

Retrieval-Augmented Generation in Biomedicine: A Survey of Technologies, Datasets, and Clinical Applications

Jiawei He

joewellhe@gmail.com

University of Geneva

Boya Zhang

University of Geneva

Hossein Rouhizadeh

University of Geneva

Yingjian Chen

Henan University

Rui Yang

National University of Singapore

Jin Lu

LvZhiDao Information Technology Co., Ltd.

Xudong Chen

Hunan City University

Nan Liu

National University of Singapore

Douglas Teodoro

University of Geneva

Research Article

Keywords: Biomedical RAG, Large Language Model, Retrieval-Augmented Generation, Information Retrieval, Natural Language Processing, Biomedical NLP

Posted Date: December 15th, 2025

DOI: <https://doi.org/10.21203/rs.3.rs-8330917/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: No competing interests reported.

1 **Retrieval-Augmented Generation in Biomedicine: A Survey of Technologies, Datasets, and** 2 **Clinical Applications**

3 Jiawei He^{1,2*}, Boya Zhang², Hossein Rouhizadeh², Yingjian Chen³, Rui Yang⁴, Jin Lu⁵, Xudong
4 Chen¹, Nan Liu⁴, Douglas Teodoro^{2*}

5 ¹Hunan City University, Yiyang, Hunan, China.

6 ²University of Geneva, Geneva, Geneva, Switzerland.

7 ³Henan University, Kaifeng, Henan, China.

8 ⁴ National University of Singapore, Singapore, Singapore.

9 ⁵LvZhiDao Information Technology Co., Ltd., Changsha, Hunan, China.

10
11 *Corresponding author(s). E-mail(s): joewellhe@gmail.com; douglas.teodoro@unige.ch;

12 13 **Abstract**

14 Large language models (LLMs) in biomedicine face a fundamental conflict between static
15 parameter knowledge and the dynamic nature of clinical evidence. Retrieval-Augmented
16 Generation (RAG) addresses this by grounding generation in external data, yet it introduces new
17 complexities in latency and architecture. This survey synthesizes the biomedical RAG landscape
18 (2020–2025), classifying systems into naive, advanced, and modular paradigms. Beyond a
19 technological taxonomy, we formalize the biomedical RAG trilemma, identifying the inherent
20 trade-offs between reasoning depth, inference latency, and data privacy that constrain current
21 clinical deployment. We analyze how recent agentic workflows enhance diagnostic reasoning
22 but risk prohibitive latency, and how privacy constraints dictate the choice between powerful
23 cloud-based models and local deployment. Finally, we outline the alignment gap in multimodal
24 RAG and propose future directions for self-correcting, verifiable clinical agents.

25
26 **Keywords:** Biomedical RAG, Large Language Model, Retrieval-Augmented Generation,
27 Information Retrieval, Natural Language Processing, Biomedical NLP

29 **1 Introduction**

30 While large language models (LLMs) demonstrate reasoning capabilities essential for clinical
31 decision support, their reliance on static parametric memory creates a fundamental misalignment
32 with the dynamic nature of biomedical evidence [1-4]. In clinical environments where guidelines
33 evolve daily, the static knowledge limitation of standard LLMs poses severe risks to patient
34 safety, manifesting as hallucinations that are plausible yet factually incorrect [5].

35 Retrieval augmented generation (RAG) aims to resolve this by decoupling reasoning (the
36 generative model) from knowledge storage (the retrieval index) [6]. This architectural separation
37 transforms the LLM from a memorizer into a reasoner, allowing for the real-time integration of
38 up-to-date clinical data and knowledge without expensive retraining, while ensuring that
39 generated claims can be grounded in verifiable citations [1-2].

40 However, standard plug-and-play RAG pipelines frequently fail in biomedical settings due to
41 the lexical-semantic gap. Unlike general domains, biomedical retrieval must navigate complex
42 ontologies (e.g., UMLS) where keyword matching fails to align lay terminology (e.g., *heart*
43 *attack*) with clinical nomenclature (e.g., *acute myocardial infarction*), necessitating specialized
44 dense retrieval and re-ranking architectures [7-9].

45 Existing reviews have predominantly focused on the *training* paradigms of medical LLMs (e.g.,
46 Liu *et al.* [10] and He *et al.* [2]), general-domain RAG taxonomies (e.g., Fan *et al.* [1] and Gao
47 *et al.* [11]), or multimodal LLMs (Xiao *et al.* [4]). These works largely treat the retrieval
48 component as a black box, overlooking specificities of sparse and dense retrieval in high-stakes
49 clinical workflows. Furthermore, they fail to address the latency vs. accuracy dilemma inherent
50 to deploying modular RAG systems in real-time hospital environments.

51 Our work differs from existing surveys by moving beyond enumeration to provide a synthesis
52 of the RAG technologies in biomedicine. The key contributions of this survey are:

- 53 1. A novel taxonomy and evolutionary timeline: We classify biomedical RAG systems into
54 distinct architectural paradigms (naive, advanced, and modular) and provide a timeline
55 chart showcasing the co-evolution of domain-specific retrievers (e.g., MedCPT) and
56 generative models.
- 57 2. Formalization of the biomedical RAG trilemma: We synthesize the architectural trade-
58 offs into a unified theoretical framework, the biomedical RAG trilemma, analyzing how
59 current systems must balance reasoning depth, inference latency, and data privacy.
- 60 3. Critical assessment of data and reliability: We analyze the implications of dataset biases,
61 including Anglocentric and demographic imbalances in sources like MIMIC-IV, and

62 discuss essential methodologies for measuring confidence and quantifying uncertainty in
63 clinical outputs.

64 4. Comparative clinical analysis: We present concrete case studies (e.g., rare disease
65 identification) that explicitly contrast conventional NLP performance with RAG-
66 enhanced workflows, demonstrating the tangible benefits in accuracy and explainability.

67

68 2 Scope and Survey Taxonomy

69 To provide a comprehensive synthesis of the field, we adopted a systematic mapping approach.
70 We queried major repositories (PubMed, arXiv, ACL Anthology) using the Boolean logic: (RAG
71 OR 'retrieval-augmented') AND (biomedical OR clinical OR medical). To distinguish our
72 analysis from general LLM surveys, we selected studies that i) explicitly couple a non-
73 parametric retriever (vector database, knowledge graph (KG), or search API) with a parametric
74 generator; ii) address specific biomedical challenges (e.g., ontology alignment, citations) rather
75 than generic question-answering (QA); and iii) prioritize peer-reviewed works or open-source
76 frameworks (e.g., MedCPT [12], Self-BioRAG [13]) that allow for architectural deconstruction.
77 We excluded works relying solely on prompt engineering or long-context windows, as these lack
78 the retrieval-based reasoning mechanism that is the focus of this review.

80 3 Problem Formulation

81 We formalize Biomedical RAG as the maximization of the generation probability $P(y|x)$ for a
82 clinical query x , conditioned on a latent set of retrieved documents Z .

83 3.1 Mathematical Framework

84 The generation process is defined as:

$$85 P(y|x) = \sum_{z \in Z} P_{\theta}(y|x, z)P_{\phi}(z|x),$$

86 where $P_{\phi}(z|x)$ represents the retrieval likelihood (optimized via dense/sparse alignment) and
87 $P_{\theta}(y|x, z)$ represents the generative reasoning capability. In modular RAG (Section 5.3), this
88 expands to include a policy π where the model iteratively decides to retrieve, reason, or
89 terminate.

90 3.2 General Execution Framework

91 To standardize the comparison of architectures discussed in Section 4, we define the canonical
92 biomedical RAG workflow in Algorithm 1.

Algorithm 1: Canonical biomedical RAG workflow

Input clinical query x , document index D , retriever R , generator G , threshold τ

Output clinical response y , citations C

1 **Step 1: retrieval**

2 $Z_{raw} \leftarrow R(x, D)$ Retrieve top- k candidates (sparse/dense)

3 **Step 2: refinement (advanced RAG)**

4 $Z_{ranked} \leftarrow \text{Rerank}(x, Z_{raw})$ Filter noise via cross-encoder

5 **Step 3: generation and verification (modular RAG)**
6 IF confidence(G, c) < τ
7 $y \leftarrow \text{self-correct}(x, Z_{\text{ranked}})$ *Iterative refinement (e.g., Self-BioRAG)*
8 ELSE
9 $y \leftarrow G(x, Z_{\text{ranked}})$
10 RETURN y

93

94

95 4 Biomedical RAG Architectures and Taxonomy

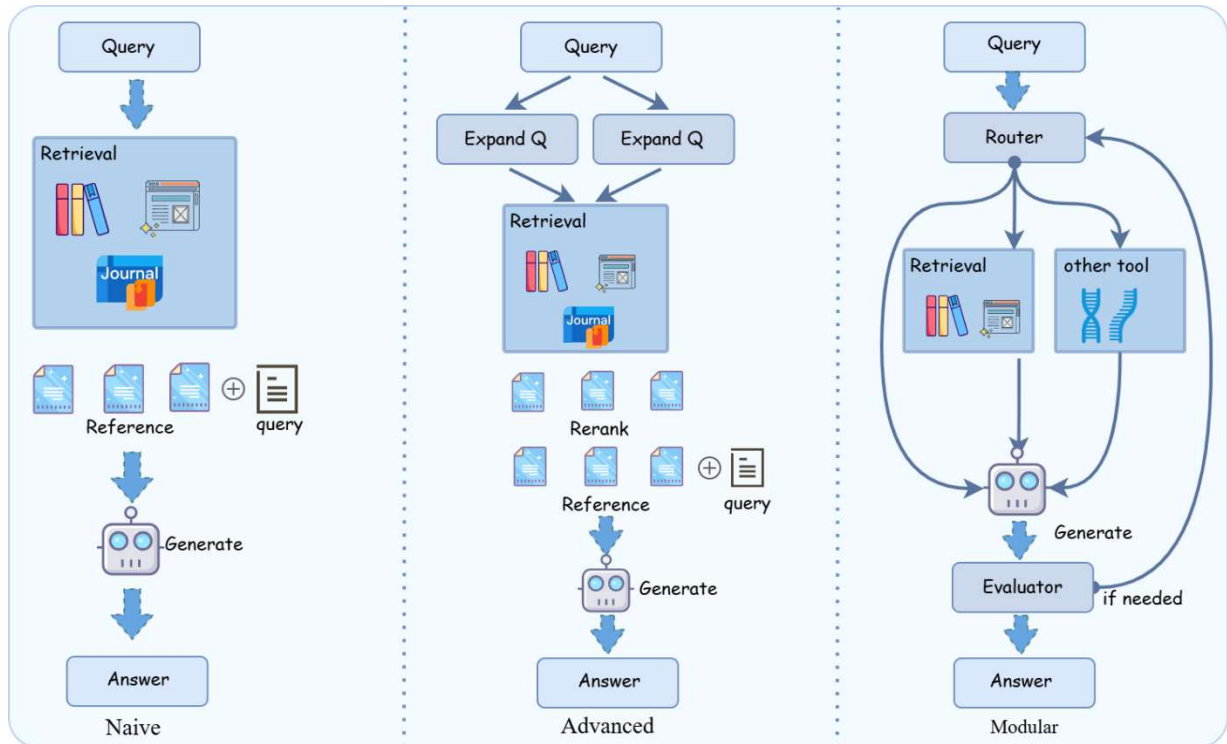
96 We classify architectures based on how they approximate and optimize this equation: *naive* RAG
97 (linear approximation), *advanced* RAG (optimization of P_ϕ), and *modular* RAG (introduction
98 of a policy π to navigate Z). These paradigms reflect a shift from simple fact retrieval to complex
99 clinical reasoning. While naive systems suffice for surface-level queries, the field is trending
100 toward modular architectures to handle multi-hop reasoning required for more complex clinical
101 scenarios, such as differential diagnosis. A comparative analysis of these biomedical RAG
102 approaches is shown in **Table 1**.

103 The foundational architecture, often termed *naive* RAG, follows a traditional retrieve-read
104 process. It assumes a single-hop linear approximation where the top- k documents retrieved by a
105 sparse or dense index are concatenated directly into the context window. This assumes that the
106 retrieval score $S_R(q, d)$ is a perfect proxy for semantic relevance. As illustrated in Figure 1, a
107 query q is passed to a retriever $R(q, K)$, which identifies a set of documents $D = \{d_1, \dots, d_n\}$
108 from a knowledge source K . These documents are concatenated directly with the query and fed
109 into the LLM to generate a response. This architecture might suffer from grounded hallucination
110 due to lexical-semantic dissonance. For example, if a retriever fetches guidelines on *viral*
111 *pneumonia* for a *bacterial pneumonia* query due to keyword overlap or semantic relatedness, the
112 LLM might generate a clinically fluent but incorrect treatment plan. Despite this risk, naive RAG
113 remains the standard where latency is prioritized over complex reasoning.

114 *Advanced* RAG introduces a post-retrieval optimization stage: re-ranking. While the initial
115 retrieval relies on computationally efficient bi-encoders ($O(N)$ via approximate nearest neighbor
116 search), advanced RAG employs cross-encoders ($O(K)$) to perform full self-attention over the
117 query-document pairs. While this might significantly improve Precision@ K by filtering or
118 downgrading irrelevant context, it introduces a latency bottleneck (often $> 500\text{ms}$) [14]. Thus,
119 while advanced RAG might be preferred for physician-facing clinical decision support systems
120 (CDSS), where accuracy is non-negotiable, response time might be prohibitive for real-time
121 conversational agents.

122 Recent developments have led to *modular* RAG, which creates flexible, non-linear workflows
123 suited for complex biomedical reasoning. Unlike the static pipelines of naive and advanced
124 RAG, modular systems employ iterative retrieval or routing mechanisms. In this sense, modular
125 systems represent a shift toward agentic RAG. Systems, such as Self-BioRAG [13] and
126 GeneGPT [15], introduce a routing policy $\pi(a|s)$, where the model iteratively decides an action
127 a , such as whether to *retrieve* from a vector store, *compute* using an external tool (e.g., a genomic

128 calculator), or *generate* a final answer, conditioned on the current context of the system s . While
 129 this allows for 'System 2' reasoning (deliberative thought) required for differential diagnosis, it
 130 introduces a latency penalty (refer to the trilemma in Section 8), subjecting the system to non-
 131 deterministic execution times.



132
 133 **Figure 1:** Basic biomedical RAG framework.

134
 135 **Table 1:** Comparative analysis of biomedical RAG paradigms (naive, advanced, modular).

Feature	Naive RAG	Advanced RAG	Modular RAG
Mathematical Model	linear approximation	optimized $P(z x)$ (re-ranking)	iterative policy $\pi(a s)$
Inference latency	low (< 2s) [16]	medium (2-5s) [17]	high (> 10s) [17]
Clinical failure mode	hallucination via irrelevant documents	latency timeout	infinite loops / API failures
Best use case	patient education / FAQ	CDSS	rare disease diagnosis / genomics

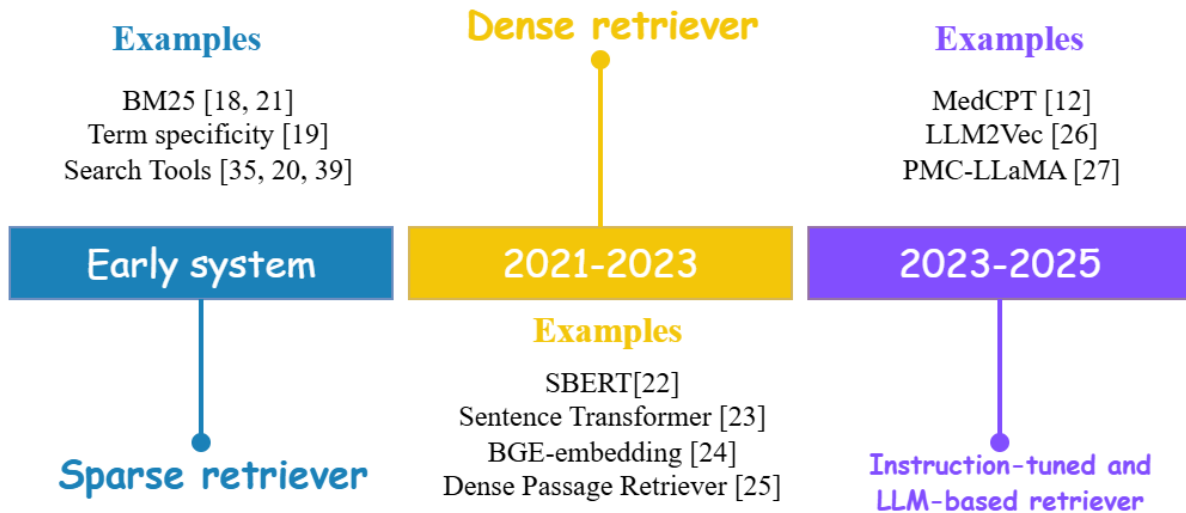
136

137

138 5 Technology

139 We trace the evolution of biomedical retrievers from lexical matching (2020) to the inclusion of
140 instruction-tuned semantic alignment (2025). This shift was necessitated by the failure of
141 keyword-based systems to capture the nuance of clinical queries (e.g., *heart failure* implying
142 *cardiac insufficiency*), leading to the adoption of complementary dense embeddings that align
143 vector space with medical ontologies.

144 5.1 Retriever



145

146 **Figure 2:** The evolutionary timeline of retrievers.

147 Retrievers function as the non-parametric memory access mechanism for the RAG system [11].
148 As illustrated by the evolutionary timeline in Figure 2, biomedical retrieval strategies have
149 shifted significantly over the surveyed period (2020–2025). Early systems predominantly relied
150 on sparse retrieval (BM25) [17-20] for precise keyword matching, e.g., of drug names and
151 symptoms. The introduction of BERT-based dense retrievers (2021–2022) enabled semantic
152 matching [21-24], while the most recent wave (2023–2025) utilizes instruction-tuned retrievers
153 like MedCPT and LLM-based embeddings to align retrieval directly with clinical reasoning
154 tasks [12, 25, 26].

