

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

Jiawei He <sup>1,\*</sup> Boya Zhang <sup>2</sup> Hossein Rouhizadeh <sup>2</sup> Yingjian Chen <sup>3</sup> Rui Yang <sup>4</sup> Jin Lu <sup>5</sup> Xudong Chen <sup>1</sup> Nan Liu <sup>4</sup> Irene Li <sup>3</sup> Douglas Teodoro <sup>2,\*</sup>

<sup>1</sup> College of Information and Electronic Engineering, Hunan City University, Yiyang, Hunan, China

<sup>2</sup> Department of Radiology and Medical Informatics, University of Geneva, Geneva, Switzerland

<sup>3</sup> Graduate School of Engineering, The University of Tokyo, Tokyo, Japan

<sup>4</sup> Centre for Quantitative Medicine, Duke-NUS Medical School, Singapore

<sup>5</sup> LvZhiDao Information Technology Co., Ltd., Changsha, Hunan, China

\*Correspondence: [joewellhe@gmail.com](mailto:joewellhe@gmail.com), [douglas.teodoro@unige.ch](mailto:douglas.teodoro@unige.ch)

## ABSTRACT

Recent advances in large language models (LLMs) have demonstrated remarkable capabilities in natural language processing tasks. However, their application in the biomedical domain presents unique challenges, particularly regarding factual accuracy and up-to-date knowledge integration. Retrieval Augmented Generation (RAG) has emerged as a promising solution to address these challenges by combining the generative capabilities of LLMs with external knowledge retrieval. This comprehensive survey examines the application of RAG in the biomedical domain, focusing on its technological components, available datasets, and clinical applications. We present a systematic analysis of retrieval methods, ranking strategies, and generation models, while also exploring the challenges and future directions in this rapidly evolving field. Our work provides researchers and practitioners with a thorough understanding of the current state of biomedical RAG systems and identifies key areas for future research and development.

## KEYWORDS

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

## INTRODUCTION

Recent years have witnessed remarkable advancements in large language models (LLMs), resulting in a profound impact on various natural language processing (NLP) tasks<sup>1,2</sup>. These models have demonstrated unprecedented capabilities in generating and understanding human-like text, achieving outstanding performance across NLP tasks, such as summarization, question answering, and information retrieval. The exceptional performance of LLMs in core NLP tasks has catalyzed their exploration in the biomedical domain, where they show promise in supporting clinical decision-making, enhancing patient care quality, and improving clinical outcomes<sup>3</sup>.

However, the application of LLMs in the biomedical domain faces two critical challenges. First, these models often generate plausible-sounding but factually incorrect responses, a phenomenon commonly known as hallucination<sup>4</sup>. This is particularly concerning in medical applications, where accuracy is paramount for patient safety. Second, the static nature of LLM training means that once the training process is complete, the model parameters remain fixed, resulting in an inability to access or incorporate up-to-date medical knowledge<sup>5</sup>. This limitation is especially problematic in the rapidly evolving field of medicine, where new research findings and treatment guidelines emerge regularly.

Retrieval Augmented Generation (RAG) has emerged as a promising solution to address these critical challenges<sup>6</sup>. By combining the generative capabilities of LLMs with dynamic access to external knowledge from trusted sources, RAG systems can: (1) provide transparent rationales for generated responses by explicitly referencing source documents; (2) maintain access to current biomedical knowledge through real-time information retrieval; (3) enhance the accuracy and reliability of medical information processing; and (4) reduce hallucination by grounding responses in the verified medical literature<sup>1,2,5</sup>.

The implementation of RAG in the biomedical domain presents unique challenges and opportunities. Biomedical information is characterized by complex terminology, specialized knowledge structures, and strict requirements for accuracy and reliability<sup>7</sup>. Fur-

thermore, the integration of RAG systems into clinical workflows requires careful consideration of factors such as information privacy, regulatory compliance, and clinical validation<sup>8,9</sup>.

Concurrent to our survey, several related surveys focused on biomedical-related LLMs. For example, Liu et al.<sup>10</sup> review medical LLMs and He et al.<sup>2</sup> discuss LLMs for healthcare. Similarly, Xiao<sup>5</sup> conducted a comprehensive overview of LLMs in medicine. Other discussed RAG systems in general domains, such as Fan et al.<sup>1</sup> and Gao et al.<sup>11</sup> focused on the architectures, training strategies, and application of RAG in general domains.

Our work differs from existing surveys in concentrating on the application of RAG in the biomedical domain, systematically reviewing its components, datasets, and challenges. Overall, this survey provides the following key contributions:

1. A systematic analysis of RAG technologies specifically adapted for biomedical applications, including specialized retrieval methods, ranking algorithms, and generation models.
2. A detailed examination of medical datasets, considering both knowledge source datasets and medical task datasets.
3. An exploration of current clinical applications and future directions, with particular attention to practical implementation challenges and emerging opportunities.
4. A structured framework for understanding the interplay between retrieval mechanisms, medical knowledge sources, and language models in biomedical RAG systems.

