

# Towards semantic reliable clinical QA: Query pipeline optimization for cancer patient question answering systems

MaoLin He<sup>♡</sup>, Rena Gao<sup>♡,♣</sup>, Mike Conway<sup>♡</sup>, Brian E. Chapman<sup>♠</sup>

<sup>♡</sup>School of Computing and Information Systems, University of Melbourne

<sup>♠</sup>Health Data Science and Biostatistics, University of Texas Southwestern Medical Center

<sup>♣</sup> Keeta AI

mlhinaus@gmail.com, wegao@student.unimelb.edu.au

mike.conway@unimelb.edu.au, Brian.Chapman@utsouthwestern.edu

## Abstract

Large Language Models (LLMs) show promise in medical Question-Answering (QA) but suffer from hallucinations that jeopardize patient safety. While Retrieval-Augmented Generation (RAG) mitigates this by grounding outputs in external evidence, existing pipelines struggle with the complex, rapidly evolving nature of oncology. We present **CoMeta**, a three-level controllable metadata-aware framework optimized for Cancer Patient QA (CPQA). We introduce Clinical Hybrid Semantic-Symbolic Document Retrieval (CHSDR), which synergizes real-time Boolean search via NCBI E-Utilities with semantic retrieval to overcome metadata blindness. Additionally, we propose **Semantic Enhanced Overlap Segmentation (SEOS)** to prevent contextual fragmentation. Our results demonstrate that CHSDR significantly improves retrieval performance, CoMeta improved the answer accuracy of Claude-3-haiku by 5.24% over chain-of-thought prompting and about 3% over a naive RAG setup. This study highlights the importance of domain-specific query optimization in realizing the full potential of RAG and provides a robust framework for building more reliable CPQA systems.

## 1 Introduction

Large Language Models (LLMs) have shifted the paradigm of online information seeking from traditional search engines to conversational agents. While promising for Question-Answering (QA), their reliance on internal parameters renders them prone to hallucination—generating fluent but factually incorrect outputs (Ji et al., 2023). This issue is acute in the **clinical domain**, where accuracy directly impacts patient safety (Umapathi et al., 2023). To mitigate this, Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) integrates external evidence to ground LLM responses. Especially in medical QA tasks where queries are knowledge intensive, LLMs excel as generators rather than

knowledge databases (Truhn et al., 2023).

To enhance medical QA, prior studies (Jeong et al., 2024; Liang et al., 2025) broadly use RAG with domain-specific adaptation of retrieval strategies. However, they overlook the query pipeline, the foundational mechanism for query-to-evidence mapping that underpins all retrieval strategies. In Cancer Patient QA (CPQA), these pipelines cause systematic retrieval failures. First, retrieval reliability face the **staleness-semantic dilemma**: the rapid evolution of oncology (Landhuis, 2016) means that standard query pipelines (e.g., Dense or BM25), built upon static and metadata-blind indexes, risk surfacing outdated and thus clinically unreliable evidence, while real-time, metadata-aware interfaces like E-Utilities (Kans, 2024) are semantically brittle to informal patient queries. Second, a **retrieval-depth paradox** stems from the asymmetric utility of publication types. We argue that reviews require full-text retrieval to capture high-level therapeutic syntheses (Wang et al., 2025), while primary research articles often benefit from abstract-only retrieval to avoid methodological noise. Most pipelines apply uniform retrieval depth across article types, either starving the model of synthesis or overwhelming it with clutter. Third, semantic representation is undermined by **contextual fragmentation**: prior encoder-agnostic segmentation (pre-defined length (Liu, 2022) or lexical (Hearst, 1997)) severs clinical qualifiers (e.g., specific mutation criteria) from therapeutic statements, yielding recommendations that appear evidence-based yet lack essential medical constraints.

We address these failures with **CoMeta**, a specialized RAG framework for CPQA that enforces controllability across three pillars: 1) Robustness against the staleness-semantic dilemma, ensuring stable performance across expert queries and informal patient narratives; 2) Evidence Metadata-awareness, enabling time filtering and adaptive retrieval depth based on publication types; 3) Rela-

tional Integrity, utilizing encoder-aware segmentation to safeguard clinical logic. To support this architecture, we curate the Cancer-related Medical QA (CMMQA) dataset and its "Clinical Narrative Variant" to simulate real-world patient interactions<sup>1</sup>. Our specific contributions are:

- We propose the Clinical Hybrid Semantic-Symbolic Document Retrieval (CHSDR) method. To our knowledge, this is the first pipeline in the clinical RAG domain to combine real-time Boolean search (via E-Utilities with LLM-rewritten queries) with MedCPT (Jin et al., 2023) semantic search.
- We introduce Semantic Enhanced Overlap Segmentation (SEOS), a novel text splitter that integrates sentence semantics and embedding model characteristics while utilizing chunk overlap to preserve critical context.
- We provide the first comparative analysis of NCBI sources for CPQA, demonstrating that PMC review articles possess significantly higher retrieval value than non-review PMC papers among clinical cancer QA datasets.

## 2 Related Work

### 2.1 Retrieval in Medical RAG

Medical RAG systems rely on their query pipelines, which map user queries to evidence. The predominant paradigm is dense retrieval, leveraging domain-specific embedding models for vector-similarity search in offline index (Miao et al., 2025). Recent advances have shifted from general-purpose embeddings to domain-specific retrievers like MedCPT (Jin et al., 2023). Advanced retrieval strategies, including hybrid search that integrates sparse methods like BM25 (Robertson et al., 2009) for lexical-semantic complementarity (Xiong et al., 2024; Xu et al., 2026), adaptive retrieval that dynamically selects retrieval timing (Jeong et al., 2024) and recursive search that iteratively refine queries through feedback loops (Liang et al., 2025), are essentially optimizations built upon static-index pipelines.

The reliance on offline indexing restricts metadata filtering (e.g., by publication date or study type), incurs substantial preprocessing and storage overhead (e.g., >400GB for PubMed indices

(Jeong et al., 2024)), and introduces timeliness gaps where newly published studies may remain inaccessible in critical decision windows. These issues are acute in oncology, where clinical standards evolve rapidly. NCBI’s E-Utilities offers a potential remedy through real-time access and metadata filtering, but its utility is hindered by a term-centric matching paradigm. Unlike BM25, E-Utilities lacks mechanisms like length normalization or term weighting, making retrieval highly sensitive to query formulation and ill-suited for natural language inputs without significant adaptation.

Thus, existing systems are implicitly forced to trade off semantic robustness against retrieval controllability. We instead explore a different design paradigm. Rather than further optimizing static-index pipelines, CHSDR mitigates these limitations by integrating E-utilities as a live, metadata-aware sparse backend. This design is orthogonal and complementary to prior RAG optimizations (e.g., MedRAG, Self-BioRAG, Adaptive RAG), enabling seamless integration with prior works.