155 5.1.1 Sparse Retrieval

156 Sparse retrieval, exemplified by algorithms like BM25 [18] and TF-IDF [19], operates on the
157 principle of lexical matching, where documents are indexed based on term occurrences rather
158 than semantic meaning. While conceptually simple, as shown in Table 2, these methods remain
159 indispensable in biomedicine due to the *lexical gap*. Neural models frequently exhibit semantic
160 drift, associating clinically distinct antonyms (e.g., *hypoglycemia* vs. *hyperglycemia*) due to their
161 shared vector space neighborhoods. In contrast, sparse retrieval enforces exact matching,

162 ensuring that specific alphanumeric variants (e.g., ICD-10 E11.9) are retrieved accurately,
 163 preventing dangerous hallucinations in downstream generation. Thus, despite the rise of neural
 164 methods, sparse retrieval remains a competitive baseline [19, 20] and an essential component in
 165 hybrid systems [27-30], combining sparse and dense retrieval techniques. Recent studies
 166 consistently demonstrate that fusing BM25 with dense vectors enhances generalization,
 167 effectively handling the lexical gap where precise medical nomenclature does not align with
 168 semantic clusters [31-33].

169

170 **Table 2:** Sparse retrieval methods in biomedical RAG systems.

Method name	Retrieval method	Task	Year
Clinfo.ai [35]	Entrez API [36]	medical QA	2023
MEDRAG toolkit [7]	Hybrid (BM25, SPECTER, Contriever, MedCPT)	medical QA	2024
SeRTS [33]	Self-Rewarding Tree Search based BM25	biomedical QA	2024
Dual RAG System [28]	Hybrid (dense+BM25 with language-specific tokenizers)	diabetes management	2024
MedExpQA [29]	BM25 and MedCPT	medical QA	2024
CliniqIR [30]	BM25 and MedCPT	diagnostic decision support	2024
VAIV Bio-Discovery [31]	Hybrid (dense+BM25)	biomedical knowledge discovery	2024
Shuangjia et. al [37]	<i>Azure AI Search</i> with keyword matching	genetic variant annotation	2025
Eun Jeong et. al [21]	BM25	medical QA	2025
Sudeshna et. al [20]	Whoosh [38]	medical QA	2025
LITURAt [39]	Entrez API [36]	scientific data analysis	2025

171

172 5.1.2 Dense Retrieval

173 Dense retrievers encode queries and documents into a continuous vector space where semantic
 174 similarity is measured via cosine distance [40]. The most effective biomedical retrievers achieve
 175 domain alignment through two primary adaptation strategies: continued pre-training on
 176 biomedical corpora (e.g., BMRETRIEVER [32]) and contrastive fine-tuning (e.g., MedCPT
 177 [12]). This training objective ensures that a query embedding q is closer to its relevant positive
 178 document d^+ than to a set of negatives D^- , mitigating the semantic disparities between general
 179 and biomedical language domains.

180

181 We classify dense retrievers into a hierarchy of utility vs. privacy (Table 3):

- 182 1. Type 1 - General PLMs: PLMs, such as BERT, often fail to capture medical nuances [21-
183 24, 38-52].
- 184 2. Type 2 - Commercial APIs: While offering superior semantic reasoning (e.g., OpenAI),
185 they pose significant data residency risks (GDPR/HIPAA), limiting their use to non-
186 patient-identifiable data [27, 53-64].
- 187 3. Type 3 - Domain-adapted retrievers: Models like MedCPT balance privacy and
188 performance but require substantial fine-tuning resources [13, 65-72].
- 189 4. Type 4 - LLM-based embeddings: The emerging state-of-the-art, allowing for
190 instruction-based retrieval (e.g., retrieve only *contraindications*), though at a higher
191 computational cost [25, 73, 74].

192

193 **Table 3:** Classification of dense retriever types in biomedical RAG systems. A more detailed
194 taxonomy is provided in Appendix Table A1.

Type	Definition	Representative example	Trade-off (Trilemma)
1 - General PLMs	off-the-shelf BERT models without domain tuning.	AskFDALabel [78]	high availability / low accuracy
2 - Commercial APIs	proprietary endpoints (e.g., OpenAI) with superior semantic alignment.	DRAGON-AI [56]	high accuracy / low privacy
3 - Domain-adapted	open-source models fine-tuned on biomedical corpora (PubMed).	MedCPT [12]	balanced privacy and accuracy / high training cost
4 - LLM-based	instruction-tuned embeddings allowing task-specific queries.	BiomedRAG [5]	high accuracy / high latency

195

196 Implementation of dense retrievers relies on vector stores like Faiss [79] and Chroma [80] (Table
197 4). However, for clinical deployment, the choice of infrastructure is dictated by privacy; local
198 instances are preferred over managed cloud solutions to ensure PHI (protected health
199 information) remains on-premise. While cloud-native stores offer auto-scaling, hospital on-
200 premise introduces an infrastructure and maintenance overhead that is often underestimated in
201 research papers.

202

203 **Table 4:** Common vector stores in biomedical RAG

Store	Description
Chroma [80]	AI-native open-source vector database with built-in functionality
Faiss [79]	Library for efficient similarity search and clustering of dense vectors
Pinecone [81]	Production-ready vector database for similarity search at scale
Weaviate [82]	Graph-based vector search engine with semantic capabilities
HNSW [51]	Hierarchical navigable small world graphs for approximate nearest neighbor search

204
 205 **The Precision-Recall Trade-off.** The analysis of current retrieval paradigms reveals a tension
 206 between *semantic recall* and *lexical precision*. While dense retrievers (e.g., MedCPT) excel at
 207 mapping symptoms to diagnoses, they frequently fail on exact alphanumeric matches (e.g., ICD-
 208 10 codes). Conversely, sparse retrieval ensures lexical safety but misses semantic nuance. This
 209 necessitates the hybrid architectures seen in recent clinical systems, where the computational
 210 cost of dual-indexing is accepted as the price for patient safety.

211
 212 **5.2 Reranker**

213 Rerankers in biomedical RAG systems serve as intermediary components that refine and
 214 prioritize the relevance of initially retrieved documents before they are processed by the
 215 generative models [83]. Positioned between retrievers and generators, rerankers apply
 216 sophisticated algorithms to re-score and re-order candidate documents based on their contextual
 217 relevance to the query, effectively functioning as a second-stage filtering mechanism that
 218 enhances the precision of information ultimately provided to the generation component.

219 Biomedical RAG rerankers can be categorized into distinct types based on their underlying
 220 methodologies:

- 221 • **Statistical rerankers** apply traditional information retrieval techniques; for instance,
 222 CliniqIR [30] and MEDRAG [73] both implement Reciprocal Rank Fusion (RRF) to
 223 combine multiple retrieval signals, while DRAGON-AI [56] employs Maximal Marginal
 224 Relevance (MMR) to balance relevance with diversity.
- 225 • **Content-based rerankers**, in contrast, focus on content similarity; notably, WeiseEule
 226 [73] introduces a keyword frequency-based ranking system specifically designed for
 227 biomedical literature retrieval.
- 228 • **Model-based rerankers** employ transformer architectures to perform contextual
 229 relevance assessment, effectively functioning as a domain adaptation strategy. For
 230 instance, MedCPT [12] trains a specialized cross-encoder on PubMed search logs, while

231 BiomedRAG [5] implements a domain-specific chunk scorer to prioritize evidence
232 containing causal medical relationships. These specialized mechanisms capture subtle
233 signals that general embedding models often miss.

234 • **LLM-based rerankers** leverage foundation models' reasoning capabilities. For example,
235 GNQA [50] employs GPT-3.5/4/4o as zero-shot binary classifiers for relevance
236 determination, and LmRaC [84] implements an LLM-based paragraph usefulness
237 assessment with a minimum threshold score of 7 (on a 10-point scale).

238 • **Hybrid approaches** combine multiple strategies to enhance performance. For example,
239 Clinfo.ai [35] combines LLM classification with BM25 to re-rank retrieved articles.

240 In practice, we observe a standard funnel architecture in clinical deployment: a computationally
241 cheap retriever (bi-encoder/BM25) fetches 50-100 candidates, which are then re-scored by an
242 expensive cross-encoder (model-based). While LLM-based reranking (e.g., GPT-4+) offers the
243 highest accuracy, its latency (>1s per query) [85] hinders its application in real-time
244 patient/physician-facing interfaces.

245
246 **The Accuracy-Latency Bottleneck.** In the context of the biomedical trilemma, rerankers
247 represent the primary constraint on inference latency. While cross-encoders (model-based)
248 maximize reasoning depth by effectively filtering noise, they increase latency by an order of
249 magnitude compared to bi-encoders. This forces a hard choice in architectural design: real-time
250 conversational agents must often bypass sophisticated reranking (sacrificing precision), whereas
251 asynchronous CDSS can afford the computational cost for maximum safety.

252

253 **5.3 Generation**

254 The generation component represents the last element in biomedical RAG systems, responsible
255 for synthesizing retrieved information into human-like responses that address user queries. In
256 biomedical contexts, this process requires balancing natural language fluency with domain-
257 specific accuracy, as generated content often informs healthcare decisions [86]. This section
258 examines various generation models employed in recent biomedical RAG studies.

259 Current generation models employed in biomedical RAG systems can be categorized into three
260 distinct classes based on their development approach and specialization: general-domain open-
261 source models, commercial proprietary models, and biomedical-specialized models. These
262 categories present different characteristics in terms of accessibility, customization potential, and
263 domain-specific performance. Table 5 showcases representative models from each category and
264 highlights recent studies that utilize these models (a more detailed taxonomy of the surveyed

265 work is provided in Appendix Table A2). It is important to note that our analysis of biomedical-
266 specialized LLMs focuses specifically on models trained directly on biomedical or healthcare
267 corpora, deliberately excluding applications developed through prompt engineering or chain-of-
268 thought (CoT) methodologies.

269 As shown in Table 5, the majority of LLMs employ decoder-only architectures, with T5 being a
270 notable exception. This architectural preference stems from several inherent advantages for
271 generative tasks, particularly the natural suitability for next token prediction and significantly
272 improved parameter efficiency [87]. Among open-source LLMs and Llama 2 [88] is the most
273 popular, with 15 biomedical applications built upon it. LLaMa 3 [89], Mistral 7B [90], and T5
274 [91] are also widely used. These open-source LLMs have various sizes, from 0.06B (T5) to 671B
275 (Deepseek R1), providing flexible options for building Biomedical RAG systems. Recently,
276 several powerful general open-source LLMs have emerged, including, Mistral Small 3.1 [92]
277 and Gemma 3 [93]. Although we have not found these newer models being used for biomedical
278 RAG systems in our survey, we believe they show promise due to their outstanding performance
279 in general domain tasks.

280 For commercial LLMs, the ChatGPT series APIs maintain a dominant position, presumably
281 attributable to their superior performance. Nevertheless, several studies have demonstrated that
282 Claude and Gemini achieve comparable efficacy when implemented in biomedical RAG
283 applications [55, 91, 92], thereby offering viable alternatives to ChatGPT. Furthermore, it is
284 noteworthy that both DeepSeek and Mistral have established dual accessibility paradigms:
285 providing open-source parameters for their base models while simultaneously offering
286 commercial API services.

287 Biomedical specialized LLMs are developed through domain adaptation strategies, including
288 continued pre-training on medical corpora (e.g., BioMistral [96], ChatENT [61]), supervised
289 fine-tuning (e.g., MEDGENIE [6]), and parameter-efficient approaches like low-rank adaptation
290 (LoRA) used in AskFDALabel [78]. Despite these specialization efforts, recent studies [72, 94,
291 95] demonstrate that within RAG frameworks, these models do not consistently outperform
292 general open-source alternatives. We attribute this to the reasoning-knowledge decoupling
293 effect. In RAG, the *knowledge* comes from the retrieved documents, so the primary role of an
294 LLM is *reasoning* and *synthesis*. Massive generalist models possess superior reasoning
295 capabilities derived from diverse training data. In contrast, smaller, domain-specific models
296 often sacrifice reasoning power for vocabulary memorization, which is counterproductive in a
297 RAG architecture.

298

299 **Table 5:** Taxonomic classification of LLMs for embedding applications.

Tier	Model class	Representative models	Reasoning capability	Privacy Safety	Ideal Use Case
1 - Cloud SOTA	proprietary API	GPT-4o, Claude 3.5 Sonnet, Gemini 1.5 Pro	very high (system 2)	low (data residency risk)	differential diagnosis, complex multi-step reasoning
2 - Local generalist	open-weights (large)	Llama 3 (70B), Qwen 2.5 (72B), DeepSeek-R1	high	high (on-premise)	hospital-side CDSS, summarization of PHI
3 - Local specialist	domain-adapted (small)	BioMistral (7B), PMC-LLaMA (13B)	medium	high (on-premise)	edge devices, specific clinical and administrative tasks (e.g., note structuring)
4 - Legacy	encoder-decoder	T5, BART	low (pattern matching)	high	text normalization, simple QA

300
 301 **The Reasoning-Privacy Conflict.** Our review identifies a divergence in generator selection
 302 driven by the privacy vs. reasoning axis of the biomedical trilemma. While commercial APIs
 303 (GPT-4) offer superior reasoning for complex differentials, they pose data residency risks.
 304 Conversely, open-source models (Llama 3, Mistral) allow for secure, on-premise deployment
 305 but often lack the "System 2" reasoning depth required in complex cases (e.g., rare disease
 306 diagnosis) without extensive fine-tuning.

307
 308 **5.4 Knowledge Graph-based RAG**

309 While vector retrieval excels at capturing implicit semantic relationships (e.g., synonyms), it
 310 struggles with explicit multi-hop reasoning (e.g., drug A treats disease B, which is caused by
 311 gene C). KG retrieval resolves this by traversing structured edges [99]. On the negative side, it
 312 introduces a high maintenance cost. Unlike vector stores, which ingest unstructured text
 313 instantly, KGs require ontology alignment, making them brittle in the face of rapidly emerging
 314 medical concepts (e.g., COVID-19 and its subsequent variants).

315 Table 6 summarizes recent RAG implementations (2024–2025) using KG. We observe a
 316 divergence in application logic: while early systems focused on static fact retrieval (QA) [69,

317 97], recent frameworks leverage KGs for predictive reasoning (e.g., diagnosis prediction) [101]
 318 and hypothesis generation [102].

319

320 **Table 6:** Recent biomedical RAG studies based on KGs.

Method	Task	KG Used	Year
KG-RAG [72]	biomedical text generation, medical question answering	SPOKE biomedical KG	2024
HEALIE [103]	personalized medical content generation	HEALIE KG (customized)	2024
KRAGEN [104]	Biomedical problem solving	Alzheimer’s KG (AlzKB)	2024
Ascle [105]	various medical text generation tasks	UMLS	2024
DALK [106]	Alzheimer’s QA	Alzheimer’s KG (AlzKB)	2024
Gilbert <i>et al.</i> [107]	medical information structuring and interlinking	Not specified	2024
NetMe 2.0 [100]	biomedical knowledge extraction	BKGs constructed using OntoTagMe and Wikidata	2024
NEKO [108]	knowledge mining in synthetic biology	PubMed-derived KGs	2025
DR.KNOWS [101]	diagnosis prediction from EHRs	UMLS	2025
ESCARGOT [102]	biomedical reasoning and knowledge retrieval	Alzheimer’s KG (AlzKB)	2025

321

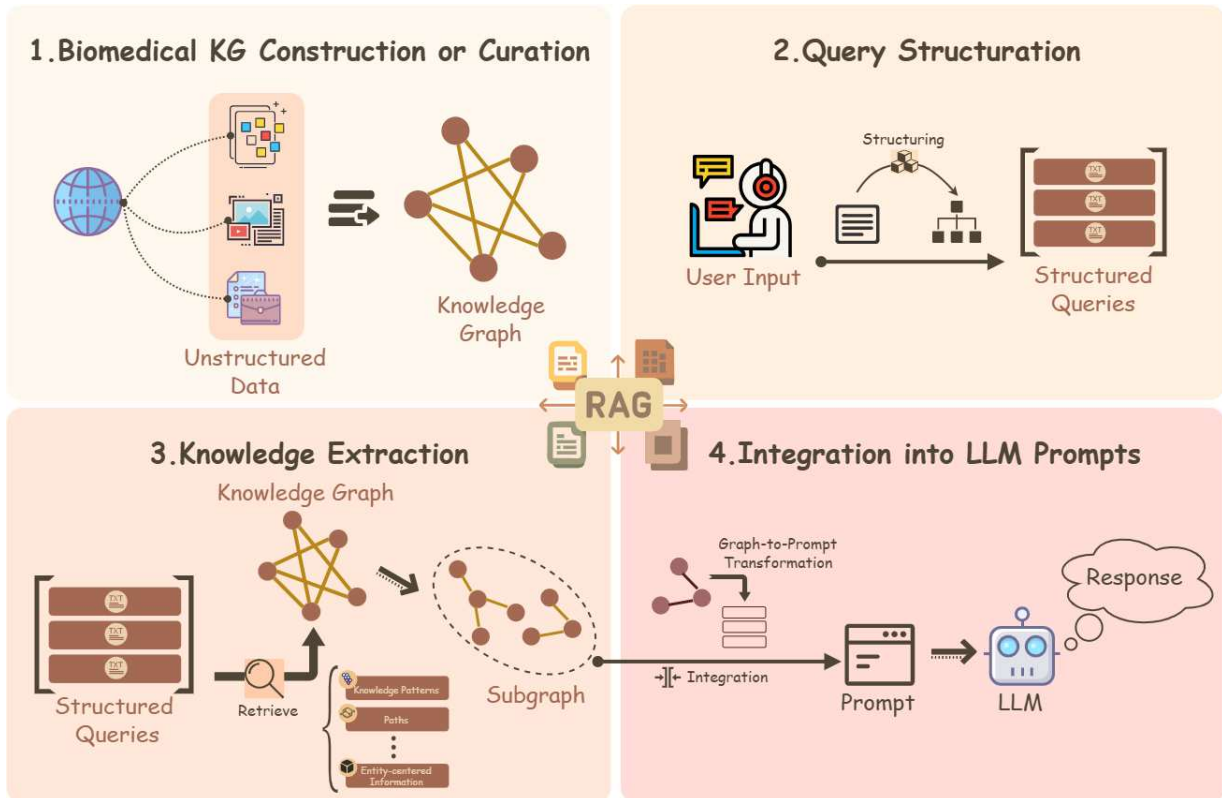
322 As shown in Figure 3, the integration of knowledge graphs with LLMs in biomedical RAG
 323 applications can be analyzed through four key components:

- 324 • **Biomedical KG construction:** We classify approaches into static authority vs. dynamic
 325 extraction. Systems like KG-RAG [72] utilize pre-curated, authoritative graphs (e.g.,
 326 SPOKE [109]) to ensure high precision but suffer from stale knowledge. In contrast, on-
 327 the-fly systems like NetMe 2.0 [100] and NEKO [108] construct dynamic subgraphs
 328 from retrieved literature using entity linking. While this resolves the knowledge currency
 329 issue, it introduces extraction noise, potentially polluting the reasoning chain with
 330 incorrect triples.
- 331 • **Query structuration:** This stage addresses the translation of natural language into
 332 formal query languages (e.g., Cypher [102]). A limitation here is the alignment gap:
 333 LLMs often generate syntactically correct but semantically invalid queries when mapped
 334 to rigid schemas like UMLS [101]. Frameworks like ESCARGOT attempt to mitigate

335 this via graph-of-thoughts (GoT) [102], while KRAGEN bypasses formal query
336 generation entirely by vectorizing the graph nodes for semantic search, trading query
337 expressivity for robustness [104].