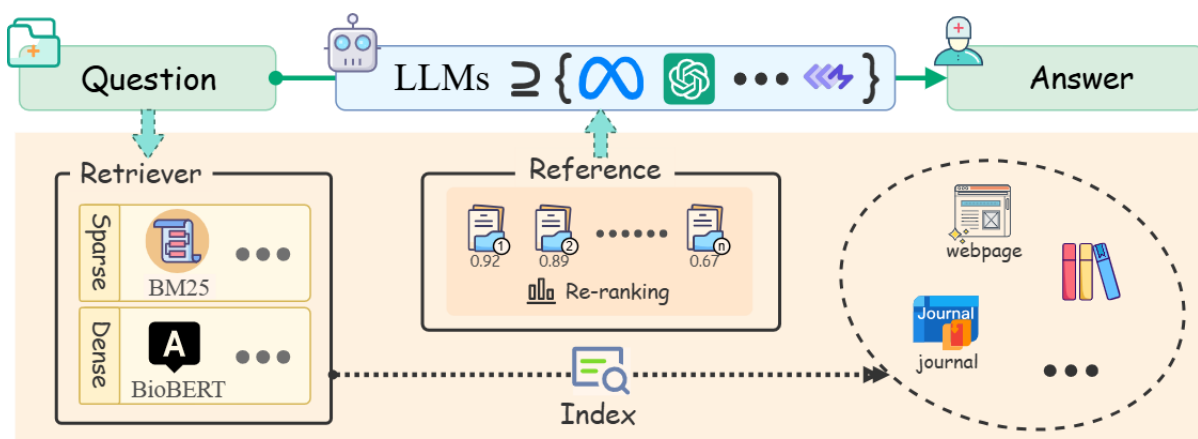
The remainder of this survey is organized as follows: Section 2 presents the background and the basic framework of biomedical RAG systems. Section 3 details the technical components, including retrieval methods, ranking algorithms, and generation models. Section 4 examines available datasets. Section 5 explores current clinical applications, and Section 6 discusses challenges and future directions.

## BACKGROUND

LLMs have disrupted NLP tasks with their ability to generate human-like text across diverse tasks<sup>12</sup>. Despite their impressive capabilities, LLMs often struggle with knowledge limitations, hallucinations, and outdated information, particularly in specialized domains requiring precise information<sup>13</sup>. RAG addresses these limitations by dynamically retrieving relevant external knowledge to supplement the LLM's parametric knowledge before generating responses<sup>11</sup>. This approach has shown particular promise in domain-specific applications, such as biomedicine, where accuracy, recency, and specialized knowledge are critical for clinical decision support, research analysis, and patient care<sup>5</sup>.

In a typical workflow, a query first passes through the retriever, which searches relevant documents from knowledge sources. The re-ranking method then sorts these references based on their relevance. Finally, the top-k most relevant references, combined with the original query, are processed by the LLM to generate the response. As Fig. 1 shows, the basic framework of biomedical RAG consists of four key components<sup>1</sup>: (1) **Retriever**: Responsible for indexing and retrieving relevant documents from various knowledge sources; Given a query  $q$  and a collection of knowledge sources  $K$ , a retriever function  $R(q, K)$  outputs a set of relevant documents  $D = \{d_1, d_2, \dots, d_n\}$ , where each  $d_i \in K$ ; (2) **Knowledge Source**:  $K$  represents a heterogeneous collection of structured and unstructured data, including but not limited to medical literature, clinical guidelines, and databases; (3) **Re-ranking Method**: Resorting and prioritizing retrieved references, A re-ranking function  $S(q, D)$  takes the query  $q$  and the retrieved document set  $D$  as input and outputs a re-ranked list of documents  $D' = [d'_1, d'_2, \dots, d'_m]$ , where  $D' \subseteq D$  and the order reflects relevance to  $q$ ; and (4) **LLM**: A LLM  $G(q, D')$  takes the original query  $q$  and the re-ranked retrieved context  $D'$  as input to generate a final response  $y$ .

RAG has been applied to a variety of biomedical applications. In clinical decision support systems, RAG powers evidence-based question-answering platforms like Clinfo.ai<sup>14</sup>, enhances specialized decision-making through expert-guided LLMs<sup>15</sup>, and facilitates rare disease diagnosis<sup>16</sup>. This technology also strengthens diagnostic prediction by connecting electronic health records (EHRs) with medical knowledge repositories<sup>17</sup> and advances genomic medicine through improved variant annotation and interpretation<sup>18</sup>. In the realm of medical imaging, RAG enhances report generation by seamlessly incorporating relevant prior studies and



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

clinical guidelines<sup>19</sup>. Across these applications, RAG leverages diverse data sources—including health records<sup>20,21</sup>, scientific literature<sup>22</sup>, clinical guidelines<sup>23</sup>, and genomic databases<sup>18</sup>—to generate evidence-based responses that significantly improve both clinical care and research outcomes. The growing body of evidence consistently demonstrates that RAG holds remarkable promise for advancing information management and decision support in the biomedical and healthcare domains.

## TECHNOLOGY

In this section, we explore the technological underpinnings of biomedical RAG systems, which integrate specialized components for retrieval, re-ranking, and generation to effectively process domain-specific information. Our examination of retrieval mechanisms encompasses a range of approaches, including sparse<sup>24,25</sup> and dense retrieval methods<sup>26,27</sup>, LLM-based retrieval<sup>4,28</sup>, knowledge graph enhancement techniques<sup>29,30</sup>, and emerging multi-modal retrieval capabilities designed to handle the diverse data formats prevalent in biomedical contexts<sup>19,31,32</sup>. We further analyze re-ranking methodologies that refine initial retrieval results to prioritize the most relevant documents<sup>33,34</sup>, help address the challenge of relevance and correctness in biomedical information retrieval, where accuracy directly impacts clinical decision-making and research outcomes. Finally, we discuss the generation component of biomedical RAG systems, comparing closed-source<sup>35–37</sup> and open-source models<sup>38–41</sup> alongside domain adaptation techniques<sup>6,42</sup> that enable these general-purpose language models to produce specialized biomedical content with appropriate terminology, and contextual understanding.

### Retriever

Retrievers are one of the core components of any RAG system, functioning as the mechanism that identifies and extracts relevant information from vast knowledge repositories in response to specific queries<sup>11</sup>. In the biomedical domain, retrievers face distinctive challenges due to the heterogeneous nature of data structures—ranging from densely technical academic papers and specialized textbooks to complex biomedical knowledge graphs containing intricate entity relationships<sup>43,44</sup>. The domain's specialized terminology, hierarchical concepts, and the life-critical nature of biomedical information make it more complex than information retrieval in general domains<sup>45</sup>. Effective biomedical retrievers must navigate these complexities while maintaining high precision to support accurate downstream analysis and generation. In this section, we describe various retrieval methodologies. We examine both standard retrieval approaches and how specific studies have modified these approaches for biomedical use, analyzing their respective strengths in handling domain-specific information structures and addressing unique retrieval challenges in the biomedical field.

