctness, LLM Completeness and NLI-based metrics while showcasing slightly inferior performance on BioClinical BERTScore F1, METEOR and LLM Faithfulness.

### D.2 LLM-as-a-judge

In Table 7 we examine the number of QA cases for each dataset that the different Judges cause disagreement with the NLI Med Contradiction metric. We observe that Qwen2.5-32B-Instruct<sup>6</sup> has far less such disagreement compared to more specialized judges. Table 8 also provides direct comparisons of LLM-as-a-judge model, demonstrating cases where LLM Correctness between such pairs is more than 0.5 and therefore high. Selene-

<sup>6</sup><https://huggingface.co/Qwen/Qwen2.5-32B-Instruct>

Dataset	# QA Pairs	Token Len. (max)	Task Types	Synthetic	Annotation	Reasoning
Medical Textbooks, Websites, Knowledge Bases						
MedQuAD	47,457	< 2K	extractive		automatic	✗
MashQA	34,808	< 2K	extractive (multi-span)		experts	✓
MedOdyssey (En.KG)	100	< 128K	extrative (graph)		model&human	✓
MedQA (USMLE)	12,723	> 1M	MC		experts	✓
Biomedical Literature						
BioMRC (LARGE)	812,707	< 2K	MC, extractive		automatic	✗
BioASQ (b)	5,729	< 8K	extractive, generative		experts	✓
PubMedQA	1,000	< 2K	MC, generative		experts	✓
Clinical Case Reports						
CliCR	104,919	< 8K	extractive (cloze)		automatic	✓
Patient Notes						
LongHealth	400	< 8K	MC	✓	experts	✓
DiSCQ	2,029	< 2K	extractive, generative		experts	✓
RadQA	6,148	< 16K	extractive		experts	✓
CliniQG4QA	8,824/1,287	< 8K	extractive		model/experts	✗
EHR-DS-QA	156,599/478	< 8K	generative	✓	model/experts	✗
EHRNoteQA	962	< 8K	MC, generative		experts	✓

Table 10: Biomedical Datasets Comparison. References for each dataset are in Table 11.

Dataset	Source
MedQA	Jin et al. (2021)
MedQuAD	Ben Abacha and Demner-Fushman (2019)
MashQA	Zhu et al. (2020)
MedOdyssey	Fan et al. (2024)
BioMRC	Pappas et al. (2020)
BioASQ	Nentidis et al. (2025)
PubMedQA	Jin et al. (2019)
CliCR	Šuster and Daelemans (2018)
LongHealth	Adams et al. (2024)
MediNote	Leong et al. (2024)
DiSCQ	Lehman et al. (2022)
RadQA	Soni et al. (2022)
CliniQG4QA	Yue et al. (2020)
EHR-DS-QA	Kotschenreuther (2024)
EHRNoteQA	Kweon et al. (2024)

Table 11: Corresponding Citations for Datasets considered in Dataset Selection of Table 10

8B <sup>7</sup> and Prometheus-8x7B-v2.0 <sup>8</sup> comparisons show that these two models are consistently more in agreement between them while Qwen2.5-32B-Instruct shows more independent judging abilities.

<sup>7</sup><https://huggingface.co/AtlaAI/Selene-1-Mini-Llama-3.1-8B>

<sup>8</sup><https://huggingface.co/prometheus-eval/prometheus-8x7b-v2.0>

### D.3 Context Size Performance

Figure 8 presents performance on open-ended QA across metrics and models.

## E Metrics

### E.1 Correlations

We examine correlations between evaluation metrics to identify metric independence. On EHRNoteQA and EHR-DS-QA across all non-exclude settings (Figure 9) and long-context high LLM Correctness non-exclude settings (Figure 10), LLM evaluation metrics exhibited strong inter-correlations and METEOR and BERTScore were highly correlated. Although the above correlations remain high, they are less pronounced for the long-context high LLM Correctness setting of Figure 10. In general all metrics are less correlated for the longer context setting.

### E.2 Analysis

Table 12 showcases a qualitative analysis on what metrics are capturing vs missing based on analyzing disagreement data.

## F Case Studies

Here we show some sample studies of our error analysis.