### 2.2 Semantic Representation in RAG

Dense retrieval efficacy hinges on precise semantic representation, which is governed by two coupled factors: document segmentation and text encoding (gao). While recent medical RAG systems extensively optimize retriever configuration (Xiong et al., 2024; Tang and Yang, 2024), document segmentation, the foundational step determining input granularity, remains critically under-explored.

Current approaches largely rely on heuristic strategies that compromise clinical context. Naive fixed-length chunking (Jeong et al., 2024) frequently truncates sentences mid-thought, severing syntactic dependencies. Even sentence-aware methods like LlamaIndex’s SentenceSplitter (Liu, 2022), while preserving sentence integrity, impose rigid window sizes that fail to adapt to the variable information density of medical narratives. Similarly, TextTiling (Hearst, 1997) infers topic boundaries via lexical overlap, but its surface-level matching struggles with the high synonymy and complex semantic shifts inherent in biomedical literature.

The fundamental limitation of these approaches is twofold. First, they function as “semantic-blind” preprocessing steps, leading contextual fragmentation. Second, they neglect the interaction between chunk size and encoder performance, ignoring evidence that optimal segmentation is highly dependent on the specific embedding model used (gao).

<sup>1</sup>All codes and datasets can be found via: [anonymous.4open.science/r/COMP90005-E51E/README.md](https://anonymous.4open.science/r/COMP90005-E51E/README.md).

### 3 Methods

#### 3.1 Dataset and Evaluation

We applied a MeSH-based filter to HealthSearchQA (Singhal et al., 2023b) and the MIRAGE benchmark (Xiong et al., 2024) that includes PubMedQA (Jin et al., 2019), BioASQ (Tsatsaronis et al., 2015), MedQA (Jin et al., 2021), MedMCQA (Pal et al., 2022), the medical subsets in MMLU (Hendrycks et al., 2020). Specifically, we use all terms and synonyms under the 'neoplasm' MeSH subtree to identify questions strictly pertaining to cancer as a disease entity. While not exhaustive, it ensures a focus on core disease-related queries. We refer to this filtered collection as the Cancer-related Medical QA Dataset (CMMQA) (Figure 1). Then, we use Llama-3-70b to rewrite the CMMQA queries into non-expert clinical narratives (mainly simulating patients), creating a derived dataset we term the "Clinical Narrative Variant". This variant serves as a challenging testbed to evaluate both retrieval stability and QA accuracy across distinct query modalities (standard vs. narrative).

| Datasets         | Original Datasets | Questions count after filter:520 | Answer Type      |
|------------------|-------------------|----------------------------------|------------------|
| CMMQA-Literature | PubMedQA          | 60                               | Yes / No / Maybe |
|                  | BioASQ            | 82                               | Yes / No         |
| CMMQA-Exam       | MMLU-Med          | 23                               | A/B/C/D          |
|                  | MedMCQA           | 93                               | A/B/C/D          |
|                  | MedQA-US          | 86                               | A/B/C/D          |
| Length-Impact-QA | HealthSearchQA    | 176                              | None             |

Figure 1: Description of filtered cancer QA datasets used in this study.

Our *document* retrieval evaluation focuses on PubMedQA and BioASQ because these two datasets provide gold-standard citations (PMIDs) but others lack source annotations. For *passage* retrieval evaluation, we utilized synthetic Query-Answer pairs generated from PubMed Abstracts, PMC Full Texts and medical textbooks (Jin et al., 2021) as a validation set. Retrieval Performance is measured at a cutoff of K=10 using three standard metrics: Hit Rate (Hit@K), representing the percentage of queries with at least one relevant document retrieved; Mean Reciprocal Rank (MRR@K), assessing the ranking quality of the first relevant document; and Recall (Recall@K), quantifying the coverage of gold-standard evidence. Besides, we report Hit0, the count of queries yielding zero results. This metric serves as a critical indicator of

retrieval robustness, highlighting system failures in handling complex or restrictive queries.

#### 3.2 Controllable Query Pipeline Design

Our target is real clinical settings that do not maintain a local PubMed mirror but do require metadata-aware, up-to-date evidence access. To ensure clinical rigor across the retrieval lifecycle, we design a hierarchical query pipeline (Figure 2) that prioritizes controllability at both the document and passage levels.

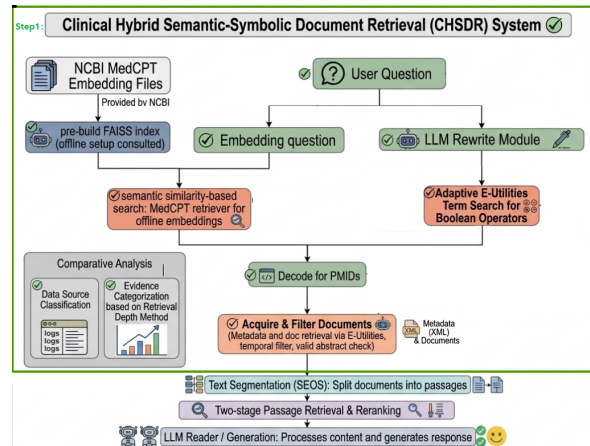


Figure 2: The CHSDR employs dual retrieval strategies, then downloads and filters candidate documents. After document Retrieval, next steps and comparative analyses are conducted

##### 3.2.1 Clinical Hybrid Semantic-Symbolic Document Retrieval (CHSDR)

While BM25 paired with dense retrieval offers strong lexical complementarity, it relies on static indices that are disconnected from real-time updates. In contrast, CHSDR (stage-1 in figure 2) enables metadata control and robust retrieval across various queries (standard vs. clinical narrative), overcoming **staleness-semantic dilemma** and **retrieval-depth paradox**. The system design comparison table are detailed in Appendix A.2.

**Boolean-constrained retrieval with adaptive query execution (Adapt-E)** Raw consumer questions contain orthographic and grammatical errors and vary in length (Lu et al., 2019; Abacha et al., 2019), which are incompatible with E-Utilities. CHSDR thus use a LLM-based rewriter that performs not only query error correction and text normalization (Lu et al., 2019), but also patient intent analysis, clinical abstraction (i.e., mapping narrative descriptions into standardized medical concepts and PICO elements), boolean expression and

temporal constraint generation. To ensure robustness, we design a *adaptive query execution*: generated outputs are executed progressively from strict boolean conjunctions to clinical abstraction of raw queries to relaxed boolean formulations, until sufficient documents are retrieved.

**Hybrid semantic-symbolic retrieval with metadata-aware control** We integrate symbolic search with semantic retrieval via Reciprocal Rank Fusion (RRF) (Cormack et al., 2009), further overcoming the term-centric limitations of E-Utilities. We adopt MedCPT (Jin et al., 2023) for semantic document retrieval, because it is trained for query-article retrieval and achieves a balance between performance and efficiency. Both retrieval streams return PubMed identifiers (PMIDs), which serve as a unified document key. All candidate documents and their associated metadata are subsequently accessed via E-Utilities.