#### Sparse Retrieval

Sparse retrieval techniques represent a foundational approach in information retrieval, particularly valuable in biomedical RAG systems which rely on terminology precision and exact keyword matching. These methods, exemplified by algorithms like BM25 (Best

Matching 25)<sup>46</sup> and TF-IDF (Term Frequency-Inverse Document Frequency)<sup>47</sup>, operate on the principle of lexical matching, where documents are indexed based on term occurrences rather than semantic meaning. Despite the rise of dense retrieval methods, sparse retrievers maintain their relevance due to their interpretability, computational efficiency, and effectiveness, when dealing with highly specialized vocabulary that characterizes biomedical literature.

Table 1 shows some recent publications that utilize sparse retrieval methods in biomedical RAG systems. These papers demonstrate various applications of sparse retrieval, from pure BM25<sup>48,49</sup> implementations to hybrid approaches<sup>50-53</sup> that combine sparse and dense retrieval techniques. Notably, sparse retrieval is frequently treated as either a baseline for comparison or strategically incorporated as a component in ensemble solutions to increase retrieval diversity but also to improve generalization<sup>54</sup>. Many of these works show that despite the attention given to newer neural approaches, traditional sparse retrieval methods remain competitive and often complementary to dense retrievers, particularly for specialized medical terminology and exact matching requirements<sup>55-57</sup>. Their continued relevance in state-of-the-art biomedical RAG systems highlights the value of integrating lexical matching with semantic understanding when dealing with the precise terminology and complex nomenclature characteristic of biomedical literature.

In this line of work, BM25 stands out as the most widely adopted algorithm, featuring in more than half the studies (n=9) either independently or within hybrid systems. Established search tools including Elasticsearch<sup>58</sup>, Whoosh<sup>59</sup>, and Entrez API<sup>60</sup> also prove popular, offering accessible implementation options without requiring extensive customization. A notable trend is the increasing development of hybrid approaches that combine BM25 with dense retrieval techniques<sup>50-53,61</sup>, leveraging both lexical precision and semantic understanding. The consistent appearance of sparse retrieval methods across all publication years, including the most recent 2025 studies, demonstrates their enduring value alongside newer neural approaches in cutting-edge biomedical information retrieval systems.

**Table 1:** Sparse retrieval methods in biomedical RAG systems

Method Name	Retrieval Method	Task	Year
Clinfo.ai <sup>14</sup>	Entrez API	Medical QA	2023
MEDRAG toolkit <sup>61</sup>	Hybrid (BM25, SPECTER, Contriever, MedCPT)	Medical QA	2024
SeRTS <sup>56</sup>	Self-Rewarding Tree Search based BM25	Biomedical QA (BioASQ-QA)	2024
Dual RAG System <sup>50</sup>	Hybrid (Dense + BM25 with language-specific tokenizers)	Diabetes management	2024
MedExpQA <sup>51</sup>	BM25 and MedCPT	Medical Question Answering	2024
CliniqIR <sup>52</sup>	BM25 and MedCPT	Diagnostic decision support	2024
VAIV Bio-Discovery <sup>53</sup>	Hybrid (Neural search + BM25)	Biomedical knowledge discovery	2024
Genomics analysis GPT <sup>18</sup>	Azure AI Search with keyword matching	Genetic variant annotation	2025
Customized GPT <sup>49</sup>	BM25 algorithm	Medical QA	2025
Two-Layer RAG <sup>48</sup>	Whoosh (Okapi BM25F ranking)	Medical question answering	2025
LITURat <sup>62</sup>	Entrez API for PubMed retrieval	Scientific data analysis	2025

Dense Retrieval

Dense retrievers represent a powerful approach to information retrieval that encodes queries and documents into continuous vector spaces to enable semantic matching based on learned representations rather than lexical overlap<sup>63</sup>. These retrievers have become important in the biomedical domain where terminology is complex and specialized, often requiring a nuanced understanding of context and meaning beyond simple keyword matching<sup>45</sup>. For instance, synonyms like 'cancer' and 'neoplasm' refer to the same concept, while hierarchical relationships exist between terms such as 'infectious disease' and its specific type 'COVID-19', further complicating accurate information retrieval. To analyze biomedical RAG systems, we categorize dense retrievers into five distinct types: Type 1: Pre-trained language models (PLMs); Type 2: Commercial embedding APIs; Type 3: Domain-adapted PLMs; Type 4: LLM-based embeddings; and Type 5: Domain-specific retrievers. In the following paragraphs, we provide a detailed analysis of each retriever type's characteristics, followed by a comprehensive review of recent works that implement these approaches in biomedical information retrieval systems.

**Type 1: PLMs** PLMs are trained on general corpora and are widely used in various information retrieval systems to create dense representations of queries and documents, so-called embeddings. There are numerous open-source general embedding models

available, which effectively capture semantic relationships and enable meaningful text matching, making them common foundational components in biomedical RAG systems. As shown in Table 2, transformer-based PLMs like BERT, SBERT<sup>64</sup>, and Sentence Transformer<sup>65</sup> have been widely adopted in recent biomedical retrieval systems. The continued popularity of these models is evident from numerous recent implementations. Additionally, newer embedding models such as BGE-embedding<sup>66</sup>, Dense Passage Retriever<sup>67</sup> and GIST Large Embedding<sup>68</sup> have gained significant adoption in recent literature, further diversifying the range of options available for biomedical retrieval tasks. Several other widely-used PLMs, including FastText<sup>69</sup>, RoBERTa<sup>70</sup>, and SimCSE<sup>71</sup>, though not specifically represented in Table 2, remain viable alternatives for developing effective biomedical RAG applications.