### EHR-DS-QA: Case 10083814 (final diagnoses).

Three discharge summaries exist across distinct admissions; each lists diagnoses. The correct target

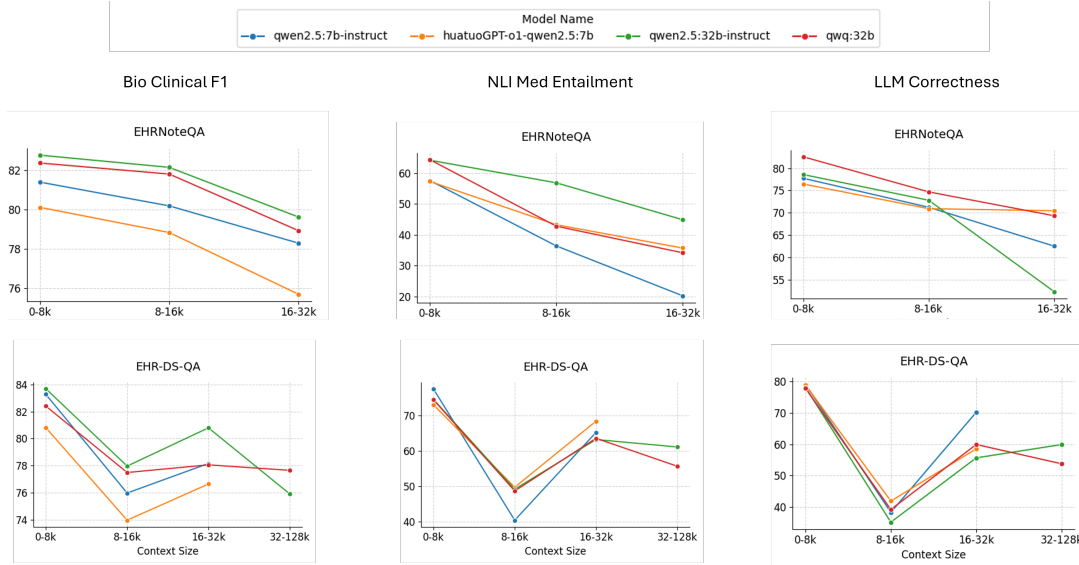


Figure 8: Model performance over context size.

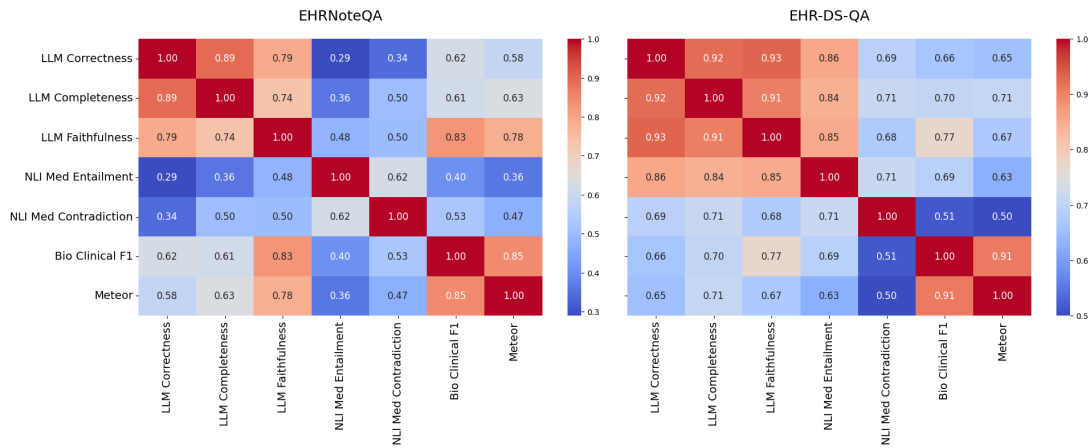


Figure 9: Metric correlations across all non-exclude settings.