- 338 • **Subgraph extraction:** Retrieving a full k -hop neighborhood around an entity often leads
339 to context explosion, overwhelming the context window of LLMs with irrelevant nodes.
340 To manage this signal-to-noise ratio, recent systems employ path-ranking algorithms [98,
341 100, 101]. For instance, DR.KNOWS [101] utilizes attention mechanisms to score and
342 prune paths, isolating only the diagnostic pathways relevant to the patient's EHR data
343 rather than the entire disease ontology.
- 344 • **Prompt integration:** The final challenge lies in graph linearization, i.e., converting
345 structured triples into natural language prompts without losing topological information.
346 While KG-RAG [72] implements token-optimized linearization to reduce costs, newer
347 approaches like NetMe 2.0 [100] integrate graph data as soft logic constraints, guiding
348 the LLM-generated content to remain faithful to the retrieved topology.

349



350
351 **Figure 3:** The framework of KG-based biomedical RAG: 1. Construction, 2. Structuration, 3.
352 Extraction, and 4. Integration.

353

354 **5.5 Multimodal RAG**

355 Multimodal biomedical RAG addresses the semantic gap between high-dimensional
 356 pixel/sensor data and low-dimensional clinical concepts [110]. Unlike general-domain models
 357 trained on billion-scale image-text pairs (e.g., CLIP [111]), biomedical systems operate in a data-
 358 scarce regime where aligned modalities (e.g., MRI + caption [109-113]) are expensive to curate.
 359 Consequently, the challenge shifts from simple retrieval to cross-modal alignment, requiring
 360 specialized encoders to map disparate inputs (e.g., histopathology patches [117]) into a shared
 361 semantic vector space compatible with text embeddings. As shown in Table 7, multimodal RAG
 362 systems have been investigated in several scenarios in biomedicine, mostly combining diverse
 363 medical imaging and textual data.

364

365 **Table 7:** Multimodal biomedical RAG systems.

Model	Task	Modalities	Year
Ranjit <i>et al.</i> [113]	chest X-Ray report generation	chest X-ray and text	2023
FactMM-RAG [112]	radiology report generation	X-Ray and text	2024
MMed-RAG [114]	medical VQA & report generation	medical images and text	2024
MEAG [118]	amblyopia diagnosis	eye tracking and text	2024
Tozuka <i>et al.</i> [119]	lung cancer staging	CT findings and TNM classifications	2024
Raminedi <i>et al.</i> [120]	radiology report generation	X-ray images and text	2024
RULE [115]	medical VQA & report generation	medical images and text	2024
Thetbanthad <i>et al.</i> [121]	prescription label identification	image (labels) and text	2025
Hu <i>et al.</i> [122]	pathology report generation	whole slide images and text	2025
STREAM [116]	chest X-ray report generation	chest X-ray and text	2025

366

367 Figure 4 illustrates the architectural workflow of a multimodal biomedical RAG system. The
 368 key distinction between multimodal biomedical RAG and traditional biomedical RAG lies in
 369 data input diversity and processing architecture. While traditional RAG systems primarily
 370 process textual queries to retrieve text-based medical knowledge, multimodal RAG incorporates
 371 and processes non-textual data (images, sensor readings, etc.) alongside text, requiring
 372 specialized encoders for each modality to meaningfully integrate these diverse inputs [110].

373

374 The basic framework of multimodal biomedical RAG typically consists of: (1) *modality-specific*
 375 *encoders* (e.g., Vision Transformers for images, specialized encoders for time-series data [123]);

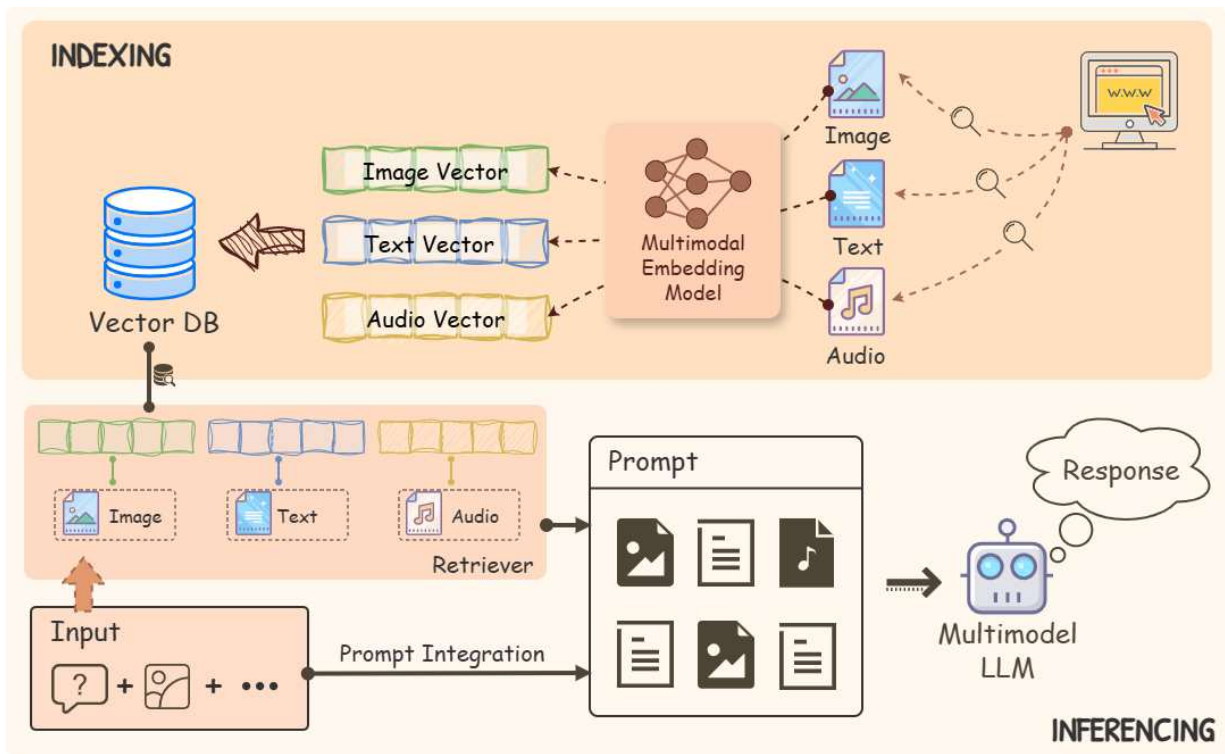
376 (2) *knowledge retrieval components* that access domain-specific databases based on multimodal
377 query representations; and (3) *generation modules* that synthesize comprehensive outputs
378 informed by both the multimodal inputs and retrieved knowledge.

379
380 We categorize encoding strategies into *symbolic conversion* vs. *neural alignment*. In symbolic
381 conversion, systems like Thetbanthad *et al.* [121] bypass alignment issues by converting images
382 to text via OCR (e.g., prescription labels) prior to retrieval. While robust, this method discards
383 visual texture information. Differently, neural alignment, which is the dominant paradigm
384 among the surveyed papers, employs contrastive learning to encode and align features from
385 different modalities. Frameworks like MMed-RAG [114] and Raminedi *et al.* [117] utilize
386 domain-adapted Vision Transformers (ViT) aligned via objectives similar to CLIP [111],
387 projecting visual features (e.g., X-rays) directly into the query space.

388
389 Retrieval mechanisms in this domain exhibit two distinct flows: visual-to-text and visual-to-
390 visual. Most systems (e.g., FactMM-RAG [112] and STREAM [116]) encode the patient's
391 medical image to retrieve textual precedents (e.g., radiology reports or guidelines). Conversely,
392 systems like Tozuka *et al.* [119] retrieve similar images (e.g., historical CT scans) to perform
393 case-based reasoning. This approach, however, struggles with intra-class variance, where, e.g.,
394 histologically identical tumors appear morphologically distinct across patients.

395
396 Finally, the generation phase faces the risk of modality hallucination, where the LLM describes
397 features present in its pre-training data but absent in the specific patient image. To mitigate this,
398 architectures like RULE [115] and Hu *et al.* [122] implement cross-attention fusion, allowing
399 the generator to attend dynamically to specific image patches (e.g., whole slide images) while
400 conditioning on retrieved text. This ensures that the generated report is fluent but also visually
401 grounded in the specific patient's pathology.

402



403

404

Figure 4: Architectural workflow of multimodal biomedical RAG system.

405

406 6 Datasets and Evaluation

407 The evaluation and development of biomedical RAG systems rely heavily on appropriate
408 datasets. In this section, we present a comprehensive overview of datasets commonly used in
409 biomedical RAG research, categorizing them into two main types: knowledge source datasets
410 and medical task datasets. This categorization reflects the dual nature of RAG systems, which
411 require both comprehensive knowledge bases and task-specific evaluation benchmarks.

412 6.1 Knowledge Source Datasets

413 The effectiveness of biomedical RAG systems considerably depends on the quality,
414 comprehensiveness, and diversity of their knowledge sources [1, 11]. These knowledge
415 repositories serve as the informational backbone from which these systems retrieve and
416 synthesize medical content. High-quality knowledge sources can improve the generated
417 responses, making them not only relevant but also reflect the current biomedical consensus and
418 best practices [13, 58, 121].

419 Biomedical knowledge spans multiple dimensions, going from basic science research to clinical
420 guidelines, from pharmaceutical data to standardized medical terminology. To address this
421 complexity, effective biomedical RAG systems integrate diverse knowledge sources that
422 complement each other in scope, specialization, and format. This integration enables systems to
423 comprehensively address the multifaceted nature of medical queries [31, 70].

424 However, we identify an important demographic and geographic bias. As shown in Table 8,
425 clinical RAG relies on MIMIC-IV. While invaluable, this dataset represents a single geographic
426 population (Boston, USA), potentially limiting the generalizability of RAG systems to
427 populations with different epidemiological profiles. Furthermore, the anglocentric dominance of
428 knowledge bases like PubMed and UMLS [125] creates a performance disparity for non-English
429 queries, as the knowledge-reasoning link breaks down when the retriever fails to find language-
430 aligned evidence. The more comprehensive sources of knowledge and relevant studies are
431 provided in Appendix Table A3.

432

433 **Table 8:** Biomedical knowledge sources.

Source	Type	Scale	Description	Studies
MIMIC-IV [126]	EHR	300K patients / 430K admissions	hospital admissions (2008-2019)	Moser et al. [127]
PubMed	literature	36M citations / 4.5B words	biomedical research citations and abstracts	Self-BioRAG [13], MEDRAG [7]

PMC	literature	8M full-text / 13.5B words	full-text biomedical and life sciences articles	Self-BioRAG [13], PodGPT [128]
GNQA [50]	literature	3,000 publications	peer-reviewed papers on aging, dementia, diabetes	GNQA [50]
UMLS [125]	knowledge base	2M entities / 900K concepts	integrated biomedical terminology	BiomedRAG [5], CliniqIR [30]
SPOKE [109]	KG	40+ databases	integrated biomedical knowledge sources	KG-RAG [72]
StatPearls [129]	clinical guide	9,330 articles	clinical decision support articles	MEDRAG [7]
Medical Textbooks [130]	educational	18 core textbooks	standard USMLE preparation texts	JMLR [3], i- MedRAG [131]
CheXpert [132]	image-report	224,316 radiographs	chest X-rays with reports	CLEAR [71]

434

435 6.2 Medical Task Datasets

436 Medical task datasets represent essential benchmarking instruments for evaluating RAG systems
437 within healthcare applications [2, 10]. These datasets simulate real clinical information
438 challenges, providing structured frameworks to assess how systems process, interpret, and
439 generate medical content. These evaluation datasets can be systematically categorized according
440 to distinct information-handling processes in the biomedical domain: biomedical information
441 extraction, entity recognition, QA, biomedical multiple-choice examination, dialogue, and
442 generation tasks. Table 9 presents some representative evaluation datasets employed across the
443 surveyed biomedical RAG systems, and a more detailed list of datasets is provided in Appendix
444 Table A4.

445 A pervasive challenge in evaluating biomedical RAG is test set contamination. Standard
446 benchmarks like MedQA and PubMedQA are frequently included in the pre-training corpora of
447 large foundation models (e.g., GPT-4, Llama 3). Consequently, it is often indistinguishable
448 whether a correct response stems from successful retrieval (i.e., RAG) or parametric
449 memorization. To rigorously evaluate the *added value* of retrieval, we recommend the adoption
450 of dynamic evaluation sets, such as private hospital cases or newly published guidelines (post-
451 training cutoff), where the LLM cannot rely on memorized knowledge.

452

453 **Table 9:** Biomedical tasks datasets.

Dataset	Task	Scale / Size	Studies
---------	------	--------------	---------

MedQA [130]	multiple-choice	61,097 questions	Self-BioRAG [13], JMLR [3]
PubMedQA [133]	multiple-choice	273.5k total: 1k expert-labeled, 61.2k unlabeled, 211.3k generated	BMRETRIEVER [32], PodGPT [128]
MIMIC-III [134]	text summarization	53,423 reports	DR.KNOWS [101], CliniqIR [30]
BioASQ [135]	info retrieval	annual challenge sets	BMRETRIEVER [32]
ChemProt corpus [136]	info extraction	2,432 abstracts	BiomedRAG [5]
MMLU [137]	multiple-choice	15,908 questions among 57 tasks	Self-BioRAG [13]

454

455 6.3 Evaluation Metrics

456 The biomedical RAG systems reviewed in this survey typically follow a three-stage pipeline:
 457 (1) retrieval, (2) reranking, and (3) generation. Consequently, evaluation methodologies must be
 458 aligned with these stages. While retrieval can be assessed using standard information retrieval
 459 (IR) metrics, the evaluation of generation in biomedicine presents specific challenges: standard
 460 NLP metrics (e.g., BLEU) often fail to capture clinical correctness (e.g., distinguishing
 461 *hypotension* from *hypertension* despite high lexical overlap). In what follows, we review the
 462 metrics used across the pipeline, highlighting the shift from n-gram matching to entity-based
 463 validation.

464

465 6.3.1 Metrics for Retrieval

466 The retrieval stage evaluates how effectively a system surfaces and orders clinically relevant
 467 documents. We review the core metrics:

- 468 • Precision@ k : Measures the proportion of relevant documents in the top- k results:

$$469 \text{Precision@}k = \frac{\sum_{i=1}^k 1[\text{doc}_i \text{ is relevant}]}{k}.$$

- 470 • Recall@ k : Quantifies the ability of the system to find all relevant cases:

$$471 \text{Recall@}k = \frac{\sum_{i=1}^k 1[\text{doc}_i \text{ is relevant}]}{[\text{Relevant Documents}]}.$$

472 In biomedical contexts, recall is often prioritized over precision, as missing a relevant
 473 guideline or patient history (false negative) carries higher safety risks than retrieving an
 474 irrelevant one (false positive). However, recent results show that providing contradictory (or

475 distractor) examples in the context has a significant negative impact on the quality of the
476 generation [138].

- 477 • MAP@ k (mean average precision): Averages precision across varying recall levels,
478 rewarding systems that place relevant documents higher in the rank.
- 479 • nDCG@ k (normalized discounted cumulative gain): Accounts for graded relevance (e.g.,
480 highly relevant vs. partially relevant), which is important, for example, when
481 distinguishing between specific clinical protocols and general medical advice.

482 6.3.2 Metrics for Reranking

483 Reranking is evaluated using the same IR metrics (Precision, MAP, nDCG) but with stricter
484 thresholds (e.g., $k=5$ instead of $k=1000$). In RAG applications, reranking utility is measured by
485 the delta in metrics, such as nDCG@ k , before and after the reranker, quantifying the ability of
486 the filter to discard distractor documents that share keywords but lack semantic alignment with
487 the clinical query.

488 6.3.3 Metrics for Generation

489 Generation metrics assess the synthesized response. We classify them into surface-level and
490 semantic-level metrics, highlighting a divergence in clinical utility.

491 N-gram Metrics - Surface Level

- 492 • ROUGE-N/L [139]: Measures n-gram overlap and longest common subsequences.
493 Widely used for summarization. ROUGE-N is given by the following equation:

$$494 \text{ROUGE-N} = \frac{\text{Overlap of N-grams}}{\text{Total N-grams in reference}}$$

495 ROUGE-L considers the length of the reference instead of the denominator and the longest
496 overlap for the numerator.

- 497 • BLEU [140]: Evaluates the precision of n-grams:

$$498 \text{BLEU} = \text{BP} \cdot \exp\left(\sum_{n=1}^N w_n \log p_n\right).$$

499 While scalable, these metrics are increasingly viewed as insufficient for biomedical RAG. A
500 model can achieve a high BLEU score by copying medical jargon while hallucinating a negation
501 (e.g., generating *the patient has pneumonia* vs. *the patient has "no" pneumonia*). Consequently,
502 reliance on BLEU/ROUGE alone is a negative signal for clinical validity.

503 Embedding-based Metrics - Semantic Level

- 504 • BERTScore [141]: Evaluates similarity by aligning token embeddings from a pre-trained
505 model, capturing synonyms that n-gram metrics miss. For a generated text X and a
506 reference text Y , the precision P and recall R are defined as

507
$$P = \frac{1}{|X|} \sum_{x \in X} \max_{y \in Y} \cos(e_x, e_y), R = \frac{1}{|Y|} \sum_{y \in Y} \max_{x \in X} \cos(e_y, e_x),$$

508 where e_x and e_y denote contextual embeddings of tokens $x \in X$ and $y \in Y$, and $|X|$ and $|Y|$
509 denote the number of generated and reference tokens, respectively. Using P and R, the F₁-
510 score can then be computed as the harmonic mean between these two metrics.