**Table 2:** Classification of dense retriever types in biomedical RAG systems

Method	Retriever	Task	Year
Type 1: PLMs			
AskFDALabel <sup>42</sup>	Sentence Transformer <sup>65</sup>	FDA drug labeling extraction	2024
CaLM <sup>26</sup>	BGE-embedding <sup>66</sup> with Chroma DB	Supporting caregivers FM	2024
RAG-HPO <sup>72</sup>	FastEmbed <sup>73</sup>	Rare genetic disorders automated deep pheno-typing	2024
RALL <sup>27</sup>	Dense Passage Retriever <sup>67</sup>	Lay language generation	2024
RT <sup>74</sup>	BERT embeddings <sup>75</sup>	Few-shot medical NER	2024
RAG-GPT <sup>76</sup>	BGE-embedding <sup>66</sup>	Breast cancer nursing care QA	2024
GNQA <sup>77</sup>	HNSW graphs <sup>78</sup>	Medical QA	2024
Klang <i>et al.</i> <sup>79</sup>	GIST Large Embedding <sup>68</sup> with FAISS	ICD-10-CM coding	2024
RAGPR <sup>80</sup>	SBERT embeddings <sup>64</sup>	Personalized physician recommendations	2025
JGCLLM <sup>81</sup>	GLuCoSE-base-ja vectors <sup>82</sup>	Genetic counseling support	2025
Type 2: Commercial APIs			
DRAGON-AI <sup>83</sup>	OpenAI Text-Embedding <sup>84</sup> with Chroma DB	Ontology generation	2024
RISE <sup>85</sup>	OpenAI Text-Embedding with FAISS	Diabetes-related inquiries	2024
Dual RAG <sup>86</sup>	UPstage API and OpenAI Text-Embedding	Diabetes management	2024
FAVOR-GPT <sup>87</sup>	Open AI Text-Embedding with Weaviate DB	Genome variant annotations	2024
Endo-chat <sup>88</sup>	Open AI Text-Embedding with Faiss	medical QA for gastrointestinal endoscopy	2024
ChatENT <sup>89</sup>	Open AI Text-Embedding	Medical QA in otolaryngology	2024
GuideGPT <sup>90</sup>	Open AI Text-Embedding	MRONJ QA on prevention, diagnosis, and treat-ment	2024
RECTIFIER <sup>91</sup>	Open AI Text-Embedding with Faiss	Clinical trial screening for heart failure patients	2024
RAP <sup>92</sup>	Amazon Titan Text Embeddings v2 <sup>93</sup>	Nutrition-related question answering	2025
InfectA-Chat <sup>94</sup>	Open AI Text-Embedding	Infectious disease monitoring and QA	2025
DCRAG <sup>95</sup>	OpenAI Text-Embedding	Auto-annotation of plant phenotype	2025
Type 3: Domain adaptation PLMs			
PM-Search <sup>24</sup>	BioBERT <sup>96</sup>	Clinical literature retrieval	2022
LADER <sup>97</sup>	PubMedBERT <sup>98</sup>	Biomedical literature retrieval	2023
CLEAR <sup>99</sup>	BioBERT	Clinical information extraction	2025
KG-RAG <sup>100</sup>	PubMedBERT	Biomedical Multiple-choice questions	2024
WeiseEule <sup>101</sup>	MedCPT, BioBERT, BioGPT <sup>102</sup> with Pinecone DB	Biomedical QA	2024
Type 4: LLM-based embedding			
BiomedRAG <sup>4</sup>	MedLLaMA 13b	Biomedical NLP tasks	2024
SurgeryLLM <sup>28</sup>	GPT4All <sup>103</sup> with Chroma DB	Surgical decision support	2024
Myers <i>et al.</i> <sup>104</sup>	LLM2Vec (for comparison) <sup>105</sup>	Clinical Information retrieval	2024
Type 5: Domain adaptation retriever			

Continued on next page

Table 2 – Continued from previous page

Method	Retriever	Task	Year
Self-BioRAG <sup>106</sup>	MedCPT <sup>33</sup>	Medical QA	2024
WeiseEule <sup>101</sup>	MedCPT, BioBERT, BioGPT with Pinecone DB	Biomedical QA	2024
CliniqIR <sup>52</sup>	MedCPT and BM25	Diagnostic decision support	2024
MEDRAG <sup>101</sup>	BM25, SPECTER <sup>22</sup> , Contriever <sup>57</sup> and MedCPT	Biomedical RAG Tool	2024
RAMIE <sup>107</sup>	MedCPT, Contriever and BMRetriever <sup>55</sup>	Biomedical IR about dietary supplements	2025

**Type 2: Commercial embedding APIs.** The proliferation of foundation models has led to the emergence of proprietary embedding solutions offered by major AI private organizations. These commercial embedding APIs have gained significant adoption in biomedical research, despite their associated costs, due to several compelling advantages: i) they typically leverage substantially larger pretraining datasets, ii) benefit from extensive computational resources during development, and iii) receive continuous refinements based on diverse user interactions. As shown in Table 2, OpenAI's embedding models dominate this category in recent studies, underscoring their enhanced semantic representation capabilities. While OpenAI's offerings dominate, alternatives are emerging. The RAP system<sup>92</sup> utilizes Amazon's Titan Text Embeddings v2<sup>93</sup> for nutrition-related question answering. Other commercial APIs such as Anthropic's Voyage-2<sup>108</sup> and Google's Vertex AI embeddings<sup>109</sup>, though not specifically represented in Table 2, are alternatives for biomedical RAG applications.