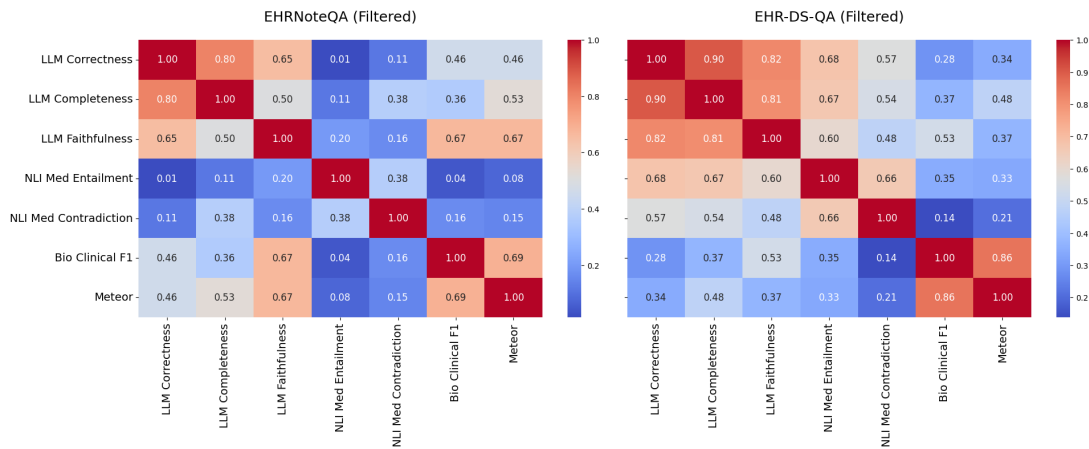


Figure 10: Metric correlations across all non-exclude settings filtered by LLM Correctness > 50 and Context >= 16K.

What they capture well	What they miss	Typical pattern	Examples (dataset #subject)
<b>LLM Correctness</b>			
Whether the main claim is right.	Unsupported add-ons that don't change the core fact.	Correctness high; Faithfulness low (adds extra).	EHRNoteQA #10043423 (core claim right, adds recommendation); EHRDSQA #10094318 (chief complaint right, extra symptom).
<b>LLM Completeness</b>			
Coverage of the requested pieces/elements.	Penalizes omissions but not ungrounded embellishments.	Completeness high; Faithfulness low (fully covers, then embellishes).	EHRNoteQA #10043423 (covers requested items + extras); EHRDSQA #19796003 (lists largely complete; minor extras elsewhere reduce faithfulness/F1).
<b>LLM Faithfulness</b>			
Grounding to evidence/gold; flags hallucinated or embellished content.	Can be over-strict on benign context or plausible but unsupported expansions.	Correct/Complete high; Faithfulness low (ungrounded detail).	EHRNoteQA #10043423 (extra recommendation not supported); EHRDSQA #10131388 (adds follow-up/psychiatric context not explicit).
<b>NLI Med Entailment / Contradiction</b>			
Sentence-level logical relation to gold (entailed vs. contradicted).	Sensitive to negation, hedging ("low risk" vs "no risk"), and phrasing templates.	Correctness high; Entailment low or Contradiction high (wording/negation mismatch).	EHRNoteQA #10494486 (mechanism phrasing hurts entailment); EHRDSQA #10010655 ("no risk" vs "low acute risk" triggers contradiction).
<b>Meteor</b>			
Surface-form overlap, good for close paraphrases.	Under-rewards long phrases, synonyms, and re-phrasings with low lexical overlap. Incorrectly rewards inaccuracies with lexical overlap.	Correctness high; Meteor low (semantic match, low surface match).	EHRDSQA #10076958 (concise "diffuse ischemic bowel" vs long gold narrative); #19796003 correct single word answer but low score; EHRNoteQA #10404814 (correct causal link, different wording/structure).
<b>BioClinical F1</b>			
Token/span overlap for clinical entities; good checklist signal for missing/extra items.	Under-rewards clinically acceptable reformulations (class vs. specific drug). Incorrectly rewards inaccuracies with lexical overlap.	Correctness high; F1 moderate/low (one entity missing or formatted differently).	EHRNoteQA #10043423 (right ideas, entity list/format differences lower F1); EHRDSQA #19796003 (near-exact meds but lower F1 for small misses).

Table 12: What each metric captures vs. misses, with typical disagreement patterns and representative examples from EHRNoteQA and EHRDSQA.

for "final diagnoses" is the most recent summary. Earlier summaries contain different diagnoses that are no longer *final* at discharge, explaining prediction-gold divergence without model hallucination.