**Metadata Utilization** Retrieved documents are augmented with metadata parsed from E-Utilities XML responses. We exploit three main fields: (1) *Publication Type* (pt), which classifies documents into PubMed abstracts (D1), PMC reviews full-text (D2), or non-review PMC papers (D3), enables empirical analysis of retrieval depth across publication types; (2) *Publication Date* (dp), which can be used to enforce recency filters and prioritizes more recent evidence when conflicts arise; (3) *Abstract Availability*, which excludes records lacking valid abstracts to ensure data quality.

### 3.2.2 Passage Retrieval and Text Splitter

**Semantic Enhanced Overlap Segmentation (SEOS)** To address the twofold limitation in section 2.2, we propose SEOS (pseudo-code in Appendix A.4). SEOS is inspired by TextTiling (Hearst, 1997), which identifies segmentation boundaries via local cohesion shifts, but extends it through following key innovations. First, We replace bag-of-words with domain-specific dense embeddings in text representation, enabling robust handling of medical terminology and discourse relations in clinical studies. Second, instead of brittle similarity thresholds, SEOS determines the optimal partition cardinality (N) from a target token budget, selecting the top-N semantic minima as breakpoints. The budget is adaptively configured to align with the effective input granularity observed for different embedding models, as prior

work has shown retrieval performance to be sensitive to chunk length (gao; Wang et al., 2019). Third, to mitigate contextual fragmentation, SEOS incorporates adaptive sentence overlap determined based on semantic continuity around each breakpoint, preserving unresolved semantic dependencies. Furthermore, neighboring chunk identifiers are explicitly stored, enabling local context recovery when evidence spans multiple segments.

**Dense retrieval and reranking for passage matching** The training objective of MedCPT is optimized for matching queries to PubMed abstract, but our passage retrieval corpus mainly is PMC full text that exhibit different discourse structures and information densities. This distributional shift from abstracts to full-text passages necessitates encoders with broader generalization capabilities. Therefore, we select top-performing embedding models from the MTEB benchmark (Muennighoff et al., 2022). For the subsequent reranking stage, we employ cross-encoder architectures to capture the fine-grained, non-linear semantic interactions that bi-encoders often miss (Jiang et al.). To evaluate the trade-off between domain specialization and general semantic reasoning, we only select representative ones: 1) bge-reranker-v2-m3 (Li et al., 2023), a general-purpose model with superior performance; and 2) MedCPT-reranker (Jin et al., 2023), a domain-specific reranker pre-trained on large-scale biomedical literature.

## 4 Results

We compare our work with two settings: 1) LLM with COT, a retrieval-free approach relying solely on parametric knowledge. Chain-of-Thought (CoT) (Wei et al., 2022) encourages the LLM to perform step-by-step reasoning to improve the quality of the generated answers. 2) LLM with naive RAG, using the whole text of the Top 5 relevant documents retrieved by the MedCPT to directly guide LLM to generate an answer.

### 4.1 Analysis in Document retrieval

#### 4.1.1 Component Ablations of CHSDR

In Table 1, we conduct the ablation study of CHSDR on retrieval performance. Results indicate: 1. Impacts of Clinical Narrative. While shifting to *Clinical Narrative* inputs universally degrades performance, the hybrid model maintains robustness. This suggests that when the input text is noisy,

the semantic understanding provided by the Hybrid system becomes the critical mechanism for retrieval resilience. **2. the Hybrid Dilemma.** The **Hybrid** model achieves best performance on BioASQ, outperforming **Adapt-E** method. This indicates that semantic retrieval (MedCPT) successfully recalls relevant documents that are not retrieved by sparse retriever. However, on PubMedQA, the Hybrid approach slightly underperforms the symbolic-only Adapt-E method in the *Standard Narrative*. We attribute this anomaly to query specificity: PubMedQA questions are often derived directly from abstract sentences, creating a bias toward exact lexical matching. **3. the "Zero-Hit" Barrier.** Across most datasets, the limitations of sparse search are evident in the **E-utils** baseline. In contrast, **Adapt-E** method (LLM Rewriter + Adaptive Fallback) improve retrieval performance across all datasets, particularly in the *Clinical Narrative* setting where it help the system to recover from the baseline’s total failure (high Hit0). This demonstrates the accuracy and robustness of rewriting module, which acts as a crucial bridge between patient queries and biomedical literature. As demonstrated by the representative failure-recovery illustration in Appendix A.3, the LLM rewriting module does not rely on a single “perfect” rewrite. Instead, it generates a hierarchy of Boolean candidates (Strict → Relaxed) and executes them sequentially via our Adaptive Fallback mechanism, enabling collect sufficient evidence even for complex queries.

#### 4.1.2 Study on Evidence Sources and Depth

To thoroughly investigate the comparative value of different NCBI literature sources for CPQA, we evaluate both their retrieval distribution (Figure 3) and their downstream impact on generation performance (Table 2). We sampled 100 questions from CMMQA via stratified sampling for this analysis.

Figure 3 reveals two key insights regarding the comparative value of NCBI sources. First, while PubMed abstracts dominate both the initial retrieval pool and the top-5 evidence due to broader coverage (23.9M indexed citations vs. 8M PMC full-texts), their share decreases in the final top-5 evidence. This suggests that full-text sources provide richer, more discriminative content for CPQA, highlighting the value of integrating complementary data sources. Second, among the full-text sources, the relative growth of PMC Reviews outpaces that of other PMC papers (from 0.1/0.28 to 0.12/0.32), indicating that review articles possess a

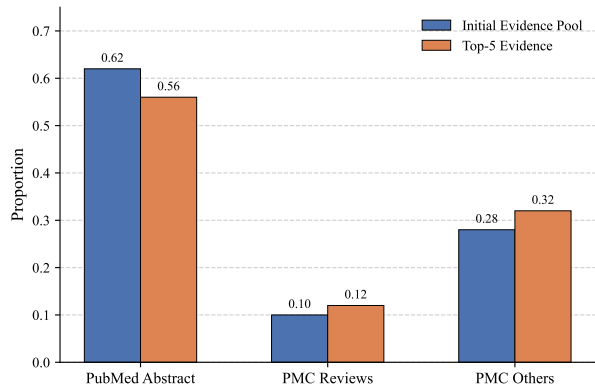


Figure 3: Distribution comparison between Initial Evidence Pool and Top-5 Evidence when CHSDR’s Retrieval Source involving PubMed Abstract, PMC Reviews and PMC Others

higher retrieval value for question answering.

| Source               | Acc          | Prec         | Rec          | F1           |
|----------------------|--------------|--------------|--------------|--------------|
| D1 (Abstracts)       | 44.00        | 41.12        | 40.12        | 39.14        |
| D1 + D2 (Reviews)    | <b>46.00</b> | <b>43.78</b> | <b>38.13</b> | <b>38.77</b> |
| D1 + D2 + D3 (Other) | 46.00        | 40.12        | 33.86        | 35.33        |

Table 2: Impact of data sources on retrieval performance (Negative Cancer QA Dataset). D1: PubMed Abstracts; D2: PMC Reviews; D3: Other PMC articles.