- 511 • BLEURT [142]: A learned metric fine-tuned on human ratings to predict fluency and
512 adequacy:

513
$$BLUERT(X, Y) = f_{\theta}(enc(X), enc(Y)),$$

514 where $enc(\cdot)$ denotes contextual embeddings, f_{θ} is a regression head trained to approximate
515 human judgments, and X and Y are the generated and reference texts, respectively.

516 **Entity-based Metrics - Clinical Level**

- 517 • RaTEScore [143]: To address the limitations of n-gram matching and standard semantic
518 similarity scores, recent work has introduced entity-centric metrics like RaTEScore. This
519 metric extracts medical entities (e.g., diseases, anatomies) and relations from both the
520 generated and reference text, computing an F₁-score based on fact matching rather than
521 word matching. This represents the state-of-the-art for validating radiology reports and
522 clinical notes.

523 **6.3.4 Task-Specific Metrics**

- 524 • Classification: For tasks like diagnostic prediction, standard metrics apply: accuracy,
525 macro-F1 (to prevent majority-class bias), and AUROC (for probability-based risk
526 scoring).
- 527 • Named entity recognition (NER): Evaluated via entity-level F₁-score, which requires
528 exact boundary matching of clinical concepts (e.g., detecting *type 2 diabetes* as a single
529 entity rather than just *diabetes*).
- 530 • Dialogue: Clinical dialogue is assessed via a hybrid of *task completion rate* [144] (did
531 the system gather the necessary symptoms?) and *safety scores* [142, 143] (did the system
532 avoid recommending dangerous contraindications?), often requiring LLM-based judges
533 or human expert review.

534

535 **7 Applications**

536 The application of RAG in the biomedical domain reflects an emerging direction in healthcare
 537 informatics, combining the strengths of LLMs with domain-specific knowledge retrieval. This
 538 integration helps mitigate challenges inherent to healthcare applications—factual accuracy,
 539 domain specificity, knowledge recency, and explainability—that conventional generative AI
 540 approaches struggle to overcome [147]. In this section, we examine the healthcare domains
 541 where biomedical RAG systems have been utilized and organize the discussion around their
 542 major application areas, including clinical decision support, clinical report generation, precision
 543 medicine, medical education, and clinical research. Building on these observations, we provide
 544 corresponding recommendations for selecting appropriate RAG architectures in Table 10.

545

546 **Table 10:** Recommended RAG architecture for healthcare applications.

Application Domain	Recommended Architecture	Primary Constraint (Trilemma)	Rationale
Patient education	naive RAG	latency	Speed is necessary for user engagement; queries are usually lexical (e.g., "what is diabetes") requiring low reasoning depth.
Real-time CDSS	advanced RAG	precision	Requires high-recall retrieval + reranking to filter irrelevant guidelines; moderate latency (2-5s) is acceptable for physician verification.
Rare disease diagnosis	modular (Agentic)	reasoning depth	Requires multi-step logic (symptoms → phenotype → gene); high latency is acceptable for asynchronous diagnostics.
Clinical coding	hybrid (sparse+dense)	precision	Must handle exact alphanumeric matching (sparse) and concept mapping (dense); hallucination is not tolerated.
Evidence synthesis	modular	completeness	Requires iterative "search → summarize → search" loops to cover vast literature (PubMed); latency is irrelevant.

547

548 **7.1 Clinical Decision Support Systems**

549 **7.1.1 Medical Question Answering**

550 Medical QA systems provide clinicians and patients with access to accurate healthcare
551 information at the point of need [72]. Unlike standard LLMs, RAG-based methodologies allow
552 patients and clinicians to verify answers against source documents. We observe a functional
553 divergence in QA systems based on their target user: patient or physician. In *patient-facing*
554 *simplification*, systems like RALL [148] prioritize health literacy, using retrieval to translate
555 complex medical jargon into lay language. The utility here is measured by accuracy, but also by
556 readability scores (e.g., Coleman-Liau [149]). Differently, in *clinician-facing verification*, the
557 priority is the source of evidence. Systems like Clinfo.ai [35] and Self-BioRAG [13] implement
558 citation-checking loops. While Self-BioRAG achieves state-of-the-art accuracy via a self-
559 reflection token (<CRITIQUE>), the iterative verification process triples the inference latency
560 compared to single-pass systems like WeiseEule [73], potentially limiting its use in time-
561 sensitive care (such as acute care).

562 **7.1.2 Diagnostic and Treatment Decision Support**

563 Diagnostic and treatment decision support systems assist clinicians in navigating complex
564 clinical pathways through evidence retrieval and synthesis [150]. They can help healthcare
565 providers identify relevant evidence, recognize patterns, and identify appropriate interventions,
566 ultimately improving diagnostic accuracy and treatment selection while reducing cognitive
567 biases and diagnostic errors. RAG enhances diagnostic workflows by bridging distinct data silos:
568 *clinical guidelines* and *EHR data*. For the former, systems like SurgeryLLM [151] and Expert-
569 Guided LLMs [94] retrieve static protocols (e.g., *ASCO*, *European Society for Medical*
570 *Oncology*, and *National Comprehensive Cancer Network* guidelines) to standardize care. The
571 primary utility here is compliance checking, ensuring treatment plans adhere to the latest
572 standards. For the latter, approaches, such as DR.KNOWS [101] and RECTIFIER [63], operate
573 on dynamic patient data. By mapping unstructured clinical notes to structured KGs (e.g.,
574 UMLS), these systems identify eligibility for clinical trials or predict diagnoses. A key insight is
575 that KG retrieval (e.g., in DR.KNOWS) outperforms vector search in capturing the complex,
576 multi-hop relationships required for differential diagnosis, albeit with higher implementation
577 complexity.

578 **7.1.3 Rare Disease Identification and Management**

579 Rare diseases represent the long tail of medicine, where standard LLMs might hallucinate due
580 to sparse training data [152]. RAG systems address this by accessing specialized external
581 databases (e.g., Orphanet) that were not present in the pre-training corpus of the model. Surveyed

582 systems focus on phenotype-genotype mapping and diagnostic recall. RAG-HPO [45] and
583 RDguru [153] demonstrate that retrieving Human Phenotype Ontology (HPO) [154] terms
584 significantly improves the precision of gene prioritization compared to standard Llama and
585 GPT-4. Particularly, in RAG-HPO F1 score for suggesting HPO codes varies from 0.10 to 0.64
586 when RAG is incorporated. Differently, for diagnostic recall, the essential metric is recall, i.e.,
587 finding the correct disease among thousands, rather than precision. Zelin *et al.* [155] show that
588 augmenting LLMs with specialized database information connected with uncommon symptom
589 patterns increases diagnostic recall for cases that physicians might miss.

Comparative Case Study: RareDxGPT vs. ChatGPT in Rare Disease Diagnosis

To evaluate the impact of external knowledge on diagnostic precision, we contrast the performance of standard ChatGPT 3.5 (parametric) with RareDxGPT (RAG-augmented) across two representative failure modes available in [152]:

1. RAG success: *Sweet Syndrome* - A skin condition requiring specific dermatological knowledge.
 - Result: RareDxGPT correctly diagnosed *Sweet Syndrome* across multiple prompting strategies. In contrast, ChatGPT 3.5 consistently misdiagnosed the case as *Gianotti-Crosti Syndrome*.
 - Insight: As noted in the study, RareDxGPT was successful in diagnosing skin conditions because the retrieved documents contained sufficient phenotypic descriptions. Having this additional external information allowed the model to correctly identify the condition, whereas the standard model (ChatGPT 3.5) misidentified it as a different dermatological syndrome (*Gianotti-Crosti*).
2. RAG failure: *Marfan Syndrome* - A relatively common rare disease with distinct features.
 - Result: ChatGPT 3.5 correctly diagnosed *Marfan Syndrome* using a standard prompt. RareDxGPT failed consistently, diagnosing *Gorlin-Chaudhry-Moss Syndrome* across all prompts.
 - Insight: This failure reveals a RAG vulnerability. The retrieval database lacked detailed phenotypes for *Marfan* but contained them for *Gorlin-Chaudhry-Moss*. To justify the incorrect retrieval, RareDxGPT fabricated a non-existent symptom (*craniosynostosis*) that was not present in the patient's case but was required for the *Gorlin* diagnosis. This proves that when retrieval quality is low, RAG can actively corrupt the reasoning process by hallucinating evidence to fit the retrieved context.

590
591
592 **7.2 Clinical Report Generation**
593 Automatic and semi-automatic report generation can help healthcare providers by significantly
594 reducing the time they spend on documentation, allowing them to dedicate more time to patient

595 care [156]. Surveyed paper for report generation focused mostly is radiology, which involves
596 image analysis, document consultation, and data evaluation [157]. This process faces multiple
597 challenges, including ensuring factual accuracy, maintaining clinical relevance, and providing
598 sufficient interpretability for clinical adoption [158]. To mitigate this, we identify two dominant
599 architectural strategies: *retrieval-based grounding* and *concept-guided generation*. In retrieval-
600 based grounding, approaches like Fact-Aware RAG [112] retrieve similar historical reports to
601 use as templates. While this ensures stylistic consistency, it risks leaking details from past
602 patients into the current report. On the other hand, in concept-guided generation, systems like
603 LaB-RAG [159] and multi-agent frameworks [160] first extract structured concepts (e.g., *pleural*
604 *effusion: absent*) and then generate text conditioned on these facts. This symbolic constraint
605 significantly reduces hallucinations compared to end-to-end generation, offering a safety rail for
606 clinical adoption.

607 **7.3 Precision Medicine Applications**

608 **7.3.1 Genomic Medicine**

609 In genomic medicine, the bottleneck is variant interpretation, i.e., classifying a genetic mutation
610 as pathogenic or benign based on shifting literature. Systems like FAVOR-GPT [59] and Lu *et*
611 *al.* [37] utilize RAG to retrieve real-time annotations from databases like ClinVar [161]. This
612 highlights the temporal advantage of RAG-based methodologies: while a fine-tuned model
613 becomes obsolete the moment a new variant is discovered, a RAG system remains current simply
614 by updating its vector index.

615 **7.3.2 Personalized Health Management**

616 For patient management, RAG shifts the focus from generic advice to context-aware CDSS.
617 Systems like RISE [58] (diabetes) and HEALIE [103] retrieve patient-specific history to tailor
618 recommendations. However, we note a privacy trade-off, as this effective personalization
619 requires indexing highly sensitive PHI, necessitating local deployment strategies (see Section
620 5.1).

621 **7.4 Healthcare Education**

622 In medical education, RAG serves two distinct pedagogical functions. Systems like EyeGPT
623 [162] and Zhao *et al.* [229] retrieve from standard textbooks to ensure answers align with board
624 exam criteria, preventing the model from offering correct but non-standard advice. Differently,
625 approaches like EyeTeacher [163] use retrieval to generate diverse synthetic patient cases. This
626 allows for unlimited case generation, providing students with exposure to rare pathologies that
627 they might not encounter during clinical rotations.

628 **7.5 Clinical Research**

629 Beyond patient care, RAG accelerates clinical research by automating evidence synthesis. Tools
630 like LITURAt [39] and PodGPT [128] allow researchers to query vast literature databases
631 (PubMed) to identify gaps or contradictions in current studies. This represents a shift from search
632 (finding papers) to synthesis (generating systematic review drafts), potentially reducing the time
633 required for meta-analyses from months to days.

634

635 **8 Discussion and Future Directions**

636 **8.1 Architectural Synthesis: The Trilemma**

637 Synthesizing the trade-offs identified in the retrieval (Section 5.1), reranking (Section 5.2), and
638 generation (Section 5.3) sections, we formally define the *biomedical RAG trilemma*. This
639 framework explains why no single architecture currently dominates the clinical landscape.

640 **Latency vs. Reasoning.** While modular architectures achieve the highest diagnostic accuracy
641 by iteratively reasoning over multiple documents, they introduce significant latency penalties.
642 Systems, such as Self-BioRAG [13] and GeneGPT [15], utilize routing policies to dynamically
643 trigger tool use or self-reflection, enabling *System 2* deliberative thought. However, this iterative
644 process often results in response times (>10s) [17] that disqualify them from real-time clinical
645 workflows, creating a dichotomy where the most intelligent systems are often too slow for the
646 point of care compared to single-pass naive RAG implementations.

647 **Privacy vs. Capability.** A clear divergence exists in LLM selection for the generative step. The
648 most capable generation models, particularly commercial APIs like GPT-4/5 [164] and Gemini
649 2.5/3 [89], dominate benchmarks but pose high data residency risks. This forces healthcare
650 institutions to choose between cloud-based intelligence (high performance, low privacy) and on-
651 premise compliance using open-source models, such as Llama 3 [165], Mistral [90] or Qwen
652 [166]. While domain-adapted models like MedAlpaca [167] attempt to bridge this gap, recent
653 studies indicate they do not consistently outperform generalist open-source models in RAG
654 settings [168], which is consistent with other recent evaluations in the literature [169]. This
655 finding suggests that reasoning capability (derived from scale) is often more important than
656 domain-specific vocabulary storage.

657 **Precision vs. Recall.** In the retrieval layer, we observe that no single method suffices. Sparse
658 retrieval (BM25) remains essential for the exact entity matching required for specific biomedical
659 terms, such as alphanumeric codes (*e.g.*, ICD-10) and drug dosages. Conversely, dense retrievers
660 like MedCPT [12] are necessary to capture semantic symptomatology that lacks keyword
661 overlap. This necessitates the hybrid architectures dominant in recent literature (*e.g.*, MEDRAG
662 [7], CliniqIR [30]), which incur higher indexing complexity to ensure patient safety by
663 combining lexical precision with semantic recall. In this sense, commercial solutions, such as
664 ElasticSearch [170], are moving towards hybrid retrieval engines.

665 **8.2 Challenges and Future Directions**

666 **Trustworthiness and citation hallucination.** A failure mode unique to RAG is citation
667 hallucination, where systems generate factually correct medical statements but attribute them to

668 irrelevant references. This illusion of verification poses severe safety risks. To mitigate this,
669 future architectures must implement citation verification modules similar to Self-BioRAG’s self-
670 reflection mechanism [13], which penalizes generation-retrieval misalignment. Furthermore,
671 robustness against misinformation remains a priority, as malicious content in web corpora can
672 contaminate knowledge sources, necessitating stricter curation of retrieval indices.

673 **Multimodality and the alignment gap.** In biomedicine, extending RAG beyond text faces the
674 alignment gap between high-dimensional pixel data and semantic concepts. While approaches,
675 such as MMed-RAG [114] and RULE [115], utilize contrastive learning to bridge this gap, they
676 struggle with high intra-class variance, where histologically identical tumors appear
677 morphologically distinct across patients. Future research must move beyond simple image-text
678 matching to Fact-Aware Multimodal Retrieval, ensuring that generated reports are grounded in
679 specific visual features rather than generic templates.

680 **Restricted resources and green AI.** The resource constraint in hospitals is not merely financial
681 but operational. State-of-the-art models like DeepSeek V3 [171] require massive high-
682 performance computing clusters that are absent in most clinical IT infrastructures. Future
683 research must prioritize small language models (SLMs) [172] (<7B parameters) optimized for
684 RAG, such as Phi-4 [173] or distilled models like QwQ [166]. Proving that a small reasoner with
685 access to a massive external memory can rival larger parametric models is essential for
686 democratizing AI in low-resource healthcare settings.

687 **Privacy and federated deployment.** Privacy concerns persist throughout the retrieval and
688 generation phases, particularly regarding the handling of PHI under frameworks like GDPR and
689 HIPAA. To address this, the field must move toward federated RAG and local deployment
690 strategies. By utilizing efficient vector stores (e.g., Faiss [79], Chroma [80]) within secure
691 hospital firewalls and leveraging portable LLMs, institutions can ensure that sensitive patient
692 data never traverses external networks, resolving the conflict between utility and confidentiality.

693

694 **9 Conclusions**

695 In conclusion, this survey advances the understanding of biomedical RAG by moving beyond a
696 technological inventory to a synthesis of architectural trade-offs. We have formalized the
697 evolution from static naive RAG to dynamic modular paradigms, identifying that while agentic
698 workflows maximize diagnostic reasoning, they introduce latency bottlenecks, which might be
699 prohibitive for real-time care. Our analysis of the biomedical RAG trilemma highlights the
700 tension between reasoning depth, inference speed, and data privacy, forcing a strategic choice
701 between high-performance cloud architectures and privacy-preserving on-premise deployments.

702 Furthermore, we critically evaluated the reasoning-knowledge decoupling, demonstrating that
703 generalist models often surpass domain-specific ones when grounded by effective retrieval.
704 Finally, by auditing dataset validity, we identified risks regarding test-set contamination and the
705 alignment gap in multimodal systems. Ultimately, we project that the field will evolve from
706 static retrieval-augmented generation to dynamic retrieval-augmented reasoning, where
707 autonomous agents actively navigate medical knowledge graphs and clinical guidelines to
708 function not just as search engines, but as verifiable clinical partners.