**Type 3: Domain adaptation PLMs.** Domain-adapted PLMs are language models that have been specifically pre-trained or fine-tuned on biomedical and clinical corpora, allowing them to develop specialized representations of medical terminology, concepts, and relationships<sup>96,98</sup>. These models mitigate a fundamental challenge in biomedical information retrieval: general-domain language models often struggle with the highly specialized vocabulary, complex semantic relationships, and unique linguistic patterns prevalent in biomedical literature<sup>110</sup>. By adapting pre-training to include PubMed articles, clinical notes, or other domain-specific corpus, these models develop a more nuanced understanding of biomedical language, resulting in significantly improved performance on retrieval tasks involving medical concepts and relationships compared to their general-domain counterparts<sup>110</sup>. As shown in Table 2, examples of these domain-adapted models include PubMedBERT<sup>98</sup>, BioBERT<sup>96</sup> and BioGPT<sup>102</sup>, which have been widely adopted in biomedical RAG systems. For example, PM-Search<sup>24</sup> utilized BioBERT for clinical literature retrieval, effectively leveraging the model's domain-specific representations to improve search precision. Similarly, PubMedBERT has been successfully deployed in LADER<sup>97</sup> for biomedical literature retrieval and in KG-RAG<sup>100</sup> for answering biomedical multiple-choice questions, demonstrating its versatility across different tasks. More recent systems like WeiseEule<sup>101</sup> combine multiple biomedical models including BioBERT, BioGPT, and MedCPT<sup>33</sup> to achieve enhanced performance in biomedical QA.

**Type 4: LLM-based embedding.** Generative LLMs have demonstrated remarkable generative capabilities across diverse NLP tasks, owing to their sophisticated understanding of language and context. With billions of parameters and training on vast corpora, these models develop rich internal representations that capture nuanced semantic relationships and contextual dependencies in text<sup>12</sup>. Recognizing this powerful language comprehension ability, researchers have begun leveraging open-source LLMs as embedding generators for retrieval tasks, extracting representations from their internal layers to create semantically meaningful document and query embeddings that outperform traditional embedding approaches. As shown in Table 2, notable implementations include BiomedRAG<sup>4</sup>, which utilizes MedLLaMA 13b<sup>111</sup>, SurgeryLLM<sup>28</sup>, which employs GPT4All<sup>103</sup>, and the comparative study by Myers *et al.*<sup>104</sup>, which evaluates LLM2Vec<sup>105</sup> for clinical information retrieval applications. However, most LLMs are decoder-only architectures that employ causal attention mechanisms, which inherently limit their ability to generate rich contextualized representations. Thus, researchers are increasingly focusing on developing more advanced LLM-based embedding methods, such as Llama2Vec<sup>112</sup>, Landmark Embedding<sup>113</sup>, which provide promising dense retrieval alternatives in biomedical RAG.

**Type 5: Domain adaptation retrievers.** Although the aforementioned off-the-shelf retrievers are readily available, the task of searching for relevant and accurate documents in the medical domain remains challenging. Consequently, customized retrievers tailored specifically to scientific and biomedical tasks have been developed. As shown in Table 2, these advanced retrievers include Con-

triever<sup>57</sup>, MedCPT<sup>33</sup>, SPECTER<sup>22</sup>, BMRetriever<sup>55</sup>, and the MEDRAG toolkit<sup>4</sup>. These specialized retrieval architectures demonstrate the importance of domain adaptation in biomedical information retrieval, significantly improving retrieval performance in several studies<sup>52,106,107</sup>. For example, Jeong *et al.*<sup>106</sup> proposed Self-BioRAG, which leverages MedCPT to improve medical question answering through a self-augmentation approach. Similarly, the MEDRAG framework<sup>101</sup> combines multiple retrieval strategies, including BM25, SPECTER, Contriever, and MedCPT to create a comprehensive biomedical RAG tool. RAMIE utilizes MedCPT, Contriever, and BMRetriever as retrievers, builds a RAG framework to extract diverse types of information about dietary supplements (DSs) from clinical records, enhanced overall accuracy in its experiments<sup>107</sup>. These implementations highlight how domain-adapted retrievers benefit biomedical RAG systems.

**Dense Retriever Infrastructure.** Vector stores provide an infrastructure element for dense retrievers in semantic search systems<sup>114</sup>. These specialized databases efficiently index, store, and search high-dimensional vector embeddings, offering optimized methods like K-nearest neighbors (KNN) for RAG developers. Vector databases enable rapid retrieval of semantically similar documents at scale, significantly enhancing the performance and response time of biomedical RAG systems. The process of embedding, storing, and searching documents involves transforming text into vector representations, organizing them in specialized data structures, and efficiently retrieving the most relevant matches. Based on our survey of biomedical RAG publications presented in Table 2, we have compiled the most commonly used vector databases in Table 3. In Table 2, it can be observed that the most popular vector stores are FAISS and Chroma, which were utilized 4 times and 3 times respectively.

Table 3: Common vector stores in biomedical RAG

Store	Description
Chroma <sup>115</sup>	AI-native open-source vector database with built-in functionality
Faiss <sup>116</sup>	Library for efficient similarity search and clustering of dense vectors
Pinecone <sup>117</sup>	Production-ready vector database for similarity search at scale
Weaviate <sup>118</sup>	Graph-based vector search engine with semantic capabilities
HNSW <sup>78</sup>	Hierarchical navigable small world graphs for approximate nearest neighbor search

Knowledge Graph Enhanced Retrieval

Knowledge graph retrievers offer significant advantages in biomedical RAG by providing structured, interlinked domain knowledge that overcomes several limitations of traditional retrieval methods. Unlike vector-based approaches, knowledge graphs capture explicit relationships between biomedical entities, enabling more precise retrieval of contextually relevant information, enhanced reasoning capabilities, and improved explainability<sup>119</sup>. Table 4 presents knowledge graph-based biomedical RAG studies published between 2024 and 2025, outlining each system’s application domains, and knowledge graph source. As we can notice, the application areas vary significantly, including medical question answering<sup>100</sup>, knowledge extraction<sup>120</sup>, and diagnosis prediction<sup>17</sup>.