Table 2 further shows these patterns. Integrating PMC reviews with PubMed abstracts (D1 + D2) improves accuracy from 44.00% to 46.00%, while further including non-review full-texts (D1+D2+D3) maintains accuracy but degrades Precision, Recall, and F1. This divergence is mechanistically coherent: review articles synthesize findings across studies, matching the broad scope of patient queries. In contrast, non-review full-texts are typically context-specific; their core clinical outcomes can be represented by their abstracts, and retrieving their full texts introduces noise that overwhelms the model.

These empirical findings validate CoMeta’s solution to the **retrieval-depth paradox**: before passage retrieval, the metadata-aware pipeline should calibrates retrieval depth based on publication type. This step is a core component of CHSDR and distinguishes CoMeta from pipelines that treat all retrieved documents uniformly.

## 4.2 Analysis in Passage Retrieval

### 4.2.1 Comparison of Text chunking Strategies

Table 3 evaluates the robustness of SEOS across passage retriever, including lexical retrieval (BM25), domain-specific dense retrieval

| Dataset  | Method         | Standard Narrative |              |              |          | Clinical Narrative |              |              |          |
|----------|----------------|--------------------|--------------|--------------|----------|--------------------|--------------|--------------|----------|
|          |                | Hit@10             | MRR@10       | Recall@10    | Hit0↓    | Hit@10             | MRR@10       | Recall@10    | Hit0↓    |
| PubMedQA | E-utils        | 41.67              | 38.06        | 41.67        | 22       | 0.00               | 0.00         | 0.00         | 55       |
|          | <b>Adapt-E</b> | <b>48.33</b>       | <b>44.72</b> | <b>48.33</b> | <b>0</b> | 8.33               | 6.04         | 8.33         | <b>1</b> |
|          | MedCPT         | 10.00              | 9.17         | 10.00        | 0        | 3.33               | 2.50         | 3.33         | 0        |
|          | Hybrid         | 46.67              | 33.25        | 46.67        | 0        | <b>10.00</b>       | <b>5.28</b>  | <b>10.00</b> | 0        |
| BioASQ   | E-utils        | 52.44              | 35.49        | 28.87        | 18       | 1.22               | 1.22         | 1.22         | 76       |
|          | <b>Adapt-E</b> | <b>65.85</b>       | <b>41.94</b> | <b>41.61</b> | <b>0</b> | 50.00              | 28.93        | 33.32        | <b>1</b> |
|          | MedCPT         | 63.41              | 42.39        | 36.50        | 0        | 41.46              | 21.43        | 21.70        | 0        |
|          | Hybrid         | <b>80.49</b>       | <b>49.74</b> | <b>51.21</b> | 0        | <b>60.98</b>       | <b>35.36</b> | <b>40.07</b> | 0        |

Table 1: Ablation study on retrieval performance. PubMedQA and BioASQ have introduced in Figure 1. **Adapt-E** refers to our Boolean-constrained retrieval with adaptive query execution. **E-utils** refer to E-utilizes.

(MedCPT), and sentence-level embedding models (PubMedBERT-matryoshka). Results demonstrates the SEOS method’s superiority in text segmentation, outperforming fixed-parameters strategies. These fixed-length baselines represent the standard chunking strategies widely adopted across current mainstream medical RAG frameworks, including MedRAG, Self-BioRAG, and LlamaIndex. Comparisons with new advanced semantic chunking strategies will be included in future work.

| Splitter Configuration | Retriever Accuracy (%) |           |           |
|------------------------|------------------------|-----------|-----------|
|                        | PubMedBERT             | BM25      | MedCPT    |
| 512 (Overlap 0)        | 46                     | 20        | 22        |
| 512 (Overlap 32)       | 52                     | 18        | 24        |
| 512 (Overlap 128)      | 42                     | 16        | 22        |
| <b>SEOS (Ours)</b>     | <b>54</b>              | <b>36</b> | <b>38</b> |

Table 3: Retrieval accuracy across different text segmentation strategies on the Passage Retrieval Evaluation Dataset. SEOS outperforms fixed-size chunking.

The SEOS method excels by preserving natural and meaningful text boundaries based on semantic similarity and its variations, which enhances the retrievers’ ability to locate relevant information. Its advantages also include sentence overlap and automatic chunk-size adjustment tailored to the embedding model. Both Pubmedbert-matryoshka and MedCPT benefited from automatic chunk-size adjustment. The corresponding chunk-size adjustment rule obtained from research indicates that 128-word chunks with 32-word overlaps optimize BERT-based models for QA tasks (Wang et al., 2019). This finding can also be shown by the 512Overlap32’s top performance among fixed-length strategies. Despite a 512 chunk size, the actual retrieval text space is roughly 128, with metadata integrated consuming the remainder of the chunk, typically around 384 tokens.

## 4.2.2 Selection of Retrievers and Rerankers

While MedCPT is suited for document retrieval due to its training on query-article pairs, it does not utilize a sentence-transformer structure, which may limit its precision for shorter texts. Therefore, to find suitable embedding-reranker pairs for finer-grained passage retrieval, we designed a systematic retrieval evaluation based on LlamaIndex (Liu, 2022), a framework for building RAG systems. Specially, we used E-Utilities to query "cancer" to retrieve PubMed abstracts and PMC full-text articles as the text corpus for building an evaluation dataset, then used LLMs to generate pairs (query, context) from each chunk of the prepared text corpus, ensuring this evaluation was suitable for all data sources. In the experiment, we evaluated retrieval performance using Hits@5 and MRR@5, which aligns with the practical constraints of RAG systems, where the limited context window of the LLM generator requires a focus on retrieving the most relevant chunks (Tang and Yang, 2024). Results are shown in Figure 4.

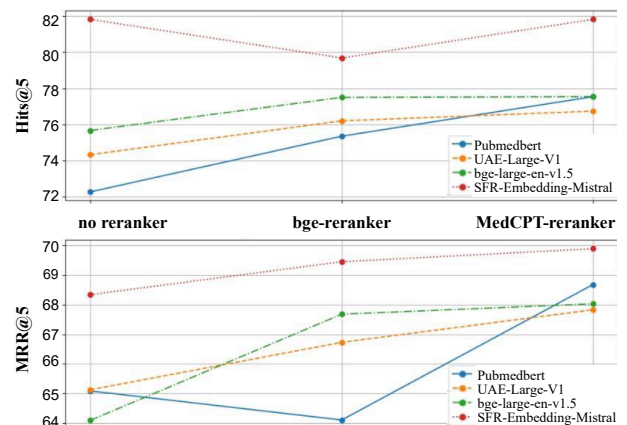


Figure 4: Performance of different reranker-retriever (reranker-embedding) pairs measured by Hit@5 (top) and MRR@5 (bottom).