709

710 **10 Funding Declarations**

711 The authors declare that no external funding was received to conduct this study.

712

713 **References**

- 714 [1] W. Fan *et al.*, “A Survey on RAG Meeting LLMs: Towards Retrieval-Augmented Large Language Models,”
715 in *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '24)*,
716 New York, NY, USA: Association for Computing Machinery, 2024, pp. 6491–6501. doi:
717 10.1145/3637528.3671470.
- 718 [2] K. He *et al.*, “A Survey of Large Language Models for Healthcare: from Data, Technology, and Applications
719 to Accountability and Ethics,” *ArXiv Prepr. ArXiv231005694*, 2023, [Online]. Available:
720 <https://arxiv.org/abs/2310.05694>
- 721 [3] J. Wang, Z. Yang, Z. Yao, and others, “JMLR: Joint Medical LLM and Retrieval Training for Enhancing
722 Reasoning and Professional Question Answering Capability,” *ArXiv Prepr. ArXiv240217887*, 2024, [Online].
723 Available: <https://arxiv.org/abs/2402.17887>
- 724 [4] H. Xiao *et al.*, “A comprehensive survey of large language models and multimodal large language models in
725 medicine,” *Inf. Fusion*, vol. 117, p. 102888, 2025, doi: <https://doi.org/10.1016/j.inffus.2024.102888>.
- 726 [5] M. Li, H. Kilicoglu, H. Xu, and R. Zhang, “BiomedRAG: A retrieval augmented large language model for
727 biomedicine,” *J. Biomed. Inform.*, vol. 162, p. 104769, Feb. 2025, doi: 10.1016/j.jbi.2024.104769.
- 728 [6] G. Frisoni, A. Cocchieri, A. Presepi, G. Moro, and Z. Meng, “To Generate or to Retrieve? On the Effectiveness
729 of Artificial Contexts for Medical Open-Domain Question Answering,” in *Proceedings of the 62nd Annual
730 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Bangkok, Thailand:
731 Association for Computational Linguistics, 2024, pp. 9878–9919.
- 732 [7] G. Xiong, Q. Jin, Z. Lu, and A. Zhang, “Benchmarking retrieval-augmented generation for medicine,” in
733 *Findings of the Association for Computational Linguistics ACL 2024*, 2024, pp. 6233–6251.
- 734 [8] Y. Li, C. Li, Z. Wang, D. Sui, and J. Yan, “A Privacy-Preserving Framework for Medical Chatbot Based on
735 LLM with Retrieval Augmented Generation,” Berlin, Heidelberg: Springer-Verlag, 2024, pp. 15–28. doi:
736 10.1007/978-981-97-9437-9_2.
- 737 [9] S. Santra, P. Kukreja, K. Saxena, S. Gandhi, and O. V. Singh, “Navigating regulatory and policy challenges for
738 AI enabled combination devices,” *Front. Med. Technol.*, vol. 6, p. 1473350, 2024, doi:
739 10.3389/fmedt.2024.1473350.
- 740 [10] L. Liu *et al.*, “A Survey on Medical Large Language Models: Technology, Application, Trustworthiness, and
741 Future Directions,” *ArXiv Prepr. ArXiv240603712*, 2024, doi: 10.48550/arXiv.2406.03712.
- 742 [11] Y. Gao, Y. Xiong, X. Gao, and others, “Retrieval-augmented generation for large language models: A survey,”
743 *ArXiv Prepr.*, vol. arXiv:2312.10997, no. 2, 2023, [Online]. Available: <https://arxiv.org/abs/2312.10997>
- 744 [12] Q. Jin *et al.*, “MedCPT: Contrastive Pre-trained Transformers with large-scale PubMed search logs for zero-
745 shot biomedical information retrieval,” *Bioinformatics*, vol. 39, no. 11, p. btad651, 2023.
- 746 [13] M. Jeong, J. Sohn, M. Sung, and J. Kang, “Improving medical reasoning through retrieval and self-reflection
747 with retrieval-augmented large language models,” *Bioinformatics*, vol. 40, no. Suppl 1, pp. i119–i129, June
748 2024, doi: 10.1093/bioinformatics/btae238.
- 749 [14] J. Song, C. Jin, W. Zhao, A. McCallum, and J.-Y. Lee, “Comparing Neighbors Together Makes it Easy: Jointly
750 Comparing Multiple Candidates for Efficient and Effective Retrieval,” in *Proceedings of the 2024 Conference
751 on Empirical Methods in Natural Language Processing*, Y. Al-Onaizan, M. Bansal, and Y.-N. Chen, Eds.,
752 Miami, Florida, USA: Association for Computational Linguistics, Nov. 2024, pp. 22255–22269. doi:
753 10.18653/v1/2024.emnlp-main.1242.
- 754 [15] Q. Jin, Y. Yang, Q. Chen, and Z. Lu, “GeneGPT: Augmenting Large Language Models with Domain Tools for
755 Improved Access to Biomedical Information,” *Bioinformatics*, vol. 40, no. 2, p. btae075, 2024, doi:
756 10.1093/bioinformatics/btae075.
- 757 [16] S. Ray *et al.*, “METIS: Fast Quality-Aware RAG Systems with Configuration Adaptation,” in *Proceedings of
758 the ACM SIGOPS 31st Symposium on Operating Systems Principles*, in SOSP '25. New York, NY, USA:
759 Association for Computing Machinery, 2025, pp. 606–622. doi: 10.1145/3731569.3764855.
- 760 [17] C.-Y. Lin *et al.*, “TeleRAG: Efficient Retrieval-Augmented Generation Inference with Lookahead Retrieval,”
761 Nov. 11, 2025, *arXiv*: arXiv:2502.20969. doi: 10.48550/arXiv.2502.20969.
- 762 [18] S. Robertson and H. Zaragoza, “The Probabilistic Relevance Framework: BM25 and Beyond,” *Found Trends
763 Inf Retr*, vol. 3, no. 4, pp. 333–389, Apr. 2009, doi: 10.1561/1500000019.
- 764 [19] K. S. Jones, “A statistical interpretation of term specificity and its application in retrieval,” *J. Doc.*, vol. 28, no.
765 1, pp. 11–21, 1972, doi: 10.1108/eb026526.
- 766 [20] S. Das *et al.*, “Two-Layer Retrieval-Augmented Generation Framework for Low-Resource Medical Question
767 Answering Using Reddit Data: Proof-of-Concept Study,” *J. Med. Internet Res.*, vol. 27, p. e66220, Jan. 2025,
768 doi: 10.2196/66220.
- 769 [21] E. J. Gong *et al.*, “The Potential Clinical Utility of the Customized Large Language Model in Gastroenterology:
770 A Pilot Study,” *Bioeng. Basel*, vol. 12, no. 1, p. 1, Dec. 2024, doi: 10.3390/bioengineering12010001.

- 771 [22] Y. Santander-Cruz *et al.*, “Semantic Feature Extraction Using SBERT for Dementia Detection,” *Brain Sci.*,
772 vol. 12, no. 2, p. 270, 2022, doi: 10.3390/brainsci12020270.
- 773 [23] N. Thakur, N. Reimers, J. Daxenberger, and I. Gurevych, “Augmented SBERT: Data Augmentation Method
774 for Improving Bi-Encoders for Pairwise Sentence Scoring Tasks,” *ArXiv Prepr. ArXiv201008240*, 2020,
775 [Online]. Available: <https://arxiv.org/abs/2010.08240>
- 776 [24] S. Xiao, Z. Liu, P. Zhang, and N. Muennighoff, “C-Pack: Packaged Resources To Advance General Chinese
777 Embedding.” 2023.
- 778 [25] V. Karpukhin *et al.*, “Dense Passage Retrieval for Open-Domain Question Answering,” in *Proceedings of the*
779 *2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 6769–6781.
- 780 [26] P. BehnamGhader, V. Adlakha, M. Mosbach, D. Bahdanau, N. Chapados, and S. Reddy, “LLM2Vec: Large
781 Language Models Are Secretly Powerful Text Encoders,” in *First Conference on Language Modeling*, 2024.
782 [Online]. Available: <https://openreview.net/forum?id=IW1PR7vEBf>
- 783 [27] C. Wu, W. Lin, X. Zhang, Y. Zhang, Y. Wang, and W. Xie, “PMC-LLaMA: Towards Building Open-source
784 Language Models for Medicine,” *ArXiv Prepr. ArXiv230414454*, 2023.
- 785 [28] J. Lee, H. Cha, Y. Hwangbo, and W. Cheon, “Enhancing Large Language Model Reliability: Minimizing
786 Hallucinations with Dual Retrieval-Augmented Generation Based on the Latest Diabetes Guidelines,” *J. Pers.*
787 *Med.*, vol. 14, no. 12, p. 1131, Nov. 2024, doi: 10.3390/jpm14121131.
- 788 [29] I. Alonso, M. Oronoz, and R. Agerri, “MedExpQA: Multilingual benchmarking of Large Language Models for
789 Medical Question Answering,” *Artif. Intell. Med.*, vol. 155, p. 102938, 2024, doi:
790 <https://doi.org/10.1016/j.artmed.2024.102938>.
- 791 [30] T. Abdullahi, L. Mercurio, R. Singh, and C. Eickhoff, “Retrieval-Based Diagnostic Decision Support: Mixed
792 Methods Study,” *JMIR Med. Inform.*, vol. 12, p. e50209, June 2024, doi: 10.2196/50209.
- 793 [31] S. Kim and J. Yoon, “VAIV bio-discovery service using transformer model and retrieval augmented
794 generation,” *BMC Bioinformatics*, vol. 25, no. 1, p. 273, Aug. 2024, doi: 10.1186/s12859-024-05903-6.
- 795 [32] R. Xu *et al.*, “BMRetriever: Tuning Large Language Models as Better Biomedical Text Retrievers,” in
796 *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Miami, Florida,
797 USA: Association for Computational Linguistics, 2024, pp. 22234–22254.
- 798 [33] M. Hu *et al.*, “SeRTS: Self-Rewarding Tree Search for Biomedical Retrieval-Augmented Generation,” in
799 *Findings of the Association for Computational Linguistics: EMNLP 2024*, Miami, Florida, USA: Association
800 for Computational Linguistics, 2024, pp. 1321–1335.
- 801 [34] G. Izacard *et al.*, “Unsupervised Dense Information Retrieval with Contrastive Learning.” 2021. doi:
802 10.48550/ARXIV.2112.09118.
- 803 [35] A. Lozano, S. L. Fleming, C. C. Chiang, and N. Shah, “Clinfo.ai: An Open-Source Retrieval-Augmented Large
804 Language Model System for Answering Medical Questions using Scientific Literature,” *Pac. Symp.*
805 *Biocomput.*, vol. 29, pp. 8–23, 2024.
- 806 [36] N. C. for B. Information (NCBI), “Entrez Programming Utilities (E-utilities) API.” 2024. [Online]. Available:
807 <https://www.ncbi.nlm.nih.gov/home/develop/api/>
- 808 [37] S. Lu and E. Cosgun, “Boosting GPT models for genomics analysis: generating trusted genetic variant
809 annotations and interpretations through RAG and Fine-tuning,” *Bioinforma. Adv.*, vol. 5, no. 1, p. vbaf019,
810 Feb. 2025, doi: 10.1093/bioadv/vbaf019.
- 811 [38] “Whoosh 2.7.4 documentation — Whoosh 2.7.4 documentation.” Accessed: Dec. 08, 2025. [Online].
812 Available: <https://whoosh.readthedocs.io/en/latest/>
- 813 [39] D. Peasley, R. Kuplicki, S. Sen, and M. Paulus, “Leveraging Large Language Models and Agent-Based
814 Systems for Scientific Data Analysis: Validation Study,” *JMIR Ment. Health*, vol. 12, p. e68135, Feb. 2025,
815 doi: 10.2196/68135.
- 816 [40] W. X. Zhao, J. Liu, R. Ren, and J.-R. Wen, “Dense Text Retrieval Based on Pretrained Language Models: A
817 Survey,” *ACM Trans Inf Syst*, vol. 42, no. 4, Feb. 2024, doi: 10.1145/3637870.
- 818 [41] A. V. Solatorio, “GISTEmbed: Guided In-sample Selection of Training Negatives for Text Embedding Fine-
819 tuning,” *ArXiv Prepr. ArXiv240216829*, 2024, [Online]. Available: <https://arxiv.org/abs/2402.16829>
- 820 [42] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, “Enriching Word Vectors with Subword Information,”
821 *Trans. Assoc. Comput. Linguist.*, vol. 5, pp. 135–146, 2017.
- 822 [43] Y. Liu *et al.*, “RoBERTa: A Robustly Optimized BERT Pretraining Approach,” *ArXiv Prepr. ArXiv190711692*,
823 2019, doi: 10.48550/arXiv.1907.11692.
- 824 [44] T. Gao, X. Yao, and D. Chen, “SimCSE: Simple Contrastive Learning of Sentence Embeddings,” in
825 *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, M.-F. Moens, X.
826 Huang, L. Specia, and S. W. Yih, Eds., Online and Punta Cana, Dominican Republic: Association for
827 Computational Linguistics, Nov. 2021, pp. 6894–6910. doi: 10.18653/v1/2021.emnlp-main.552.
- 828 [45] B. T. Garcia *et al.*, “Improving Automated Deep Phenotyping Through Large Language Models Using
829 Retrieval Augmented Generation,” *MedRxiv Prepr.*, p. 2024.12.01.24318253, Dec. 2024, doi:
830 10.1101/2024.12.01.24318253.

- 831 [46] Q. Team, “fastembed: A blazing fast embedding library.” 2023. [Online]. Available:
832 <https://github.com/qdrant/fastembed>
- 833 [47] M. Li, H. Zhou, H. Yang, and R. Zhang, “RT: a Retrieving and Chain-of-Thought framework for few-shot
834 medical named entity recognition,” *J. Am. Med. Inform. Assoc.*, vol. 31, no. 9, pp. 1929–1938, Sept. 2024, doi:
835 10.1093/jamia/ocae095.
- 836 [48] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for
837 language understanding,” in *Proceedings of the 2019 conference of the North American chapter of the*
838 *association for computational linguistics: human language technologies, volume 1 (long and short papers)*,
839 2019, pp. 4171–4186.
- 840 [49] R. Xu, Y. Hong, F. Zhang, and H. Xu, “Evaluation of the integration of retrieval-augmented generation in large
841 language model for breast cancer nursing care responses,” *Sci. Rep.*, vol. 14, no. 1, p. 30794, Dec. 2024, doi:
842 10.1038/s41598-024-81052-3.
- 843 [50] S. S. Darnell *et al.*, “Creating a biomedical knowledge base by addressing GPT inaccurate responses and
844 benchmarking context,” *bioRxiv*, 2024, doi: 10.1101/2024.10.16.618663.
- 845 [51] Y. A. Malkov and D. A. Yashunin, “Efficient and Robust Approximate Nearest Neighbor Search Using
846 Hierarchical Navigable Small World Graphs,” *CoRR*, vol. abs/1603.09320, 2016, [Online]. Available:
847 <https://arxiv.org/abs/1603.09320>
- 848 [52] E. Klang *et al.*, “Assessing Retrieval-Augmented Large Language Model Performance in Emergency
849 Department ICD-10-CM Coding Compared to Human Coders,” *MedRxiv Prepr.*, Oct. 2024, doi:
850 10.1101/2024.10.15.24315526.
- 851 [53] Y. Zheng *et al.*, “Integrating retrieval-augmented generation for enhanced personalized physician
852 recommendations in web-based medical services: model development study,” *Front. Public Health*, vol. 13, p.
853 1501408, Jan. 2025, doi: 10.3389/fpubh.2025.1501408.
- 854 [54] T. Fukushima *et al.*, “Evaluating and Enhancing Japanese Large Language Models for Genetic Counseling
855 Support: Comparative Study of Domain Adaptation and the Development of an Expert-Evaluated Dataset,”
856 *JMIR Med. Inform.*, vol. 13, p. e65047, Jan. 2025, doi: 10.2196/65047.
- 857 [55] A. Fukuchi, Y. Hoshino, and Y. Watanabe, “GLuCoSE (General Luke-based Contrastive Sentence
858 Embedding).” 2023. [Online]. Available: <https://huggingface.co/pkshatech/GLuCoSE-base-ja>
- 859 [56] S. Toro *et al.*, “Dynamic Retrieval Augmented Generation of Ontologies using Artificial Intelligence
860 (DRAGON-AI),” *J. Biomed. Semant.*, vol. 15, no. 1, p. 19, Oct. 2024, doi: 10.1186/s13326-024-00320-3.
- 861 [57] OpenAI, “Vector Embeddings.” Accessed: Feb. 19, 2025. [Online]. Available:
862 <https://platform.openai.com/docs/guides/embeddings>
- 863 [58] D. Wang *et al.*, “Enhancement of the Performance of Large Language Models in Diabetes Education through
864 Retrieval-Augmented Generation: Comparative Study,” *J. Med. Internet Res.*, vol. 26, p. e58041, Nov. 2024,
865 doi: 10.2196/58041.
- 866 [59] T. C. Li *et al.*, “FAVOR-GPT: a generative natural language interface to whole genome variant functional
867 annotations,” *Bioinforma. Adv.*, vol. 4, no. 1, p. vbae143, Sept. 2024, doi: 10.1093/bioadv/vbae143.
- 868 [60] Z. Fu *et al.*, “Application of large language model combined with retrieval enhanced generation technology in
869 digestive endoscopic nursing,” *Front. Med. Lausanne*, vol. 11, p. 1500258, Nov. 2024, doi:
870 10.3389/fmed.2024.1500258.
- 871 [61] C. Long *et al.*, “ChatENT: Augmented Large Language Model for Expert Knowledge Retrieval in
872 Otolaryngology-Head and Neck Surgery,” *Otolaryngol.-Head Neck Surg.*, vol. 171, no. 4, pp. 1042–1051, Oct.
873 2024, doi: 10.1002/ohn.864.
- 874 [62] D. Steybe *et al.*, “Evaluation of a context-aware chatbot using retrieval-augmented generation for answering
875 clinical questions on medication-related osteonecrosis of the jaw,” *J Craniomaxillofac Surg*, pp. S1010-
876 5182(24)00341-X, Jan. 2025, doi: 10.1016/j.jcms.2024.12.009.
- 877 [63] O. Unlu *et al.*, “Retrieval Augmented Generation Enabled Generative Pre-Trained Transformer 4 (GPT-4)
878 Performance for Clinical Trial Screening,” *MedRxiv Prepr.*, p. 2024.02.08.24302376, Feb. 2024, doi:
879 10.1101/2024.02.08.24302376.
- 880 [64] I. Azimi, M. Qi, L. Wang, A. M. Rahmani, and Y. Li, “Evaluation of LLMs accuracy and consistency in the
881 registered dietitian exam through prompt engineering and knowledge retrieval,” *Sci Rep*, vol. 15, no. 1, p.
882 1506, Jan. 2025, doi: 10.1038/s41598-024-85003-w.
- 883 [65] Amazon Web Services, “Amazon Titan Embeddings Models - User Guide.” 2024. [Online]. Available:
884 <https://docs.aws.amazon.com/bedrock/latest/userguide/titan-embedding-models.html>
- 885 [66] Y. Selcuk, E. Kim, and I. Ahn, “InfectA-Chat, an Arabic Large Language Model for Infectious Diseases:
886 Comparative Analysis,” *JMIR Med Inf.*, vol. 13, p. e63881, Feb. 2025, doi: 10.2196/63881.
- 887 [67] D. Kainer, “The effectiveness of large language models with RAG for auto-annotating trait and phenotype
888 descriptions,” *Biol Methods Protoc*, vol. 10, no. 1, p. bpaf016, Feb. 2025, doi: 10.1093/biomet/bpaf016.
- 889 [68] J. Lee *et al.*, “BioBERT: a pre-trained biomedical language representation model for biomedical text mining,”
890 *Bioinformatics*, Sept. 2019, doi: 10.1093/bioinformatics/btz682.