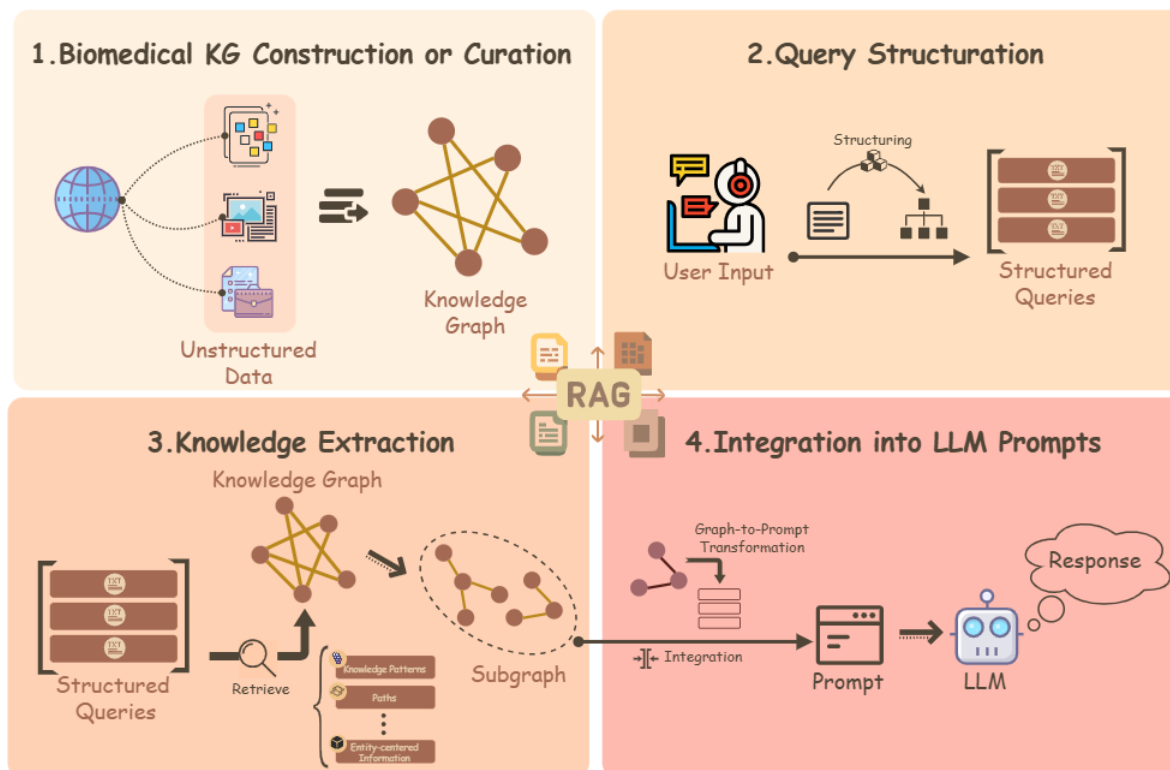
Table 4: Recent biomedical RAG studies based on knowledge graphs

Method	Task	Knowledge Graph Used	Year
KG-RAG <sup>100</sup>	Biomedical text generation, medical question answering	SPOKE biomedical knowledge graph	2024
HEALIE <sup>29</sup>	Personalized medical content generation	HEALIE KG (customized)	2024
KRAGEN <sup>121</sup>	Biomedical problem solving	Alzheimer’s knowledge graph (AlzKB)	2024
Ascle <sup>30</sup>	Various medical text generation tasks	UMLS	2024
DALK <sup>122</sup>	Alzheimer’s QA	Alzheimer’s knowledge graph (AlzKB)	2024
Gilbert <i>et al.</i> <sup>123</sup>	Medical information structuring and interlinking	Not specified	2024
NetMe 2.0 <sup>120</sup>	Biomedical knowledge extraction	BKGs constructed using OntoTagMe and Wikidata	2024
NEKO <sup>124</sup>	Knowledge mining in synthetic biology	PubMed-derived knowledge graphs	2025
DR.KNOWS <sup>17</sup>	Diagnosis prediction from EHRs	UMLS	2025
ESCARGOT <sup>125</sup>	Biomedical reasoning and knowledge retrieval	Alzheimer’s knowledge graph (AlzKB)	2025

As shown in Fig 2, the integration of knowledge graphs with LLMs in biomedical RAG applications can be analyzed through four key components:

- **Biomedical knowledge graph construction or curation** involves the systematic process of building, updating, or adapting





**Figure 2:** The Framework of Knowledge Graphs Based Biomedical RAG

knowledge graphs from domain-specific ontologies, literature, and other data sources to create structured biomedical knowledge repositories. NetMe 2.0<sup>120</sup> constructs biomedical knowledge graphs on-the-fly with OntoTagMe for entity linking while NEKO<sup>124</sup> mines PubMed literature for synthetic biology knowledge graphs. Differently, KG-RAG<sup>100</sup> utilizes the pre-existing SPOKE biomedical knowledge graph.