**Domain Specialization Matters:** Pubmedbert-matryoshka, despite its smaller size and absence from the MTEB leaderboard, achieved the second-best performance when paired with the MedCPT-reranker. This suggests that the size of the embedding model is not the only determinant of effectiveness and that domain-specific fine-tuning or training can significantly improve performance by leveraging domain-specific understanding (Xiong et al., 2024). The importance of domain-specific features is also demonstrated by the MedCPT-reranker outperforming the general domain reranker in enhancing retrieval across all embedding models.

**Model Compatibility Matters:** PubMedBERT-Matryoshka performed poorly without a reranker but benefited substantially from the MedCPT Reranker. This synergy likely stems from the MedCPT Reranker being trained on negative samples from the MedCPT retriever (derived from Pubmedbert), aligning the reranker more effectively with PubMedBERT’s embedding space and enabling it to capture complementary information. Meanwhile, the BGE Reranker enhances the performance of BGE Embedding in terms of hits, also suggesting the importance of compatibility and complementarity between the embedding model and the reranker. However, the observed performance decline when BGE Reranker is paired with incompatible embeddings highlights the risks of mismatched reranker-embedding combinations. If a reranker cannot align with the embedding space or provide complementary semantic insights, it can fail to capture semantic nuances or introduce noise, leading to performance degradation. In conclusion, while rerankers can enhance retrieval, selecting compatible reranker-embedding combinations is crucial.

### 4.3 Discussion

**Overall Performance.** As shown in Table 4, CoMeta outperforms Naive RAG by 2.91% on average and LLM+CoT by 5.24% on CMMQA. These gains demonstrate the effectiveness of our query pipeline optimization, which includes two complementary retrieval improvements: (i) CHSDR tackle **staleness-semantic dilemma** and **retrieval-depth paradox** through LLM-based query rewriting, adaptive query execution, and metadata-aware hybrid search. (ii) SEOS and two-stage retrieval reduces semantic noise, ensuring that the model attends to most relevant, semantically segmented passages.

| Method    | MMLU         | MedQA        | MedMC        | PMQA         | BioASQ       | Avg          |
|-----------|--------------|--------------|--------------|--------------|--------------|--------------|
| LLM + CoT | 78.26        | 68.60        | 65.59        | 45.00        | 80.49        | 67.15        |
| Naive RAG | <b>82.61</b> | 67.44        | 65.59        | 56.67        | 81.71        | 69.48        |
| CoMeta    | <b>82.61</b> | <b>69.77</b> | <b>68.82</b> | <b>65.00</b> | <b>81.71</b> | <b>72.39</b> |

Table 4: Performance comparison on CMMQA datasets.

**Why Average Accuracy Understates CoMeta’s Contribution.** The average gain of 2.91% over Naive RAG understates our contribution due to three limitations: (1) *Blindness to ceiling effect of retrieval impact* The identical performance on MMLU and BioASQ does not imply a lack of improvement. For example, as shown in Table 1, the hybrid retriever (used by CoMeta) over MedCPT alone (63.41 → 80.49 under Standard Narrative), yet this is not reflected in end-task accuracy due to the ceiling effect, where unoptimized retrieval already provides sufficient context for correct generation. In contrast, on PubMedQA, which requires precise study-level matching, Hit@10 improvement (MedCPT → Hybrid: 10.00 → 46.67 under Standard Narrative) is accompanied by a +8.33% end-task accuracy gain, suggesting that retrieval quality is a key limiting factor under this setting. Together, these results suggest that the impact of retrieval is conditional on whether it constitutes the bottleneck. (2) *Blindness to retrieval reliability.* Accuracy metrics cannot reflect retrieval robustness. As shown in Table 5, E-Utilities fails on most narrative queries (PubMedQA: 55/60; BioASQ: 76/82), while Adapt-E reduces this to zero. For non-expert users, this constitutes a qualitative reliability shift that is invisible to accuracy-based evaluation. (3) *Blindness to evidence staleness.* A system may achieve high accuracy by retrieving outdated studies that remain correct for benchmark questions but are unsafe in real-world settings. CoMeta addresses this via metadata-aware filtering in CHSDR, but this safety dimension is orthogonal to accuracy and systematically overlooked by existing benchmarks.

| Dataset-Setting      | E-utils | Adapt-E (Ours) |
|----------------------|---------|----------------|
| PubMedQA – Standard  | 22/60   | 0/60           |
| PubMedQA – Narrative | 55/60   | 0/60           |
| BioASQ – Standard    | 18/82   | 0/82           |
| BioASQ – Narrative   | 76/82   | 0/82           |

Table 5: Zero-Hit failure Proportion.

**Adaptive Retrieval Necessity** Table 4 reveals that Naive RAG underperforms LLM+CoT on MedQA (67.44 vs. 68.60), despite having access to retrieved evidence. This is consistent with prior findings

that RAG can negatively impact the original outcome (Wang et al., 2023; Asai et al., 2023). Additionally, Table 1 also shows the optimal retrieval method should be query-specific. While the Hybrid retriever dominates in most settings because it captures complementary evidence missed by symbolic search, it is surpassed by Adapt-E on the PubMedQA Standard Narrative (48.33 vs. 46.67 Hit@10). This exception suggests that for queries derived directly from abstracts, precise lexical matching can be more effective than semantic blending, which might introduce 'semantic noise.' Together, these observations highlight the necessity of adaptive retrieval mechanisms, where the system dynamically determines whether and how to retrieve, rather than a one-size-fits-all solution.

**Generalizability of CoMeta** we evaluate Llama-3-8 on both CMMQA and MIRAGE-subset, with Claude-3-Haiku and GPT-4 on CMMQA (Table 6). We observed two trends: First, CoMeta consistently outperforms Naive RAG across all backbones, with larger gains for smaller models, consistent with the intuition that weaker parametric knowledge increases reliance on retrieval quality. This suggests CoMeta is particularly valuable for smaller open-source models (ideal choices in healthcare due to privacy and cost constraints) while still providing meaningful gains for stronger proprietary models. Second, For Llama-3-8B, the gain (+4.05) on MIRAGE-subset confirms CoMeta’s generalizability beyond oncology.

| LLM            | Method      | CMMQA                | MIRAGE-subset        |
|----------------|-------------|----------------------|----------------------|
| Llama-3-8B     | Naive RAG   | 65.70                | 62.57                |
|                | <b>Ours</b> | <b>70.80 (+5.10)</b> | <b>66.62 (+4.05)</b> |
| Claude-3-Haiku | Naive RAG   | 69.48                | –                    |
|                | <b>Ours</b> | <b>72.39 (+2.91)</b> | –                    |
| GPT-4          | Naive RAG   | 75.29                | –                    |
|                | <b>Ours</b> | <b>77.90 (+2.61)</b> | –                    |

Table 6: Performance across multiple LLM backbones. MIRAGE-subset: a stratified 10% subset of MIRAGE (764 questions,  $\geq 50$  per sub-dataset). “–”: omitted due to API budget constraints. We note that MIRAGE-subset results are indicative rather than conclusive given the smaller evaluation scale

**Latency and Cost Analysis.** E-Utilities latency is 1–2s; semantic search is milliseconds. The main overhead is LLM rewriting. In a "Patient Consultant" scenario, 10–20s latency is acceptable for a reliable, evidence-backed answer. Regarding infrastructure cost (See the system design comparison table in Appendix), CoMeta requires only a

static embedding file (~40GB) plus API access, versus a full BM25+MedCPT index (~400GB). Furthermore, LLM generation cost also decrease because metadata filtering and SEOS segmentation reduce the reader’s context window.