- 891 [69] Q. Jin, A. Shin, and Z. Lu, "LADER: Log-Augmented DENSE Retrieval for Biomedical Literature Search," in
892 *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information*
893 *Retrieval*, in SIGIR '23. New York, NY, USA: Association for Computing Machinery, 2023, pp. 2092–2097.
894 doi: 10.1145/3539618.3592005.
- 895 [70] Y. Gu *et al.*, "Domain-Specific Language Model Pretraining for Biomedical Natural Language Processing."
896 2020.
- 897 [71] I. Lopez *et al.*, "Clinical entity augmented retrieval for clinical information extraction," *NPJ Digit. Med.*, vol.
898 8, no. 1, p. 45, Jan. 2025, doi: 10.1038/s41746-024-01377-1.
- 899 [72] K. Soman *et al.*, "Biomedical knowledge graph-optimized prompt generation for large language models,"
900 *Bioinformatics*, vol. 40, no. 9, p. btac560, Sept. 2024, doi: 10.1093/bioinformatics/btac560.
- 901 [73] W. Aftab, Z. Apostolou, K. Bouazoune, and T. Straub, "Optimizing biomedical information retrieval with a
902 keyword frequency-driven prompt enhancement strategy," *BMC Bioinformatics*, vol. 25, no. 1, p. 281, Aug.
903 2024, doi: 10.1186/s12859-024-05902-7.
- 904 [74] R. Luo *et al.*, "BioGPT: Generative Pre-trained Transformer for Biomedical Text Generation and Mining,"
905 *Brief. Bioinform.*, vol. 23, no. 6, p. bbac409, Nov. 2022, doi: 10.1093/bib/bbac409.
- 906 [75] Z. Zhan, S. Zhou, M. Li, and R. Zhang, "RAMIE: retrieval-augmented multi-task information extraction with
907 large language models on dietary supplements," *J. Am. Med. Inform. Assoc.*, vol. 32, no. 3, pp. 545–554, Mar.
908 2025, doi: 10.1093/jamia/ocaf002.
- 909 [76] Y. Anand, Z. Nussbaum, B. Duderstadt, and B. Schmidt, "GPT4All: Training an Assistant-style Chatbot with
910 Large Scale Data Distillation from GPT-3.5-Turbo." 2023. [Online]. Available: [https://github.com/nomic-](https://github.com/nomic-ai/gpt4all)
911 [ai/gpt4all](https://github.com/nomic-ai/gpt4all)
- 912 [77] S. Myers *et al.*, "Lessons learned on information retrieval in electronic health records: a comparison of
913 embedding models and pooling strategies," *J. Am. Med. Inform. Assoc.*, vol. 32, no. 2, pp. 357–364, Feb. 2025,
914 doi: 10.1093/jamia/ocae308.
- 915 [78] L. Wu *et al.*, "A framework enabling LLMs into regulatory environment for transparency and trustworthiness
916 and its application to drug labeling document," *Regul. Toxicol. Pharmacol.*, vol. 149, p. 105613, 2024, doi:
917 <https://doi.org/10.1016/j.yrtph.2024.105613>.
- 918 [79] Facebook Research, "FAISS: A Library for Efficient Similarity Search." Accessed: Feb. 19, 2025. [Online].
919 Available: <https://github.com/facebookresearch/faiss>
- 920 [80] Chroma, "Getting Started with Chroma." Accessed: Feb. 19, 2025. [Online]. Available:
921 <https://docs.trychroma.com/docs/overview/getting-started>
- 922 [81] Pinecone, "Pinecone Overview." 2024. [Online]. Available: [https://docs.pinecone.io/guides/getting-](https://docs.pinecone.io/guides/getting-started/overview)
923 [started/overview](https://docs.pinecone.io/guides/getting-started/overview)
- 924 [82] Weaviate, "Weaviate Developers Documentation." 2024. [Online]. Available:
925 <https://weaviate.io/developers/weaviate>
- 926 [83] S. Sivarajkumar *et al.*, "Clinical Information Retrieval: A Literature Review," *J. Healthc. Inform. Res.*, vol. 8,
927 no. 2, pp. 313–352, 2024, doi: 10.1007/s41666-024-00159-4.
- 928 [84] D. B. Craig and S. Drăghici, "LmRaC: a functionally extensible tool for LLM interrogation of user
929 experimental results," *Bioinformatics*, vol. 40, no. 12, p. btac679, Nov. 2024, doi:
930 10.1093/bioinformatics/btac679.
- 931 [85] J. Wu, A. Shrivastava, J. Zhu, A. Samuel, A. Kumar, and D. Liu, "LLM Optimization Unlocks Real-Time
932 Pairwise Reranking," Nov. 10, 2025, *arXiv*: arXiv:2511.07555. doi: 10.48550/arXiv.2511.07555.
- 933 [86] J. Lee, L. H. Pham, and Ö. Uzuner, "Enhancing Consumer Health Question Reformulation: Chain-of-Thought
934 Prompting Integrating Focus, Type, and User Knowledge Level," in *Proceedings of the First Workshop on*
935 *Patient-Oriented Language Processing (CL4Health) @ LREC-COLING 2024*, Torino, Italia, pp. 220–228.
- 936 [87] W. X. Zhao, K. Zhou, J. Li, and others, "A Survey of Large Language Models," *ArXiv Prepr.*, vol.
937 arXiv:2303.18223, no. 1(2), 2023, [Online]. Available: <https://arxiv.org/abs/2303.18223>
- 938 [88] H. Touvron *et al.*, "LLaMA 2: Open Foundation and Fine-Tuned Chat Models," *ArXiv Prepr.*, 2023, [Online].
939 Available: <https://arxiv.org/abs/2307.09288>
- 940 [89] Google AI, "Gemini API: Text generation guide." Accessed: Feb. 19, 2025. [Online]. Available:
941 <https://ai.google.dev/gemini-api/docs/text-generation>
- 942 [90] Mistral AI, "Mistral-7B-Instruct-v0.3." 2024. [Online]. Available: [https://huggingface.co/mistralai/Mistral-](https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3)
943 [7B-Instruct-v0.3](https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3)
- 944 [91] C. Raffel *et al.*, "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer," *J. Mach.*
945 *Learn. Res.*, vol. 21, no. 140, pp. 1–67, 2020.
- 946 [92] Mistral AI, "Mistral-Small-3.1-24B-Instruct-2503." 2024. [Online]. Available:
947 <https://huggingface.co/mistralai/Mistral-Small-3.1-24B-Instruct-2503>
- 948 [93] Google, "Gemma 3 27B Instruct - Hugging Face." 2024. [Online]. Available:
949 <https://huggingface.co/google/gemma-3-27b-it>
- 950 [94] J. Lammert *et al.*, "Expert-Guided Large Language Models for Clinical Decision Support in Precision
951 Oncology," *JCO Precis. Oncol.*, vol. 8, p. e2400478, Oct. 2024, doi: 10.1200/PO-24-00478.

- 952 [95] K. Hewitt *et al.*, “Large language models as a diagnostic support tool in neuropathology,” *J. Pathol. Clin. Res.*,
953 vol. 10, no. 6, p. e70009, Nov. 2024, doi: 10.1002/2056-4538.70009.
- 954 [96] Y. Labrak, A. Bazoge, E. Morin, P.-A. Gourraud, M. Rouvier, and R. Dufour, “BioMistral: A Collection of
955 Open-Source Pretrained Large Language Models for Medical Domains.” Association for Computational
956 Linguistics, 2024.
- 957 [97] S. Ozaki *et al.*, “Understanding the Impact of Confidence in Retrieval Augmented Generation: A Case Study
958 in the Medical Domain.” 2024. [Online]. Available: <https://arxiv.org/abs/2412.17895>
- 959 [98] M. Alkaeed *et al.*, “Open Foundation Models in Healthcare: Challenges, Paradoxes, and Opportunities with
960 GenAI Driven Personalized Prescription.” 2025. [Online]. Available: <https://arxiv.org/abs/2502.04356>
- 961 [99] B. Peng *et al.*, “Graph Retrieval-Augmented Generation: A Survey.” 2024. [Online]. Available:
962 <https://arxiv.org/abs/2408.08092>
- 963 [100] A. Di Maria *et al.*, “NetMe 2.0: a web-based platform for extracting and modeling knowledge from
964 biomedical literature as a labeled graph,” *Bioinformatics*, vol. 40, no. 5, p. btae194, May 2024, doi:
965 10.1093/bioinformatics/btae194.
- 966 [101] Y. Gao *et al.*, “Leveraging Medical Knowledge Graphs Into Large Language Models for Diagnosis
967 Prediction: Design and Application Study,” *JMIR AI*, vol. 4, p. e58670, Feb. 2025, doi: 10.2196/58670.
- 968 [102] N. Matsumoto *et al.*, “ESCARGOT: an AI agent leveraging large language models, dynamic graph of
969 thoughts, and biomedical knowledge graphs for enhanced reasoning,” *Bioinformatics*, vol. 41, no. 2, p. btaf031,
970 Feb. 2025, doi: 10.1093/bioinformatics/btaf031.
- 971 [103] C. Kakalou, C. Karamanidou, T. Dalamagas, and M. Koubarakis, “Enhancing Patient Empowerment and
972 Health Literacy: Integrating Knowledge Graphs with Language Models for Personalized Health Content
973 Delivery,” in *Studies in Health Technology and Informatics*, Aug. 2024, pp. 1018–1022. doi:
974 10.3233/SHTI240582.
- 975 [104] N. Matsumoto *et al.*, “KRAGEN: a knowledge graph-enhanced RAG framework for biomedical problem
976 solving using large language models,” *Bioinformatics*, vol. 40, no. 6, p. btae353, June 2024, doi:
977 10.1093/bioinformatics/btae353.
- 978 [105] R. Yang *et al.*, “Ascle-A Python Natural Language Processing Toolkit for Medical Text Generation:
979 Development and Evaluation Study,” *J. Med. Internet Res.*, vol. 26, p. e60601, Oct. 2024, doi: 10.2196/60601.
- 980 [106] D. Li *et al.*, “DALK: Dynamic Co-Augmentation of LLMs and KG to answer Alzheimer’s Disease
981 Questions with Scientific Literature,” in *Findings of the Association for Computational Linguistics: EMNLP*
982 2024, Miami, Florida, USA: Association for Computational Linguistics, 2024, pp. 2187–2205.
- 983 [107] S. Gilbert, J. N. Kather, and A. Hogan, “Augmented non-hallucinating large language models as medical
984 information curators,” *NPJ Digit. Med.*, vol. 7, no. 1, p. 100, Apr. 2024, doi: 10.1038/s41746-024-01081-0.
- 985 [108] Z. Xiao, H. B. Pakrasi, Y. Chen, and Y. J. Tang, “Network for knowledge Organization (NEKO): An AI
986 knowledge mining workflow for synthetic biology research,” *Metab. Eng.*, vol. 87, pp. 60–67, 2025, doi:
987 <https://doi.org/10.1016/j.ymben.2024.11.006>.
- 988 [109] J. H. Morris *et al.*, “The scalable precision medicine open knowledge engine (SPOKE): a massive
989 knowledge graph of biomedical information,” *Bioinformatics*, vol. 39, no. 2, p. btad080, 2023, doi:
990 10.1093/bioinformatics/btad080.
- 991 [110] M. M. Abootorabi *et al.*, “Ask in Any Modality: A Comprehensive Survey on Multimodal Retrieval-
992 Augmented Generation.” 2025. [Online]. Available: <https://arxiv.org/abs/2502.08826>
- 993 [111] A. Radford *et al.*, “Learning Transferable Visual Models From Natural Language Supervision,” Feb. 26,
994 2021, *arXiv: arXiv:2103.00020*. doi: 10.48550/arXiv.2103.00020.
- 995 [112] L. Sun, J. Zhao, M. Han, and C. Xiong, “Fact-aware multimodal retrieval augmentation for accurate
996 medical radiology report generation,” *ArXiv Prepr. ArXiv240715268*, 2024.
- 997 [113] M. Ranjit, G. Ganapathy, R. Manuel, and T. Ganu, “Retrieval augmented chest x-ray report generation
998 using openai gpt models,” in *Machine Learning for Healthcare Conference*, PMLR, 2023, pp. 650–666.
- 999 [114] P. Xia *et al.*, “MMed-RAG: Versatile Multimodal RAG System for Medical Vision Language Models,”
1000 *ArXiv Prepr. ArXiv241013085*, 2024.
- 1001 [115] P. Xia *et al.*, “RULE: Reliable Multimodal RAG for Factuality in Medical Vision Language Models,” in
1002 *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Miami, Florida,
1003 USA: Association for Computational Linguistics, 2024, pp. 1081–1093.
- 1004 [116] Y. Yang and others, “Spatio-Temporal and Retrieval-Augmented Modelling for Chest X-Ray Report
1005 Generation,” *IEEE Trans. Med. Imaging*, 2025, doi: 10.1109/TMI.2025.3554498.
- 1006 [117] M. Omar, V. Ullanat, M. Loda, L. Marchionni, and R. Umeton, “ChatGPT for digital pathology research,”
1007 *Lancet Digit. Health*, vol. 6, no. 8, pp. e595–e600, Aug. 2024, doi: 10.1016/S2589-7500(24)00114-6.
- 1008 [118] D. P. Upadhyaya *et al.*, “A 360-degree View for Large Language Models: Early Detection of Amblyopia
1009 in Children using Multi-View Eye Movement Recordings,” *MedRxiv Prepr.*, p. 2024.05.03.24306688, May
1010 2024, doi: 10.1101/2024.05.03.24306688.
- 1011 [119] R. Tozuka *et al.*, “Application of NotebookLM, a large language model with retrieval-augmented
1012 generation, for lung cancer staging,” *Jpn. J. Radiol.*, Nov. 2024, doi: 10.1007/s11604-024-01705-1.

- 1013 [120] S. Ramedini, S. Shridevi, and D. Won, “Multi-modal transformer architecture for medical image analysis
1014 and automated report generation,” *Sci. Rep.*, vol. 14, no. 1, p. 19281, Aug. 2024, doi: 10.1038/s41598-024-
1015 69981-5.
- 1016 [121] P. Thetbanthad, B. Sathanarugsawait, and P. Praneetpolgrang, “Application of Generative Artificial
1017 Intelligence Models for Accurate Prescription Label Identification and Information Retrieval for the Elderly in
1018 Northern East of Thailand,” *J Imaging*, vol. 11, no. 1, p. 11, Jan. 2025, doi: 10.3390/jimaging11010011.
- 1019 [122] D. Hu *et al.*, “Pathology report generation from whole slide images with knowledge retrieval and multi-
1020 level regional feature selection,” *Comput Methods Programs Biomed*, vol. 263, p. 108677, May 2025, doi:
1021 10.1016/j.cmpb.2025.108677.
- 1022 [123] X. Chen, C.-J. Hsieh, and B. Gong, “When Vision Transformers Outperform ResNets without Pretraining
1023 or Strong Data Augmentations,” *ArXiv Prepr. ArXiv210601548*, 2021.
- 1024 [124] Z. Huang *et al.*, “Tool Calling: Enhancing Medication Consultation via Retrieval-Augmented Large
1025 Language Models.” 2024. [Online]. Available: <https://arxiv.org/abs/2404.17897>
- 1026 [125] O. Bodenreider, “The unified medical language system (umls): integrating biomedical terminology,”
1027 *Nucleic Acids Res.*, vol. 32, no. suppl_1, pp. D267–D270, 2004.
- 1028 [126] A. E. Johnson *et al.*, “MIMIC-IV, a freely accessible electronic health record dataset,” *Sci. Data*, vol. 10,
1029 no. 1, p. 1, 2023.
- 1030 [127] D. Moser, M. Bender, and M. Sariyar, “Generating Synthetic Healthcare Dialogues in Emergency
1031 Medicine Using Large Language Models,” *Stud. Health Technol. Inform.*, vol. 321, pp. 235–239, Nov. 2024,
1032 doi: 10.3233/SHTI241099.
- 1033 [128] S. Jia *et al.*, “PodGPT: An audio-augmented large language model for research and education.” Nov. 2024.
1034 doi: 10.1101/2024.07.11.24310304.
- 1035 [129] K. Bashir and H. Bukumiric, *Basic Metabolic Panel*. Treasure Island (FL): StatPearls Publishing, 2024.
- 1036 [130] D. Jin, E. Pan, N. Oufattole, W.-H. Weng, H. Fang, and P. Szolovits, “What disease does this patient have?
1037 a large-scale open domain question answering dataset from medical exams,” *Appl. Sci.*, vol. 11, no. 14, p. 6421,
1038 2021.
- 1039 [131] G. Xiong, Q. Jin, X. Wang, M. Zhang, Z. Lu, and A. Zhang, “Improving Retrieval-Augmented Generation
1040 in Medicine with Iterative Follow-up Questions,” *Pac. Symp. Biocomput.*, vol. 30, pp. 199–214, 2025.
- 1041 [132] J. Irvin *et al.*, “CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison,”
1042 in *Proceedings of the AAAI conference on artificial intelligence*, 2019, pp. 590–597.
- 1043 [133] Q. Jin, B. Dhingra, Z. Liu, W. Cohen, and X. Lu, “PubMedQA: A dataset for biomedical research question
1044 answering,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*
1045 *and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Association
1046 for Computational Linguistics, Nov. 2019, pp. 2567–2577.
- 1047 [134] A. E. Johnson *et al.*, “MIMIC-III, a freely accessible critical care database,” *Sci. Data*, vol. 3, no. 1, pp.
1048 1–9, 2016.
- 1049 [135] G. Tsatsaronis, G. Balikas, P. Malakasiotis, and others, “An overview of the BIOASQ large-scale
1050 biomedical semantic indexing and question answering competition,” *BMC Bioinformatics*, vol. 16, pp. 1–28,
1051 2015, doi: 10.1186/s12859-015-0564-6.
- 1052 [136] O. Taboureau *et al.*, “ChemProt: a disease chemical biology database,” *Nucleic Acids Res.*, vol. 39, no.
1053 suppl_1, pp. D367–D372, 2010, doi: 10.1093/nar/gkq870.
- 1054 [137] D. Hendrycks *et al.*, “Measuring Massive Multitask Language Understanding,” *Proc. Int. Conf. Learn.*
1055 *Represent. ICLR*, 2021.
- 1056 [138] D. Teodoro, N. Naderi, A. Yazdani, B. Zhang, and A. Bornet, “A scoping review of artificial intelligence
1057 applications in clinical trial risk assessment,” *Npj Digit. Med.*, vol. 8, no. 1, p. 486, July 2025, doi:
1058 10.1038/s41746-025-01886-7.
- 1059 [139] C.-Y. Lin, “ROUGE: A Package for Automatic Evaluation of Summaries,” in *Annual Meeting of the*
1060 *Association for Computational Linguistics*, 2004. [Online]. Available:
1061 <https://api.semanticscholar.org/CorpusID:964287>
- 1062 [140] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, “BLEU: a Method for Automatic Evaluation of Machine
1063 Translation,” in *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*,
1064 Philadelphia, Pennsylvania, USA: Association for Computational Linguistics, 2002, pp. 311–318. doi:
1065 10.3115/1073083.1073135.
- 1066 [141] T. Zhang*, V. Kishore*, F. Wu*, K. Q. Weinberger, and Y. Artzi, “BERTScore: Evaluating Text Generation
1067 with BERT,” in *International Conference on Learning Representations*, 2020. [Online]. Available:
1068 <https://openreview.net/forum?id=SkeHuCVFDr>
- 1069 [142] T. Sellam, D. Das, and A. Parikh, “BLEURT: Learning Robust Metrics for Text Generation,” in
1070 *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, D. Jurafsky, J. Chai,
1071 N. Schluter, and J. Tetreault, Eds., Online: Association for Computational Linguistics, July 2020, pp. 7881–
1072 7892. doi: 10.18653/v1/2020.acl-main.704.