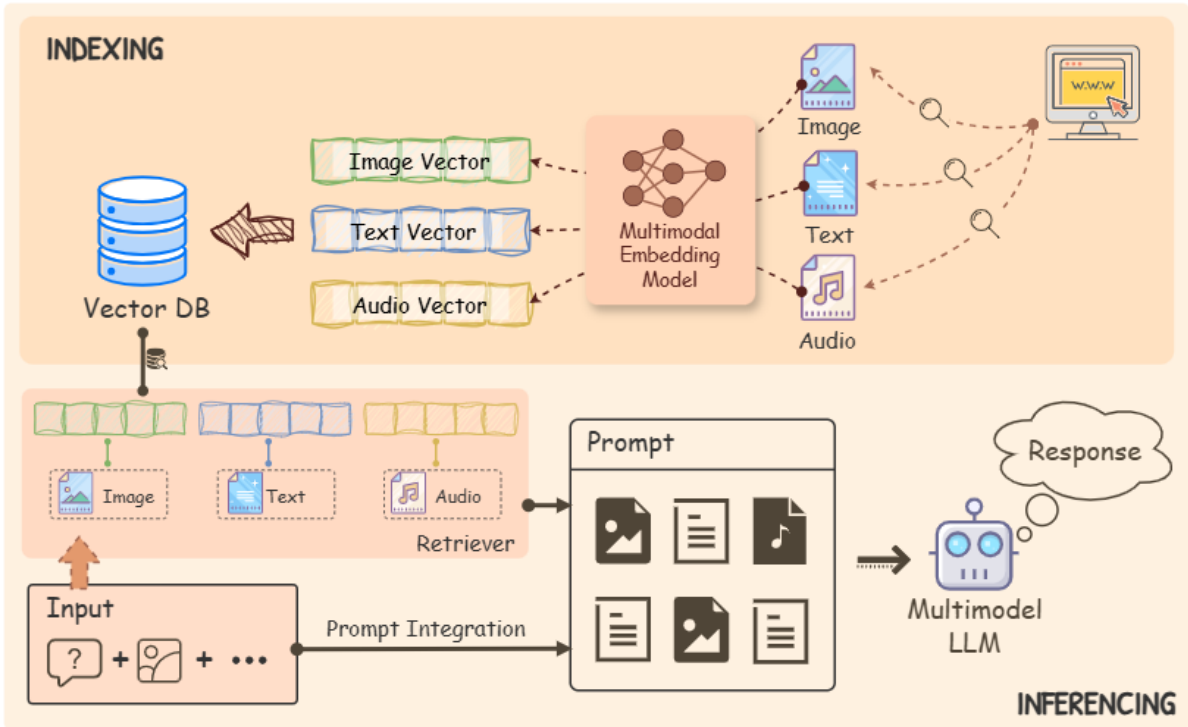
- **Query structuration** encompasses the conversion of natural language queries or user inputs into structured queries that efficiently retrieve relevant subgraphs from the knowledge graph. ESCARGOT<sup>125</sup> employs Cypher queries for precise retrieval, KRAGEN<sup>121</sup> vectorizes knowledge graphs for semantic searching, and DR.KNOWS<sup>17</sup> transforms queries to retrieve UMLS-based diagnostic pathways aligned with patient data.
- **Path-based or entity-centric knowledge extraction** refers to the process of extracting relevant knowledge patterns, paths, or entity-centered information from the knowledge graph based on the transformed query. DR.KNOWS<sup>17</sup> utilizes attention-based path ranking for diagnostic reasoning, KRAGEN<sup>121</sup> applies graph-of-thoughts for multi-hop relationships, and HEALIE<sup>29</sup> extracts paths incorporating both medical knowledge and socio-economic factors.
- **Strategic integration into LLM prompts** involves methods for incorporating retrieved knowledge into prompts for LLMs in ways that maximize factual accuracy while optimizing token usage and computational efficiency. KG-RAG<sup>100</sup> implements token-optimized retrieval to reduce consumption, ESCARGOT<sup>125</sup> selectively integrates knowledge through dynamic Graph of Thoughts, and NetMe 2.0<sup>120</sup> transforms graph data into natural language prompts via its Graph-RAG module.

### Multi-modal Retrieval

Multimodal biomedical RAG represents an advanced approach to healthcare AI that integrates LLMs with multiple types of data inputs beyond text, such as medical images, electrocardiogram (ECG), and other clinical measurements. These systems can enhance diagnostic accuracy and clinical decision support by leveraging the complementary strengths of different data modalities, retrieving relevant information from specialized knowledge bases, and generating comprehensive insights that would be impossible to



achieve with single-modality approaches<sup>126</sup>. As shown in Table 5, multimodal biomedical RAG systems have been successfully deployed for amblyopia detection<sup>32</sup>, prescription label interpretation<sup>127</sup>, lung cancer staging<sup>128</sup>, digital pathology analysis<sup>129</sup>, and automated radiology report generation<sup>31,130-133</sup>.



**Figure 3:** Architectural workflow of multimodal biomedical RAG system

**Table 5:** Multimodal Biomedical RAG Systems

Model	Task	Modalities	Year
Ranjit <i>et al.</i> <sup>130</sup>	Chest X-Ray report generation	Chest X-ray and text	2023
FactMM-RAG <sup>31</sup>	Radiology report generation	X-Ray and text	2024
MMed-RAG <sup>131</sup>	Medical VQA & Report generation	Medical images and text	2024
MEAG <sup>32</sup>	Amblyopia diagnosis	Eye tracking and text	2024
Tozuka <i>et al.</i> <sup>128</sup>	Lung cancer staging	CT findings and TNM classifications	2024
GPT4DFCI-RAG <sup>129</sup>	Digital pathology QA	Pathology images and text	2024
Raminedi <i>et al.</i> <sup>19</sup>	Radiology report generation	X-ray images and text	2024
RULE <sup>132</sup>	Medical VQA & Report generation	Medical images and text	2024
Thetbanthad <i>et al.</i> <sup>127</sup>	Prescription label identification	Image (labels) and text	2025
Hu <i>et al.</i> <sup>134</sup>	Pathology report generation	Whole slide images and text	2025
STREAM <sup>133</sup>	Chest X-ray report generation	Chest X-ray and text	2025

Figure 3 illustrates the architectural workflow of a multimodal biomedical RAG system. The key distinction between multimodal biomedical RAG and traditional biomedical RAG lies in data input diversity and processing architecture. While traditional RAG systems primarily process textual queries to retrieve text-based medical knowledge, multimodal RAG incorporates and processes non-textual data (images, sensor readings, etc.) alongside text, requiring specialized encoders for each modality to meaningfully integrate these diverse inputs<sup>126</sup>. The basic framework of multimodal biomedical RAG typically consists of: (1) **modality-specific encoders** (e.g., Vision Transformers<sup>135</sup> for images, specialized encoders for time-series data<sup>135</sup>); (2) **knowledge retrieval components** that access domain-specific databases based on multimodal query representations; and (3) **generation modules** that synthesize comprehensive outputs informed by both the multimodal inputs and retrieved knowledge. For instance, Ranjit *et al.*<sup>130</sup>, FactMM-RAG<sup>31</sup> and STREAM<sup>133</sup> use multimodal embeddings for retrieval of relevant candidate radiology text to generate Chest X-Ray report. MMed-RAG<sup>131</sup> and RULE<sup>132</sup> build a multimodal RAG system for Medical Visual Question Answering (Medical VQA) and report generation. The MEAG<sup>32</sup> system uses specialized encoders for eye movement data alongside text, Raminedi *et al.*<sup>129</sup> em-

ploy Vision Transformers as encoders for radiological images with cross-attention mechanisms for fusion with text, while Thetbanthad *et al.*<sup>127</sup> first converts prescription label images to text through OCR before retrieving relevant pharmaceutical information from DrugBank.