## 5 Conclusion

In this study, we proposed CoMeta, a framework designed for Cancer Patient QA (CPQA). Our work systematically resolves three critical bottlenecks in medical RAG. First, CHSDR overcomes the staleness-semantic dilemma by synergizing real-time, metadata-aware Boolean search (via Adpat-E) with domain-specific semantic retrieval (via MedCPT). Powered by Adapt-E, CHSDR eliminates the “Zero-Hit” barrier of term-based search, demonstrating exceptional robustness when processing noisy, informal patient narratives. Second, through metadata utilization, CoMeta resolves the retrieval-depth paradox, empirically analyzing the comparative retrieval value of different NCBI sources. Third, SEOS tackles contextual fragmentation by preserving natural clinical dependencies, offering a domain-agnostic chunking strategy that outperforms traditional fixed-length methods. Beyond the observed retrieval performance and end-task accuracy gains (Table 1 and Table 4), our results highlight CoMeta’s particular value in empowering the smaller open-source model (like LLama-3-8B in Table 6). Future work will include exploring advanced semantic segmentation and integrating adaptive retrieval mechanisms as mentioned before. Furthermore, recognizing that oncology care is inherently longitudinal, future work will also investigate evolving CoMeta into a memory-augmented agent system (Jiang et al., 2026) deployed on commodity mobile devices (Wan et al., 2026), enabling privacy-preserving, long-term health tracking.

## Limitations

While our evaluation provides robust insights, there are three limitations need to be addressed. First, the reliance on multiple-choice-style evaluation may introduce guessing bias, where random guesses can artificially boost accuracy. Second, the Clinical Narrative queries in CMMQA are synthetic. Although designed to rigorously stress-test the system against worst-case conditions (see Narrative queries Construction Prompt in Appendix A.6), potential distribution shifts relative to real patient queries cannot be entirely ruled out. Finally, ex-

isting benchmarks do not explicitly test temporal validity and metadata awareness. Consequently, the real-world safety advantages of our pipeline are underrepresented. We view our current scores as a conservative lower bound on CHSDR’s actual impact in fast-evolving domains like oncology.

## Acknowledgments

We would like to thank the anonymous reviewers and the Area Chairs for their constructive feedback, which has significantly improved the quality of this manuscript. The funding of Keeta AI for supporting the travel of this work.

## References

- Retrieval-augmented generation for large language models: A survey.
- Asma Ben Abacha, Yassine Mrabet, Mark Sharp, Travis R Goodwin, Sonya E Shooshan, and Dina Demner-Fushman. 2019. Bridging the gap between consumers’ medication questions and trusted answers. In *MEDINFO 2019: Health and Wellbeing e-Networks for All*, pages 25–29. IOS Press.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations*.
- Gordon V Cormack, Charles LA Clarke, and Stefan Buettcher. 2009. Reciprocal rank fusion outperforms condorcet and individual rank learning methods. In *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*, pages 758–759.
- Marti A Hearst. 1997. Text tiling: Segmenting text into multi-paragraph subtopic passages. *Computational linguistics*, 23(1):33–64.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.
- 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. *Bioinformatics*, 40(Supplement\_1):i119–i129.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38.
- Dongming Jiang, Yi Li, Songtao Wei, Jinxin Yang, Ayushi Kishore, Alysa Zhao, Dingyi Kang, Xu Hu, Feng Chen, Qiannan Li, and 1 others. 2026. Anatomy of agentic memory: Taxonomy and empirical analysis of evaluation and system limitations. *arXiv preprint arXiv:2602.19320*.
- Fan Jiang, Qionikai Xu, Tom Drummond, and Trevor Cohn. Boot and switch: Alternating distillation for zero-shot dense retrieval. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Di Jin, Eileen Pan, Nassim Oufattole, Weng Wei-Hung, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.
- Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2567–2577.
- 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. *Bioinformatics*, 39(11):btad651.
- Jonathan Kans. 2024. Entrez direct: E-utilities on the unix command line. In *Entrez programming utilities help [Internet]*. National Center for Biotechnology Information (US).
- Esther Landhuis. 2016. Scientific literature: Information overload. *Nature*, 535(7612):457–458.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Chaofan Li, Zheng Liu, Shitao Xiao, and Yingxia Shao. 2023. Making large language models a better foundation for dense retrieval. *Preprint*, arXiv:2312.15503.
- Sichu Liang, Linhai Zhang, Hongyu Zhu, Wenwen Wang, Yulan He, and Deyu Zhou. 2025. Rgar: Recurrence generation-augmented retrieval for factual-aware medical question answering. *arXiv preprint arXiv:2502.13361*.
- Valentin Liévin, Christoffer Egeberg Hother, Andreas Geert Motzfeldt, and Ole Winther. 2024. Can large language models reason about medical questions? *Patterns*, 5(3).