**EHRNoteQA: Case 15877599 (cause of AKI).** The note supports a causal chain: gastroenteritis → increased ostomy output → severe dehydration (prerenal) → acute kidney injury. The gold captures the underlying illness; a model response may describe the immediate mechanism. Both are clinically coherent parts of the same sequence.

**EHR-DS-QA: Case 10023117 (blood pressure at discharge).** Predictions expressed as exact values versus ranges produce different behaviors across metrics. Ranges can be faithful to notes yet fail entailment or completeness thresholds defined against a single gold value.

## G RAG vs FC: Detailed Analysis

We summarize our insights about which RAG or FC setting tends to be advantageous across datasets and query types in Table 13. We also provide examples along each of our main findings below:

*FC performs better for questions requiring comprehensive understanding* of the entire document or when answers are located in structured sections, such as summaries, lists, or temporal sequences. Examples include:

- Questions like "What is the patient's discharge condition?" or "What were the patient's discharge diagnoses?"
- Temporal sequence questions, such as "What surgeries has the patient undergone and in what order?"

*RAG outperforms FC for questions that require retrieving specific details* or synthesizing information scattered across multiple parts of the document.

Examples include:

- Questions like "What family history does the patient have?" or "How many tablets of dilaudid did the patient receive?"
- Synthesis tasks, such as "What were the patient's postoperative course details?"
- Asking about specific dates like "What was the outcome of the patient's colonoscopy as described in the discharge summary from the stay starting on 2113-09-30?" or "What was the patient's diagnosis for the hospital admission on 2154-01-28..." and other examples.

*RAG tends to perform better overall in datasets with complex or lengthy documents.*

For example:

- In the EHRNoteQA dataset, RAG consistently outperformed FC for questions needing specific details from notes or summaries, such as "What was the outcome of the patient's colonoscopy?"

1581 *There is mixed evidence suggesting FC might per-*  
1582 *form better for tasks involving inferential reasoning*  
1583 *or identifying the absence of information.*

1584 For example:

- 1585 • Questions like “Were there any complications  
1586 during the procedure?” where RAG retrieves  
1587 statements like “No complications,” poten-  
1588 tially diminishing FC’s advantage.
- 1589 • Subtle inference tasks, such as “Does the pa-  
1590 tient have any psychological issues?” where  
1591 FC occasionally performs better, though in-  
1592 consistently.