- 1073 [143] W. Zhao, C. Wu, X. Zhang, Y. Zhang, Y. Wang, and W. Xie, “RaTEScore: A Metric for Radiology Report
1074 Generation,” in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*,
1075 2024, pp. 15004–15019.
- 1076 [144] X. Li, W. Wu, L. Qin, and Q. Yin, “How to Evaluate Your Dialogue Models: A Review of Approaches,”
1077 *ArXiv Prepr. ArXiv210801369*, 2021, [Online]. Available: <https://arxiv.org/abs/2108.01369>
- 1078 [145] E. Croxford, Y. Gao, E. First, and others, “Evaluating clinical AI summaries with large language models
1079 as judges,” *Npj Digit. Med.*, vol. 8, p. 640, 2025, doi: 10.1038/s41746-025-02005-2.
- 1080 [146] A. Ben Abacha, W. Yim, G. Michalopoulos, and T. Lin, “An Investigation of Evaluation Methods in
1081 Automatic Medical Note Generation,” in *Findings of the Association for Computational Linguistics: ACL*
1082 *2023*, A. Rogers, J. Boyd-Graber, and N. Okazaki, Eds., Toronto, Canada: Association for Computational
1083 Linguistics, July 2023, pp. 2575–2588. doi: 10.18653/v1/2023.findings-acl.161.
- 1084 [147] S. Liu, A. B. McCoy, and A. Wright, “Improving large language model applications in biomedicine with
1085 retrieval-augmented generation: a systematic review, meta-analysis, and clinical development guidelines,” *J.*
1086 *Am. Med. Inform. Assoc.*, vol. 32, no. 4, pp. 605–615, Apr. 2025, doi: 10.1093/jamia/ocaf008.
- 1087 [148] Y. Guo, W. Qiu, G. Leroy, S. Wang, and T. Cohen, “Retrieval augmentation of large language models for
1088 lay language generation,” *J. Biomed. Inform.*, vol. 149, p. 104580, Jan. 2024, doi: 10.1016/j.jbi.2023.104580.
- 1089 [149] M. Coleman and T. L. Liau, “A computer readability formula designed for machine scoring,” *J. Appl.*
1090 *Psychol.*, vol. 60, no. 2, pp. 283–284, 1975, doi: 10.1037/h0076540.
- 1091 [150] R. T. Sutton, D. Pincock, D. C. Baumgart, D. C. Sadowski, R. N. Fedorak, and K. I. Kroeker, “An overview
1092 of clinical decision support systems: benefits, risks, and strategies for success,” *NPJ Digit. Med.*, vol. 3, p. 17,
1093 Feb. 2020, doi: 10.1038/s41746-020-0221-y.
- 1094 [151] C. S. Ong, N. T. Obey, Y. Zheng, A. Cohan, and E. B. Schneider, “SurgeryLLM: a retrieval-augmented
1095 generation large language model framework for surgical decision support and workflow enhancement,” *NPJ*
1096 *Digit. Med.*, vol. 7, no. 1, p. 364, Dec. 2024, doi: 10.1038/s41746-024-01391-3.
- 1097 [152] C. Martínez-deMiguel, I. Segura-Bedmar, E. Chacón-Solano, and S. Guerrero-Aspizua, “The RareDis
1098 corpus: A corpus annotated with rare diseases, their signs and symptoms,” *J. Biomed. Inform.*, vol. 125, p.
1099 103961, 2022, doi: 10.1016/j.jbi.2021.103961.
- 1100 [153] J. Yang, L. Shu, H. Duan, and H. Li, “RDguru: A Conversational Intelligent Agent for Rare Diseases,”
1101 *IEEE J. Biomed. Health Inform.*, vol. PP, 2024, doi: 10.1109/JBHI.2024.3464555.
- 1102 [154] John Snow Labs, “Entity Resolver for Human Phenotype Ontology.” 2024. [Online]. Available:
1103 https://nlp.johnsnowlabs.com/2021/05/16/sbiobertresolve_HPO_en.html
- 1104 [155] C. Zelin, W. K. Chung, M. Jeanne, G. Zhang, and C. Weng, “Rare disease diagnosis using knowledge
1105 guided retrieval augmentation for ChatGPT,” *J. Biomed. Inform.*, vol. 157, p. 104702, Sept. 2024, doi:
1106 10.1016/j.jbi.2024.104702.
- 1107 [156] P. Sloan, P. Clatworthy, E. Simpson, and M. Mirmehdi, “Automated radiology report generation: A review
1108 of recent advances,” *IEEE Rev. Biomed. Eng.*, 2024.
- 1109 [157] I. E. Hamamci, S. Er, and B. Menze, “Ct2rep: Automated radiology report generation for 3d medical
1110 imaging,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*,
1111 Springer, 2024, pp. 476–486.
- 1112 [158] O. Alfarghaly, R. Khaled, A. Elkorany, M. Helal, and A. Fahmy, “Automated radiology report generation
1113 using conditioned transformers,” *Inform. Med. Unlocked*, vol. 24, p. 100557, 2021, doi:
1114 10.1016/j.imu.2021.100557.
- 1115 [159] S. Song, A. Subramanyam, I. Madejski, and R. L. Grossman, “LaB-RAG: Label Boosted Retrieval
1116 Augmented Generation for Radiology Report Generation,” *ArXiv Prepr. ArXiv241116523*, 2024.
- 1117 [160] H. M. T. Alam, D. Srivastav, M. A. Kadir, and D. Sonntag, “Towards Interpretable Radiology Report
1118 Generation via Concept Bottlenecks using a Multi-Agentic RAG,” *ArXiv Prepr. ArXiv241216086*, 2024.
- 1119 [161] M. J. Landrum *et al.*, “ClinVar: improving access to variant interpretations and supporting evidence,”
1120 *Nucleic Acids Res.*, vol. 46, no. D1, pp. D1062–D1067, Jan. 2018, doi: 10.1093/nar/gkx1153.
- 1121 [162] X. Chen *et al.*, “EyeGPT for Patient Inquiries and Medical Education: Development and Validation of an
1122 Ophthalmology Large Language Model,” *J. Med. Internet Res.*, vol. 26, p. e60063, 2024, doi: 10.2196/60063.
- 1123 [163] M. Sevgi, F. Antaki, and P. A. Keane, “Medical education with large language models in ophthalmology:
1124 custom instructions and enhanced retrieval capabilities,” *Br. J. Ophthalmol.*, vol. 108, no. 10, pp. 1354–1361,
1125 2024, doi: 10.1136/bjo-2023-325046.
- 1126 [164] OpenAI, “Text generation - OpenAI API.” Accessed: Feb. 19, 2025. [Online]. Available:
1127 <https://platform.openai.com/docs/guides/text-generation>
- 1128 [165] M. AI, “Llama 3: 70B high-performance open language models,” *Meta AI Res.*, 2024, [Online]. Available:
1129 <https://ai.meta.com/llama/papers-and-research/>
- 1130 [166] Qwen Team, “Qwen/QwQ-32B.” 2024. [Online]. Available: <https://huggingface.co/Qwen/QwQ-32B>
- 1131 [167] T. Han *et al.*, “MedAlpaca—An Open-Source Collection of Medical Conversational AI Models and Training
1132 Data,” *ArXiv Prepr. ArXiv230408247*, 2023.

- 1133 [168] B. Pingua *et al.*, “Medical LLMs: Fine-Tuning vs. Retrieval-Augmented Generation,” *Bioeng. Basel*
1134 *Switz.*, vol. 12, no. 7, p. 687, June 2025, doi: 10.3390/bioengineering12070687.
- 1135 [169] F. J. Dorfner *et al.*, “Evaluating the effectiveness of biomedical fine-tuning for large language models on
1136 clinical tasks,” *J. Am. Med. Inform. Assoc. JAMIA*, vol. 32, no. 6, pp. 1015–1024, June 2025, doi:
1137 10.1093/jamia/ocaf045.
- 1138 [170] Elastic, “Elasticsearch Documentation: Getting Started Guide.” 2024. [Online]. Available:
1139 <https://www.elastic.co/guide/en/elasticsearch/reference/current/getting-started.html>
- 1140 [171] DeepSeek-AI Team, “DeepSeek-AI/DeepSeek-R1.” 2024. [Online]. Available:
1141 <https://huggingface.co/deepseek-ai/DeepSeek-R1>
- 1142 [172] M. Garg, S. Raza, S. Rayana, X. Liu, and S. Sohn, “The Rise of Small Language Models in Healthcare: A
1143 Comprehensive Survey,” Apr. 25, 2025, *arXiv*: arXiv:2504.17119. doi: 10.48550/arXiv.2504.17119.
- 1144 [173] M. Abdin *et al.*, “Phi-4 Technical Report,” Dec. 12, 2024, *arXiv*: arXiv:2412.08905. doi:
1145 10.48550/arXiv.2412.08905.
- 1146 [174] B. Parmanto *et al.*, “A Reliable and Accessible Caregiving Language Model (CaLM) to Support Tools for
1147 Caregivers: Development and Evaluation Study,” *JMIR Form. Res.*, vol. 8, p. e54633, July 2024, doi:
1148 10.2196/54633.
- 1149 [175] Q. Jin *et al.*, “State-of-the-Art Evidence Retriever for Precision Medicine: Algorithm Development and
1150 Validation,” *JMIR Med. Inform.*, vol. 10, no. 12, p. e40743, Dec. 2022, doi: 10.2196/40743.
- 1151 [176] A. Cohan, S. Feldman, I. Beltagy, D. Downey, and D. S. Weld, “SPECTER: Document-level
1152 Representation Learning using Citation-informed Transformers,” in *ACL*, 2020.
- 1153 [177] T. GLM *et al.*, “ChatGLM: A Family of Large Language Models from GLM-130B to GLM-4 All Tools.”
1154 2024.
- 1155 [178] L. Zhao, Y. Wang, X. Wang, R. Lin, Z. Gang, and B. Xu, “COPD-ChatGLM: A Chronic Obstructive
1156 Pulmonary Disease Diagnostic Model,” in *2024 IEEE International Conference on Bioinformatics and*
1157 *Biomedicine (BIBM)*, 2024, pp. 2965–2970. doi: 10.1109/BIBM62325.2024.10822841.
- 1158 [179] M. J. Luo *et al.*, “Development and Evaluation of a Retrieval-Augmented Large Language Model
1159 Framework for Ophthalmology,” *JAMA Ophthalmol.*, vol. 142, no. 9, pp. 798–805, 2024, doi:
1160 10.1001/jamaophthalmol.2024.2513.
- 1161 [180] K. Kreimeyer, J. V. Canzoniero, M. Fattah, V. Anagnostou, and T. Botsis, “Using Retrieval-Augmented
1162 Generation to Capture Molecularly-Driven Treatment Relationships for Precision Oncology,” *Stud. Health*
1163 *Technol. Inform.*, vol. 316, pp. 983–987, Aug. 2024, doi: 10.3233/SHTI240575.
- 1164 [181] Z. Zhan, J. Wang, S. Zhou, J. Deng, and R. Zhang, “MMRAG: Multi-Mode Retrieval-Augmented
1165 Generation with Large Language Models for Biomedical In-Context Learning.” 2025. [Online]. Available:
1166 <https://arxiv.org/abs/2502.15954>
- 1167 [182] H. Yu, P. Guo, and A. Sano, “Zero-Shot ECG Diagnosis with Large Language Models and Retrieval-
1168 Augmented Generation,” in *Proceedings of the 3rd Machine Learning for Health Symposium*, S. Hegselmann,
1169 A. Parziale, D. Shanmugam, S. Tang, M. N. Asiedu, S. Chang, T. Hartvigsen, and H. Singh, Eds., in
1170 *Proceedings of Machine Learning Research*, vol. 225. PMLR, Dec. 2023, pp. 650–663. [Online]. Available:
1171 <https://proceedings.mlr.press/v225/yu23b.html>
- 1172 [183] X. Du *et al.*, “Enhancing Early Detection of Cognitive Decline in the Elderly: A Comparative Study
1173 Utilizing Large Language Models in Clinical Notes,” *medRxiv*, p. 2024.04.03.24305298, 2024, doi:
1174 10.1101/2024.04.03.24305298.
- 1175 [184] A. Fatharani and A. Alsayegh, “Pharmacogenomics meets generative AI: transforming clinical trial design
1176 with large language models,” *J. Pharmacol. Pharmacother.*, p. 0976500X251321885, 2025, doi:
1177 10.1177/0976500X251321885.
- 1178 [185] J. J. Woo *et al.*, “Custom Large Language Models Improve Accuracy: Comparing Retrieval Augmented
1179 Generation and Artificial Intelligence Agents to Noncustom Models for Evidence-Based Medicine,”
1180 *Arthroscopy*, vol. 41, no. 3, pp. 565–573.e6, Mar. 2025, doi: 10.1016/j.arthro.2024.10.042.
- 1181 [186] Microsoft, “Phi-3-mini-4k-instruct.” 2024. [Online]. Available: [https://huggingface.co/microsoft/Phi-3-](https://huggingface.co/microsoft/Phi-3-mini-4k-instruct)
1182 [mini-4k-instruct](https://huggingface.co/microsoft/Phi-3-mini-4k-instruct)
- 1183 [187] M. N. Kamel Boulos and R. Dellavalle, “NVIDIA’s ‘Chat with RTX’ Custom Large Language Model and
1184 Personalized AI Chatbot Augments the Value of Electronic Dermatology Reference Material,” *JMIR*
1185 *Dermatol.*, vol. 7, p. e58396, July 2024, doi: 10.2196/58396.
- 1186 [188] B. Wang *et al.*, “Baichuan-M1: Pushing the Medical Capability of Large Language Models.” 2025.
1187 [Online]. Available: <https://arxiv.org/abs/2502.09298>
- 1188 [189] E. Almazrouei *et al.*, “The Falcon Series of Open Language Models,” *ArXiv Prepr. ArXiv231116867*, 2023,
1189 [Online]. Available: <https://arxiv.org/abs/2311.16867>
- 1190 [190] L. Tunstall *et al.*, “Zephyr: Direct Distillation of LM Alignment.” 2023.
- 1191 [191] H. Huang *et al.*, “AceGPT, Localizing Large Language Models in Arabic,” in *Proceedings of the 2024*
1192 *Conference of the North American Chapter of the Association for Computational Linguistics: Human*