- Jerry Liu. 2022. Llamaindex. [Acceso el](#), 6.
- Chris J Lu, Alan R Aronson, Sonya E Shooshan, and Dina Demner-Fushman. 2019. Spell checker for consumer language (cspell). *Journal of the American Medical Informatics Association*, 26(3):211–218.
- Yiqun Miao, Yuhan Zhao, Yuan Luo, Huiying Wang, and Ying Wu. 2025. Improving large language model applications in the medical and nursing domains with retrieval-augmented generation: Scoping review. *Journal of Medical Internet Research*, 27:e80557.
- Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2022. [Mteb: Massive text embedding benchmark](#). [arXiv preprint arXiv:2210.07316](#).
- Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, and 1 others. 2023. Can generalist foundation models outcompete special-purpose tuning? case study in medicine. *Medicine*, 84(88.3):77–3.
- Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Conference on health, inference, and learning*, pages 248–260. PMLR.
- Stephen Robertson, Hugo Zaragoza, and 1 others. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, and 1 others. 2023a. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180.
- Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, and 1 others. 2023b. Towards expert-level medical question answering with large language models. [arXiv preprint arXiv:2305.09617](#).
- Yixuan Tang and Yi Yang. 2024. Multihop-rag: Benchmarking retrieval-augmented generation for multi-hop queries. [arXiv e-prints](#), pages arXiv–2401.
- Daniel Truhn, Jorge S Reis-Filho, and Jakob Nikolas Kather. 2023. Large language models should be used as scientific reasoning engines, not knowledge databases. *Nature medicine*, 29(12):2983–2984.
- George Tsatsaronis, Georgios Balikas, Prodromos Malakasiotis, Ioannis Partalas, Matthias Zschunke, Michael R Alvers, Dirk Weissenborn, Anastasia Krithara, Sergios Petridis, Dimitris Polychronopoulos, and 1 others. 2015. An overview of the bioasq large-scale biomedical semantic indexing and question answering competition. *BMC bioinformatics*, 16(1):138.
- Logesh Kumar Umapathi, Ankit Pal, and Malaikannan Sankarasubbu. 2023. Med-halt: Medical domain hallucination test for large language models. [arXiv preprint arXiv:2307.15343](#).
- Zhenchen Wan, Ce Chen, Runqi Lin, Jiaxin Huang, Tianxi Chen, Yanwu Xu, Tongliang Liu, and Mingming Gong. 2026. Mobile-vton: High-fidelity on-device virtual try-on. [arXiv e-prints](#), pages arXiv–2603.
- Yile Wang, Peng Li, Maosong Sun, and Yang Liu. 2023. Self-knowledge guided retrieval augmentation for large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10303–10315.
- Yinuo Wang, Baiyang Wang, Robert Mercer, Frank Rudzicz, Sudipta Singha Roy, Pengjie Ren, Zhumin Chen, and Xindi Wang. 2025. Trustworthy medical question answering: An evaluation-centric survey. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 27477–27490.
- Zhiguo Wang, Patrick Ng, Xiaofei Ma, Ramesh Nallapati, and Bing Xiang. 2019. Multi-passage bert: A globally normalized bert model for open-domain question answering. [arXiv preprint arXiv:1908.08167](#).
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. 2024. [Benchmarking retrieval-augmented generation for medicine](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 6233–6251, Bangkok, Thailand. Association for Computational Linguistics.
- Shijia Xu, Zhou Wu, Xiaolong Jia, Yu Wang, Kai Liu, and April Xiaowen Dong. 2026. [Self-correcting rag: Enhancing faithfulness via mmkp context selection and nli-guided mcts](#). [Preprint](#), arXiv:2604.10734.

## A Example Appendix

### A.1 Test Data Details

We construct evaluation dataset by applying a MeSH-based filter on six widely used medical datasets:

- PubMedQA (Jin et al., 2019) and BioASQ (Xiong et al., 2024) are QA datasets with answers from biomedical literature. PubMedQA (500 questions) is characterized by **single-document grounding**, where each manually annotated question is paired with a single corresponding PMID. In contrast, **BioASQ** (Xiong et al., 2024) (618 questions) features **multi-document grounding**, providing a list of multiple relevant PMIDs for each question. They are suitable for document retrieval evaluation due to varying evidence densities and the explicit annotation of gold-truth PMIDs.
- MedQA (Jin et al., 2021), MedMCQA (Pal et al., 2022), and the medical subsets in MMLU (Hendrycks et al., 2020) collected questions from medical exams. Specifically, MedQA (Jin et al., 2021) contains 1273 questions from the US Medical Licensing Examination (USMLE), focusing on complex clinical decision-making scenarios faced by professionals. MedMCQA is a large-scale dataset covering 21 medical subjects with 4183-question development set, focusing on evaluating across diverse healthcare topics. MMLU-med is a subset of the MMLU benchmark comprising 1,089 questions across six subtasks (e.g., professional medicine, human genetics), serving as a testbed for general biomedical reasoning.
- HealthSearchQA (Singhal et al., 2023b), comprising general consumer search queries, lacks definitive answers and is therefore unsuitable for accuracy assessment. However, its question lengths make it suitable for investigating the impact of question length on document retrieval methods. (Nori et al., 2023; Singhal et al., 2023a; Liévin et al., 2024).

First five datasets are multiple-choice QA datasets, because using multiple-choice questions can simplify evaluation, eliminate biases from text similarity computations or human annotation, and align with large-scale medical QA systems evaluations.

### A.2 The system design comparison Table

BM25 and E-Utilities address fundamentally different problems within the RAG pipeline; they are not simply interchangeable sparse retrievers. Actually, our target deployment is real-world clinical settings that require metadata-aware, up-to-date evidence access without maintaining massive local PubMed mirrors. Thus, a static BM25+MedCPT index is not a suitable baseline. Table 7 shows the key distinctions in indexing paradigms.

| Feature                             | BM25 + MedCPT | E-Utills + MedCPT      |
|-------------------------------------|---------------|------------------------|
| <b>Lexical Complementarity</b>      | ✓ (Strong)    | △ (Moderate)           |
| <b>Query Expansion</b> (ATM → MeSH) | ×             | ✓                      |
| <b>Metadata-Awareness</b>           | ×             | ✓                      |
| <b>Real-time Access</b>             | ×             | ✓                      |
| <b>Infrastructure</b>               | ~400GB index  | ~40GB index + Live API |

Table 7: System design comparison.

### A.3 The representative failure-recovery illustration

| Phase                        | Content / Query Candidates  | Hits |
|------------------------------|---|------|
| <b>Patient Query</b>         | "...The doc thinks she might have a tumor in her pancreas and they're gonna do a biopsy tomorrow. My sister and I are her caregivers and we think it's best if we don't tell her if it's canc'r..." | 0    |
| <b>Candidate 1</b> (Strict)  | "pancreatic tumor" AND "biopsy" AND "truth-telling" AND "caregivers"  | 0    |
| <b>Candidate 2</b>           | "pancreatic tumor" AND "biopsy" AND ("truth-telling" OR "disclosure")   | 10   |
| <b>Candidate 3</b>           | "pancreatic tumor" AND "biopsy" AND "caregivers"  | 10   |
| <b>Candidate 4</b>           | "pancreatic tumor" AND "caregivers" AND ("truth-telling" OR "bad news")   | 10   |
| <b>Candidate 5</b> (Relaxed) | "pancreatic tumor" OR "biopsy" OR "truth-telling"   | 32   |

Table 8: A representative failure-recovery case demonstrating the Adaptive Fallback mechanism. The system sequentially relaxes the Boolean logic until sufficient evidence is retrieved, efficiently overcoming the "Zero-Hit" barrier typical of raw sparse retrieval.

To illustrate how the Adapt-E module achieves a near-zero Hit0 rate on complex patient narratives, Table 8 presents a representative failure-recovery scenario. When a standard query fails, the system does not fail outright. Instead, it utilizes the LLM to generate a hierarchy of Boolean search strings, progressing from strict clinical constraints to relaxed conceptual matching. The system executes these candidates sequentially and triggers an early-stopping mechanism once sufficient documents are retrieved (e.g., stopping at Candidate 2 or 3). This mechanism avoids unnecessary relaxation and preserves retrieval precision while ensuring high robustness against noisy inputs.

#### A.4 The pseudo-code of SEOS

Algorithm 1 outlines the process of Semantic Enhanced Overlap Segmentation (SEOS). This method ensures that the document is split at natural semantic inflection points (where major topic shifts occur), while our sentence-level overlapping strategy guarantees that logical continuity is preserved across chunks.