Insight	Favored	Explanation	Supporting Examples
Specific Fact Retrieval	RAG	RAG excels at extracting precise, well-defined medical facts (dates, medications, lab values, procedures) that are typically documented in structured sections of medical records.	<p><b>EHRNoteQA:</b>  <i>Subject 15036658:</i> Colonoscopy outcome from specific date (RAG: 0.632 vs. FC: 0.399)  <i>Subject 11049732:</i> Medication changes (RAG: 0.829 vs. FC: 0.352)  <i>Subject 17818938:</i> Surgical procedure for erectile dysfunction (RAG: 0.918 vs. FC: 0.479)</p> <p><b>EHR-DS-QA:</b>  <i>Subject 10131388:</i> Dilaudid tablet count (RAG: 0.982 vs. FC: 0.778)  <i>Subject 19926045:</i> DVT medication (RAG: 0.901 vs. FC: 0.564)  <i>Subject 10090787:</i> Discharge medications (RAG: 0.425 vs. FC: 0.190)</p>
Temporal Information Processing	Mixed	RAG excels at explicit temporal facts (specific dates, temporal relationships) while FC is better at temporal reasoning (sequencing, duration calculation, recognizing absence of temporal information).	<p><b>RAG Advantage:</b>  <i>EHRNoteQA 15036658:</i> Specific date anchoring  <i>EHRNoteQA 18467824:</i> Temporal relationship between admissions  <i>EHR-DS-QA 10264949:</i> Nausea/vomiting timing</p> <p><b>FC Advantage:</b>  <i>EHRNoteQA 11552479:</i> Temporal sequencing (FC: 0.696 vs. RAG: 0.236)  <i>EHR-DS-QA 10751849:</i> Duration calculation (FC: 0.541 vs. RAG: 0.320)  <i>EHR-DS-QA 19397212:</i> Absence of temporal information (FC: 0.426 vs. RAG: 0.180)</p>
Medical Terminology and Technical Content	RAG	RAG performs better with specialized medical terminology, complex procedures, and technical test results due to its ability to locate and interpret specific sections containing this information.	<p><b>EHRNoteQA:</b>  <i>Subject 17445067:</i> Diagnosis and surgical procedure details (RAG: 0.708 vs. FC: 0.328)  <i>Subject 18122852:</i> MRI and EMG test findings (RAG: 0.750 vs. FC: 0.454)  <i>Subject 16313269:</i> Brain mass pathological diagnosis (RAG: 0.915 vs. FC: 0.372)</p> <p><b>EHR-DS-QA:</b>  <i>Subject 10044189:</i> Necrotic ulcer treatment (RAG: 0.762 vs. FC: 0.465)  <i>Subject 19401508:</i> Treatment for hyponatremia (RAG: 0.479 vs. FC: 0.262)</p>
Discharge Planning and Instructions	RAG	RAG performs better on Discharge information (instructions, medications, condition) that is usually well-structured in specific sections.	<p><b>EHRNoteQA:</b>  <i>Subject 11690633:</i> Discharge condition and instructions (RAG: 0.762 vs. FC: 0.349)  <i>Subject 11863782:</i> Discharge disposition and medications (RAG: 0.749 vs. FC: 0.327)</p> <p><b>EHR-DS-QA:</b>  <i>Subject 10921250:</i> Discharge condition (RAG: 0.856 vs. FC: 0.466)  <i>Subject 10940920:</i> Discharge instructions (RAG: 0.650 vs. FC: 0.253)  <i>Subject 10064678:</i> Discharge instructions (RAG: 0.352 vs. FC: 0.094)</p>
Cause–Effect and Relationship Understanding	RAG	RAG is better at understanding relationships (symptom–procedure, medication–outcome, test result–action).	<p><b>EHRNoteQA:</b>  <i>Subject 13032648:</i> Causes of leg pain and surgical procedure (RAG: 0.730 vs. FC: 0.293)  <i>Subject 15748482:</i> Flomax usage reason and outcome (RAG: 0.819 vs. FC: 0.495)  <i>Subject 17436366:</i> Blood/urine culture results and actions (RAG: 0.718 vs. FC: 0.445)</p> <p><b>EHR-DS-QA:</b>  <i>Subject 19926045:</i> Lovenox symptom management (RAG: 0.545 vs. FC: 0.334)  <i>Subject 19401508:</i> Admission cause and treatment (RAG: 0.372 vs. FC: 0.108)</p>
Holistic Patient Understanding	FC	FC excels when questions require synthesis of information across multiple document sections to build a complete picture of the patient’s status, multiple diagnoses.	<p><b>EHRNoteQA:</b>  <i>Subject 18753609:</i> Therapeutic interventions for leg pain (FC: 0.675 vs. RAG: 0.450)</p> <p><b>EHR-DS-QA:</b>  <i>Subject 10262565:</i> Discharge condition (FC: 0.882 vs. RAG: 0.681)  <i>Subject 10049941:</i> Discharge diagnoses (FC: 0.579 vs. RAG: 0.188)  <i>Subject 10978236:</i> Discharge condition — unusual circumstance (FC: 0.321 vs. RAG: 0.093)</p>
Absence or Negative Information Recognition	FC	FC is better at recognizing when information is absent or when negative findings are documented, as it can assess the entire document context.	<p><b>EHR-DS-QA:</b>  <i>Subject 19397212:</i> Absence of symptom presentation timing (FC: 0.426 vs. RAG: 0.180)  <i>Subject 10751849:</i> Absence of major procedures (FC: 0.243 vs. RAG: 0.037)  <i>Subject 10264949:</i> Absence of social factors (FC: 0.497 vs. RAG: 0.149)  <i>Subject 19397212:</i> Absence of age information (FC: 0.725 vs. RAG: 0.452)</p>

Table 13: RAG vs FC Detailed Analysis.