- 1193 *Language Technologies (Volume I: Long Papers)*, Mexico City, Mexico: Association for Computational
1194 Linguistics, 2024, pp. 8139–8163.
- 1195 [192] T. G. T. Mesnard, C. Hardin, and others, “Gemma: Open models based on Gemini research and
1196 technology,” *ArXiv Prepr. ArXiv240308295*, 2024.
- 1197 [193] Y. Feng, J. Wang, R. He, L. Zhou, and Y. Li, “A Retrieval-Augmented Knowledge Mining Method with
1198 Deep Thinking LLMs for Biomedical Research and Clinical Support,” *ArXiv Prepr. ArXiv250323029*, 2025.
- 1199 [194] S. Rau *et al.*, “A retrieval-augmented chatbot based on GPT-4 provides appropriate differential diagnosis
1200 in gastrointestinal radiology: a proof of concept study,” *Eur. Radiol. Exp.*, vol. 8, no. 1, p. 60, 2024, doi:
1201 10.1186/s41747-024-00457-x.
- 1202 [195] J. R. Tan, D. Y. Z. Lim, Q. Le, and others, “ChatGPT performance in assessing musculoskeletal MRI scan
1203 appropriateness based on ACR appropriateness criteria,” *Sci. Rep.*, vol. 15, p. 7140, 2025, doi: 10.1038/s41598-
1204 025-88925-1.
- 1205 [196] S. Puts, C. M. L. Zegers, A. Dekker, and I. Bermejo, “Developing an ICD-10 Coding Assistant: Pilot Study
1206 Using RoBERTa and GPT-4 for Term Extraction and Description-Based Code Selection,” *JMIR Form. Res.*,
1207 vol. 9, p. e60095, 2025, doi: 10.2196/60095.
- 1208 [197] H. Choi, D. Lee, and Y. Kang, “Empowering PET Imaging Reporting with Retrieval-Augmented Large
1209 Language Models and Reading Reports Database: A Pilot Single Center Study,” *medRxiv*, 2024, doi:
1210 10.1101/2024.05.13.24307312.
- 1211 [198] N. Markey, I. El-Mansouri, G. Rensonnet, C. van Langen, and C. Meier, “From RAGs to riches: Utilizing
1212 large language models to write documents for clinical trials,” *Clin. Trials*, Feb. 2025, doi:
1213 10.1177/17407745251320806.
- 1214 [199] Anthropic, “Getting Started - Anthropic API Documentation.” Accessed: Feb. 19, 2025. [Online].
1215 Available: <https://docs.anthropic.com/en/api/getting-started>
- 1216 [200] I. Buhnla, A. Sinha, and M. Constant, “Retrieve, Generate, Evaluate: A Case Study for Medical
1217 Paraphrases Generation with Small Language Models.” 2024. [Online]. Available:
1218 <https://arxiv.org/abs/2407.16565>
- 1219 [201] Z. Chen *et al.*, “MEDITRON-70B: Scaling Medical Pretraining for Large Language Models,” Nov. 27,
1220 2023, *arXiv: arXiv:2311.16079*. doi: 10.48550/arXiv.2311.16079.
- 1221 [202] H. Kim *et al.*, “Small language models learn enhanced reasoning skills from medical textbooks,” *ArXiv*
1222 *Prepr. ArXiv240400376*, 2024.
- 1223 [203] J. Sohn *et al.*, “Rationale-Guided Retrieval Augmented Generation for Medical Question Answering.”
1224 2024. [Online]. Available: <https://arxiv.org/abs/2411.00300>
- 1225 [204] D. Soong *et al.*, “Improving Accuracy of GPT-3/4 Results on Biomedical Data Using a Retrieval-
1226 Augmented Language Model,” *PLOS Digit. Health*, vol. 3, no. 8, p. e0000568, 2024, doi:
1227 10.1371/journal.pdig.0000568.
- 1228 [205] N. Rekabsaz, O. Lesota, M. Schedl, J. Brassey, and C. Eickhoff, “TripClick: The Log Files of a Large
1229 Health Web Search Engine,” in *Proceedings of the 44th International ACM SIGIR Conference on Research*
1230 *and Development in Information Retrieval (SIGIR '21)*, New York, NY, USA: Association for Computing
1231 Machinery, 2021, pp. 2507–2513. doi: 10.1145/3404835.3463242.
- 1232 [206] Centers for Medicare & Medicaid Services, “ICD-10-PCS Official Guidelines for Coding and Reporting
1233 2023.” 2023. [Online]. Available: [https://www.cms.gov/files/document/2023-official-icd-10-pcs-coding-
1234 guidelines.pdf](https://www.cms.gov/files/document/2023-official-icd-10-pcs-coding-guidelines.pdf)
- 1235 [207] A. Albayrak, Y. Xiao, P. Mukherjee, S. S. Barnett, C. A. Marcou, and S. N. Hart, “Enhancing Human
1236 Phenotype Ontology Term Extraction Through Synthetic Case Reports and Embedding-Based Retrieval: A
1237 Novel Approach for Improved Biomedical Data Annotation,” *J. Pathol. Inform.*, vol. 16, p. 100409, Nov. 2024,
1238 doi: 10.1016/j.jpi.2024.100409.
- 1239 [208] J. D. Romano *et al.*, “The Alzheimer’s Knowledge Base: A Knowledge Graph for Alzheimer Disease
1240 Research,” *J. Med. Internet Res.*, vol. 26, no. 1, p. e46777, Apr. 2024, doi: 10.2196/46777.
- 1241 [209] D. S. Wishart *et al.*, “DrugBank: A comprehensive resource for in silico drug discovery and exploration,”
1242 *Nucleic Acids Res.*, vol. 34, no. suppl 1, pp. D668–D672, 2006, doi: 10.1093/nar/gkj067.
- 1243 [210] I. Segura-Bedmar, P. Martínez Fernández, and M. Herrero Zazo, “Semeval-2013 Task 9: Extraction of
1244 Drug-Drug Interactions from Biomedical Texts (DDIExtraction 2013),” in *Proceedings of the Seventh*
1245 *International Workshop on Semantic Evaluation (SemEval 2013)*, Atlanta, Georgia, USA: Association for
1246 Computational Linguistics, 2013, pp. 341–350.
- 1247 [211] R. I. Doğan, R. Leaman, and Z. Lu, “NCBI disease corpus: a resource for disease name recognition and
1248 concept normalization,” *J. Biomed. Inform.*, vol. 47, pp. 1–10, 2014, doi: 10.1016/j.jbi.2013.12.006.
- 1249 [212] J. Li *et al.*, “BioCreative V CDR Task Corpus: A resource for chemical disease relation extraction,”
1250 *Database*, vol. 2016, 2016, doi: 10.1093/database/baw068.
- 1251 [213] C. Eickhoff, F. Gmehlin, A. V. Patel, J. Boullier, and H. Fraser, “DC3 – A Diagnostic Case Challenge
1252 Collection for Clinical Decision Support,” in *Proceedings of the 2019 ACM SIGIR International Conference*

1253 *on Theory of Information Retrieval (ICTIR '19)*, Santa Clara, CA, USA: Association for Computing Machinery,
1254 2019. doi: 10.1145/3341981.3344239.

1255 [214] A. E. Johnson *et al.*, “MIMIC-CXR, a de-identified publicly available database of chest radiographs with
1256 free-text reports,” *Sci. Data*, vol. 6, no. 1, p. 317, 2019.

1257 [215] A. B. Abacha and D. Demner-Fushman, “A question-entailment approach to question answering,” *BMC*
1258 *Bioinformatics*, vol. 20, pp. 1–23, 2019.

1259 [216] V. Liévin, C. E. Hother, A. G. Motzfeldt, and O. Winther, “Can large language models reason about medical
1260 questions?,” *Patterns N*, vol. 5, no. 3, p. 100943, Mar. 2024, doi: 10.1016/j.patter.2024.100943.

1261 [217] A. Pal, L. K. Umapathi, and M. Sankarasubbu, “MedMCQA: A large-scale multi-subject multi-choice
1262 dataset for medical domain question answering,” in *Proceedings of the Conference on Health, Inference, and*
1263 *Learning*, PMLR, 2022, pp. 248–260.

1264 [218] Y. Li, Z. Li, K. Zhang, R. Dan, S. Jiang, and Y. Zhang, “ChatDoctor: A Medical Chat Model Fine-Tuned
1265 on a Large Language Model Meta-AI (LLaMA) Using Medical Domain Knowledge,” *Cureus*, vol. 15, no. 6,
1266 2023.

1267 [219] W. Hou and Z. Ji, “GeneTuring Tests GPT Models in Genomics.” 2023. [Online]. Available:
1268 <https://doi.org/10.1101/2023.03.11.532238>

1269 [220] M. Li, M. Chen, H. Zhou, and R. Zhang, “Petaylor: Improving large language model by tailored chunk
1270 scorer in biomedical triple extraction,” *ArXiv Prepr. ArXiv231018463*, 2023, [Online]. Available:
1271 <https://arxiv.org/abs/2310.18463>

1272

1273

1274 **Appendix**1275 **Table A1:** Classification of dense retriever types in biomedical RAG systems.

Method	Retriever	Task	Year
Type 1: PLMs			
AskFDALabel [78]	Sentence Transformer [23]	FDA drug labeling extraction	2024
CaLM [174]	BGE-embedding [24] with Chroma DB	supporting caregivers FM	2024
RAG-HPO [45]	FastEmbed [46]	rare genetic disorders automated deep phenotyping	2024
RALL [148]	Dense Passage Retriever [25]	lay language generation	2024
RT [47]	BERT embeddings [48]	few-shot medical NER	2024
RAG-GPT [49]	BGE-embedding [24]	breast cancer nursing care QA	2024
GNQA [50]	HNSW graphs [51]	medical QA	2024
Klang <i>et al.</i> [52]	GIST Large Embedding [41] with FAISS	ICD-10-CM coding	2024
RAGPR [53]	SBERT embeddings [22]	personalized physician recommendations	2025
JGCLLM [54]	GLuCoSE-base-ja vectors [55][genetic counseling support	2025
Type 2: Commercial APIs			
DRAGON-AI [56]	OpenAI Text-Embedding [57]	ontology generation with Chroma DB	2024
RISE [58]	OpenAI Text-Embedding with FAISS	diabetes-related inquiries	2024
Dual RAG [28]	UPstage API and OpenAI Text Embedding	diabetes management	2024
FAVOR-GPT [59]	Open AI Text-Embedding with Weaviate DB	genome variant annotations	2024
Endo-chat [60]	Open AI Text-Embedding with Faiss	medical QA for gastrointestinal endoscopy	2024
ChatENT [61]	Open AI Text-Embedding	medical QA in otolaryngology	2024
GuideGPT [62]	Open AI Text-Embedding	MRONJ QA on prevention, diagnosis, and treatment	2024
RECTIFIER [63]	Open AI Text-Embedding with Faiss	clinical trial screening for heart failure patients	2024
RAP [64]	Amazon Titan Text Embeddings v2 [65]	nutrition-related question answering	2025
InfectA-Chat [66]	Open AI Text-Embedding	infectious disease monitoring and QA	2025
DCRAG [67]	OpenAI Text-Embedding	auto-annotation of plant phenotype	2025
Type 3: Domain adaptation PLMs			
PM-Search [175]	BioBERT [68]	clinical literature retrieval	2022
LADER [69]	PubMedBERT [70]	biomedical literature retrieval	2023
CLEAR [71]	BioBERT	clinical information extraction	2025

KG-RAG [72]	PubMedBERT	biomedical multiple-choice questions	2024
WeiseEule [73]	MedCPT, BioBERT, BioGPT [74] with Pinecone DB	biomedical QA	2024
Self-BioRAG [13]	MedCPT [12]	medical QA	2024
CliniqIR [30]	MedCPT and BM25	diagnostic decision support	2024
MEDRAG [73]	BM25, SPECTER [176], Contriever [34] and MedCPT	biomedical RAG Tool	2024
RAMIE [75]	MedCPT, Contriever and BMR retriever [32]	biomedical IR about dietary supplements	2025

Type 4: LLM-based embedding

BiomedRAG [5]	MedLLaMA 13b	biomedical NLP tasks	2024
SurgeryLLM [151]	GPT4All [76] with Chroma DB	surgical decision support	2024
Myers <i>et al.</i> [77]	LLM2Vec (for comparison) [26]	clinical information retrieval	2024

1276

Table A2: Taxonomic classification of LLMs for embedding applications.

Model	Parameters	Architecture	Studies
Open-Source LLMs			
T5 [91]	0.06-11B	encoder-decoder	VAIV [31], CLEAR [71], DR.KNOWS [101]
ChatGLM3 [177]	6B	decoder-only	COPD [178], ChatZOC [179]
LLaMA 2 [88]	7-70B	decoder-only	Guo <i>et al.</i> [148], KREIMEYER <i>et al.</i> [180], CaLM [174], MMRAG [181], Yu <i>et al.</i> [182], RAMIE [75], EyeGPT [162], KG-RAG [72], JMLR [3], BiomedRAG [5], Self-BioRAG [13], CaLM [174], Du <i>et al.</i> [183], SurgeryLLM [151], Kreimeyer <i>et al.</i> [180]
LLaMA 3 [165]	8-70B	decoder-only	MMRAG [181], Fatharanihttps <i>et al.</i> [184], RAMIE [75], Woo <i>et al.</i> [185], RAGHPO [45], i-MedRAG [131], Hewitt <i>et al.</i> [95], PodGPT [128], SurgeryLLM [151]
Phi-3 Mini [186]	3.8B	decoder-only	Fatharani <i>et al.</i> [184]
Mistral 7B [90]	7B	decoder-only	Boulos <i>et al.</i> [187], RAMIE [75], Woo <i>et al.</i> [185], RAGPR [53], LITURAt [39], Kreimeyer <i>et al.</i> [180]
QWen [166]	32B	decoder-only	NEKO [108], Thetbanthad <i>et al.</i> [121]
Baichuan [188]	7-13B	decoder-only	ChatZOC [179]
Falcon [189]	7-180B	decoder-only	AskFDALabel [78], CaLM [174]
Zephyr [190]	7B	decoder-only	Moser <i>et al.</i> [127]
AceGPT [191]	7-13B	decoder-only	InfectA-Chat [66]
Gemma [192]	1-27B	decoder-only	PodGPT [128]
Deepseek R1 [171]	671B	decoder-only	Feng <i>et al.</i> [193]
Commercial LLMs			
ChatGPT3.5/4/4o [164]	Proprietary	decoder-only	Woo <i>et al.</i> [185], RAGPR [53], ChatZOC [179], KGRAG [72], Clinfo.ai [35], Yu <i>et al.</i> [182], ChatENT [61], CaLM [174], Rau <i>et al.</i> [194], Endo-chat [60], BiomedRAG [5], Lu <i>et al.</i> [37], Tan <i>et al.</i> [195], GNQA [50], Puts <i>et al.</i> [196], DRAGON-AI [56], Choi <i>et al.</i> [197], RISE [58], Du <i>et al.</i> [183], i-MedRAG [131], Hewitt <i>et al.</i> [95], RECTIFIER [63], Kainer <i>et al.</i> [67], Gong <i>et al.</i> [21], VAIA [31], RareDxGPT [155], InfectA-Chat, [66], FAVOR-GPT [59], Markey <i>et al.</i> [198]
Claude-3.5 [199]	Proprietary	decoder-only	Woo <i>et al.</i> [185], RISE [58], Hewitt <i>et al.</i> [95]
Gemini [89]	Proprietary	multimodal	Upadhyaya <i>et al.</i> [118], Tozuka <i>et al.</i> [119], MEREDITH [94]

Biomedical Specialized LLMs			
MedAlpaca [167]	7-13B	LLaMA-based	RAMIE [75]
PMC-LLaMA [27]	7B	LLaMA-based	RAMIE [75], BiomedRAG [5], Ozaki <i>et al.</i> [97]
BioMistral [96]	7B	Mistral-based	RAMIE [75], pRAGe [200]
MEDITRON [201]	7-70B	LLaMA-based	Ozaki <i>et al.</i> [97], Alkaeed <i>et al.</i> [98]
Meerkat [202]	7B	Mistral-based	Sohn <i>et al.</i> [203]

1278

1279

1280 **Table A3:** Comprehensive biomedical knowledge sources.

Source	Type	Description	Studies
MIMIC-IV [126]	EHR	decade of hospital admissions between 2008 and 2019 with data from about 300K patients and 430K admissions	Moser <i>et al.</i> [127]
PubMed	literature	repository of over 36 million biomedical research citations and abstracts, containing about 4.5B words	Clinifo.ai [35], MEDRAG [7], Self-BioRAG [13], JMLR [3], Ascle [105], RISE [58], MEREDITH [94], LITURAt [39], DR.KNOWS [101], Aftab [73], CliniqIR [30], PM-Search [108], Kreimeyer [180], Kim [31]
PMC	literature	8 million full-text biomedical and life sciences articles with free access, about 13.5B words	Self-BioRAG [13], Soong <i>et al.</i> [204], PodGPT [128]
TripClick [205]	literature	biomedical literature search and retrieval dataset	LADER [69]
GNQA [50]	literature	3000 peer reviewed publications on aging, dementia, Alzheimer’s and diabetes	GNQA [50]
UMLS [125]	knowledge base	2 million entities for 900K biomedical concepts	BiomedRAG [5], CLEAR [71], Guo <i>et al.</i> [148], CliniqIR [30]
ICD-10 [206]	knowledge base	international statistical classification of diseases and related health problems	Puts <i>et al.</i> [196]
HPO [154]	knowledge base	human phenotype ontology	Albayrak <i>et al.</i> [207], RAG-HPO [45], Kainer [67]
MeSH	taxonomy	hierarchical organization of biomedical and health-related topics	BMRETRIEVER [32]
AlzKB [208]	KG	knowledge graph for Alzheimer’s disease research	ESCARGOT [102], KRAGEN [104]
DrugBank [209]	database	comprehensive database integrating detailed drug data with drug target information	BMRETRIEVER [32], Kim [31]
SPOKE [109]	KG	integrating more than 40 publicly available biomedical knowledge sources of separate domains	KG-RAG [72]

COD [179]	knowledge base	comprehensive ophthalmic dataset	ChatZOC [179]
HEALIE KG [103]	KG	drawing from various medical ontologies, resources, and insights from domain experts	KAKALOU <i>et al.</i> [103]
ADRD [174]	corpus	collection of resources supporting Alzheimer’s disease caregivers	CaLM [174]
StatPearls [129]	clinical guide	collection of 9,330 clinical decision support articles available through NCBI Bookshelf	MEDRAG [7], i-MedRAG [131]
Medical Textbooks [130]	Educational	18 core medical textbooks commonly used in USMLE preparation	MEDRAG [7], Self-BioRAG [13], JMLR [3], i-MedRAG [131]
Ophthalmology textbooks [162]	educational	14 specialized ophthalmology textbooks	EyeGPT [162]
CheXpert [132]	image-report	224,316 chest radiographs with associated reports	CLEAR [71]

1281

1282

1283 **Table A4:** Comprehensive biomedical tasks datasets.

Dataset	Task	Studies	Year
ChemProt [136]	information extraction	BiomedRAG [5], ChemProt	2010
DDI [210]	information extraction	BiomedRAG [5], ChemProt	2013
NCBI [211]	entity recognition	GeneGPT [15], NetMe [100], RT [47]	2014
BioASQ [135]	information retrieval	BMRETRIEVER [32], SeRTS [33]	2015
MIMIC-III [134]	text summarization	DR.KNOWS [101], CliniqIR [30]	2016
BC5CDR [212]	entity recognition	RT [47]	2016
DC3 [213]	diagnostic classification	CliniqIR [30]	2019
MIMIC-CXR [214]	text summarization	Ascle [105]	2019
MedQuAD [215]	QA	Ascle [105]	2019
PubMedQA [133]	multiple-choice	BMRETRIEVER [32], Liévin <i>et al.</i> [216], PodGPT [128]	2019
MMLU [137]	multiple-choice	Self-BioRAG [13], JMLR [3], i-MedRAG [131], InfectA-Chat [66], PodGPT [128]	2021
MedQA [130]	multiple-choice	Self-BioRAG [13], JMLR [3], Ascle [105], Liévin <i>et al.</i> [216], i-MedRAG [131], PodGPT [128]	2021
MedMCQA [217]	multiple-choice	Self-BioRAG [13], JMLR [3], Ascle [105], Liévin <i>et al.</i> [216], EyeGPT [162], PodGPT [128]	2022
ChatDoctor [218]	dialogue	EyeGPT [162]	2023
MedAlpaca [167]	dialogue	EyeGPT [162]	2023
GeneTuring [219]	genomics QA	GeneGPT [15]	2023
GIT [220]	information extraction	BiomedRAG [5]	2023
MIRAGE [7]	multiple-choice & QA	MEDRAG [7]	2024
MedExpQA [29]	multilingual medical QA	MedExpQA [29], PodGPT [128]	2024
PubMedRS-200 [35]	medical QA	Clinifo.ai [35]	2024
MedicineQA [124]	multi-round dialogue	RagPULSE [124]	2024
CELLS [148]	lay language generation	Guo <i>et al.</i> [148]	2024

1284