---

#### Algorithm 1 Semantic Enhanced Overlap Segmentation

---

- 1: **procedure** SEOS(*document*, *k*, *embed\_model*)
  - 2:   Split *document* into sentences using *sentence\_splitter*
  - 3:   Merge each sentence with its *k* neighboring sentences
  - 4:   Get each sentence group's embedding by *embed\_model*
  - 5:   Use cosine dissimilarity change between adjacent groups to identify similarity inflection (semantic shift) points
  - 6:   Adjust *chunk\_size* based on *embed\_model*
  - 7:    $n \leftarrow \text{len}(\textit{document}) / \textit{chunk\_size}$
  - 8:   Breakpoints:  
     top *n* points with largest semantic shift
  - 9:   Split *document* into chunks at breakpoints
  - 10:   Overlap is added to per chunk by including the first sentence of the next chunk to preserve sequence order
  - 11:   **return** Chunks
  - 12: **end procedure**
- 

#### A.5 Retrieval Examples

#### Retrieval Examples (Output)

"step\_by\_step\_thinking":  
 "The question asks whether pneumologists' skepticism about chemotherapy and survival in advanced non-smallcell lung carcinoma is justified. The context information suggests that chemotherapy can prolong survival and improve quality of life for patients with advanced non-small cell lung cancer, especially those with good performance status. Several studies have shown that platinum-based chemotherapy regimens are the standard of care and can lead to improved survival compared to supportive care alone.

However, the context also suggests that there may be some skepticism among pneumologists about the benefits of chemotherapy. The context indicates that while chemotherapy can prolong survival, the benefits are moderate, with median and 1-year survival of 8--10 months and 30--40%, respectively. Additionally, the survival benefits vary depending on factors like disease stage, patient characteristics, and chemotherapy regimen.

Overall, the context information suggests that chemotherapy can provide meaningful survival and quality of life benefits for many patients with advanced non-small cell lung cancer, but the benefits are not dramatic and must be weighed against the potential for toxicity. This may explain some of the skepticism among pneumologists.

Based on the information provided, I believe the answer is B: no, the pneumologists' skepticism is not justified. While the benefits of chemotherapy are not as large as some might hope, the context information indicates that it can still provide important survival and quality of life benefits for many patients, and the risks are generally manageable. The skepticism may be understandable, but the evidence does not seem to fully support it."

"answer\_choice": "B"

#### A.6 LLM Prompts

### Prompt: Rewrite to Clinical Narrative Queries

```
prompt_searchapi = ""
### Role
You are a simulator for a Consumer Cancer QA System. Real users have no medical background,
including cancer patients and their family members.

### Input Data
RAW QUESTION: {question}

### Task
Identify the medical scenario, You need to transform RAW QUESTION into a a first-person query from
a real user.

### Requirements
1. Use layperson language (e.g., "spread" instead of "metastasize").
2. Include 1-2 minor typos (e.g., "cancr" for "cancer") or grammar slips.
3. Considering the real situation of non-expert users, rewrite from patient perspective.
4. Keep the core medical doubt unchanged so it can still be answered.

### Output Format
Return ONLY the rewritten string. No explanations.

""
```

### LLM with CoT + RAG

```
prompt_template = ""
You are a helpful medical expert.

QUESTION: {question}
POTENTIAL CHOICES: {options}
EVIDENCE: {context}

Task: answer a multi-choice medical question with provided evidence.

Instructions:
1.think step-by-step then choose the answer for QUESTION from the POTENTIAL CHOICES, using the
EVIDENCE.
2.Organize your output in a json formatted as Dict>{"step_by_step_thinking":Str(explanation),
"answer_choice": Str{{A/B/C/...}}}.

Note:
1.Only Dict in output is needed. No need for other text.
2.Your responses will be used for research purposes only, so please have a definite answer.
3.Keep your answer ground in the facts of the EVIDENCE.

""
```

## Prompt: Rewrite for boolean queries and set time constraints

```
prompt_searchapi = ""
```

You are a helpful expert in information extraction.

```
QUERY: {query}
```

Task: generate boolean search queries for QUERY

Instructions:

1. Preprocessing Stage:

- **Correct** any orthographic or grammatical errors in the QUERY.
- **Analyze** the Question Type and User intent to identify core concepts and **SEPARATE** the "Noise" from the "Core Dilemma".
- **Extract** the core terms. Especially identify **PICO** (Population, Intervention, Outcome)
- **Abstraction** specific scenario details into general medical concepts (e.g. "Don't tell" -> "Truth-telling").
- **Time**: If time constraints (like "recent") appears, set reasonable year\_range (e.g., 5). Otherwise None.

2. Boolean Queries Generation Stage:

- Provide 5 possible boolean search queries using core terms.
- Start with the most specific query (Strict AND of all terms).
- Progressively relax the strict query based on your Intent Analysis in stage 1:  
[Strategy A (Synonyms): If a term is a CORE concept, relax it by adding synonyms];  
[Strategy B (Pruning): If a term is a SPECIFIC but NOT CORE detail (e.g., "in Italy"), remove]  
[Strategy C (Logical Relaxation)]: Switch AND to OR between distinct core concepts

Examples:

# Example 1:

Input: "prevalence of concurrent diabetes and cardiovascular disease in women in Beijing"

Output: (

```
[
    "prevalence" AND "diabetes" AND "cardiovascular disease" AND "women" AND "Beijing"',
    "prevalence" AND "diabetes" AND "heart disease" AND "women"',
    "epidemiology" AND "concurrent" AND "diabetes" AND "cardiovascular disease"',
    "diabetes" AND "cardiovascular disease" AND "comorbidity"',
    "prevalence" AND ("diabetes" OR "cardiovascular disease")'
```

```
],
```

None

```
)
```

# Example 2:

# Abstractions: Stroke->End of life; Disconnect->Withdrawal; Son/Daughter->Family Conflict

Input: "75-year-old man with severe stroke and ventilator. Daughter wants to disconnect the machine but son disagrees."

Output: (

```
[
    "withdrawal of life support" AND "family conflict" AND "surrogate decision making"',
    "withdrawal of life support" AND ("family conflict" OR "disagreement") AND "decision making"',
    "withdrawal of life support" AND "family conflict"',
    "surrogate decision making" AND "family conflict"',
    "withdrawal of life support" OR "end of life care"'
```

```
],
```

None

```
)
```

# Example 3: Time Sensitive

Input: "latest treatments for glioblastoma in the last 8 years"

Output: (

```
[
    "Glioblastoma"[MeSH] AND "Therapeutics"[MeSH]',
    "Glioblastoma" AND "Treatment" AND "Novel"',
    "Brain Neoplasms" AND "Therapy"',
    "Glioblastoma" AND "Immunotherapy"',
    "Glioblastoma" OR "Brain Tumor"'
```

```
],
```

```
8
```

```
)
```

Note: Return ONLY the Python tuple (list\_of\_strings, year\_range\_int). No other text.

```
"""
```