## Reranker

Rerankers in biomedical RAG systems serve as intermediary components that refine and prioritize the relevance of initially retrieved documents before they are processed by the generative models<sup>45</sup>. Positioned between retrievers and generators, rerankers apply sophisticated algorithms to re-score and re-order candidate documents based on their contextual relevance to the query, effectively functioning as a second-stage filtering mechanism that enhances the precision of information ultimately provided to the generation component. By prioritizing the most pertinent information, rerankers help to minimize hallucinations and improve the veracity of generated outputs.

Biomedical RAG rerankers can be categorized into distinct types based on their underlying methodologies. **Statistical rerankers** apply traditional information retrieval techniques; for instance, ClinIQIR<sup>52</sup> and MEDRAG<sup>101</sup> both implement Reciprocal Rank Fusion (RRF) to combine multiple retrieval signals, while DRAGON-AI<sup>83</sup> employs Maximal Marginal Relevance (MMR) to balance relevance with diversity. **Content-based rerankers**, in contrast, focus on content similarity; notably, WeiseEule<sup>101</sup> introduces a keyword frequency-based ranking system specifically designed for biomedical literature retrieval. **Model-based rerankers** employ transformer architectures to perform contextual relevance assessment; specifically, PodGPT<sup>34</sup> uses a general embedding model (bge-reranker-large) for scientific literature retrieval, whereas the CHA2DS2-VASc risk factor extraction system<sup>136</sup> utilizes a specialized cross-encoder model known as BGE M3-Embedding<sup>137</sup> to enhance clinical document relevance scoring. Within this category, some domain adaptation retrievers incorporate a reranker as a component of their model; for example, MedCPT<sup>33</sup> trains a specific reranker denoted as CrossEnc. Additionally, **LLM-based rerankers** leverage foundation models' reasoning capabilities; GNQA<sup>77</sup> employs GPT-3.5/4/4o as zero-shot binary classifiers for relevance determination, and LmRaC<sup>138</sup> implements an LLM-based paragraph usefulness assessment with a minimum threshold score of 7. (on a 10 point scale). Finally, **hybrid approaches** combine multiple strategies to enhance performance. For example, Clinfo.ai<sup>14</sup> combines LLM classification with BM25 to re-rank retrieved articles.

## Generation

The generation component represents the last element in biomedical RAG systems, responsible for synthesizing retrieved information into human-like responses that address user queries. In biomedical contexts, this process requires balancing natural language fluency with domain-specific accuracy, as generated content often informs healthcare decisions<sup>139</sup>. This section examines various generation models employed in recent biomedical RAG studies.

Current generation models employed in biomedical RAG systems can be categorized into three distinct classes based on their development approach and specialization: general-domain open-source models, commercial proprietary models, and biomedical-specialized models. These categories present different trade-offs in terms of accessibility, customization potential, and domain-specific performance. Table 6 showcases representative models from each category and highlights recent studies that utilize these models. It is important to note that our analysis of biomedical-specialized LLMs focuses specifically on models trained directly on biomedical or healthcare corpora, deliberately excluding applications developed through prompt engineering or Chain-of-Thought (CoT) methodologies.

As shown in Table 6, the majority of LLMs employ decoder-only architectures, with T5 being a notable exception. This architectural preference stems from several inherent advantages for generative tasks, particularly the natural suitability for next-token prediction and significantly improved parameter efficiency<sup>12</sup>. Among open-source LLMs, LLaMa 2<sup>140</sup> is the most popular, with 15 biomedical applications built upon it. LLaMa 3<sup>38</sup>, Mistral 7B<sup>141</sup>, and T5<sup>142</sup> are also widely used. These open-source LLMs have various sizes, from 0.06B (T5) to 671B (Deepseek R1), providing flexible options for building Biomedical RAG systems. Recently, several powerful general open-source LLMs have emerged, including, Mistral Small 3.1<sup>40</sup> and Gemma 3<sup>143</sup>. Although we have not found these newer models being used for biomedical RAG systems in our survey, we believe they show promise due to their outstanding performance in general domain tasks.

For commercial LLMs, the ChatGPT series APIs maintain a dominant position, presumably attributable to their superior performance. Nevertheless, several studies have demonstrated that Claude and Gemini achieve comparable efficacy when implemented in biomedical RAG applications<sup>15,85,144</sup>, thereby offering viable alternatives to ChatGPT. Furthermore, it is noteworthy that both DeepSeek and Mistral have established dual accessibility paradigms: providing open-source parameters for their base models while simultaneously offering commercial API services.

In contrast to open-source and commercial LLMs, biomedical specialized LLMs have been relatively underutilized in biomedical RAG applications. Recent studies<sup>107,145,146</sup> have demonstrated that within RAG frameworks, biomedical specialized LLMs do not consistently outperform open-source alternatives with statistical significance. Thus, the more effective methodologies for leveraging biomedical specialized LLMs within retrieval-augmented generation frameworks remain an area requiring substantial further investigation.

**Table 6:** Taxonomic Classification of LLMs for Embedding Applications

Model	Parameters	Architecture	Studies
Open-Source LLMs			
T5 <sup>142</sup>	0.06-11B	Encoder-Decoder	VAIV <sup>53</sup> , CLEAR <sup>99</sup> , DR.KNOWS <sup>17</sup>
ChatGLM3 <sup>147</sup>	6B	Decoder-only	COPD <sup>148</sup> , ChatZOC <sup>149</sup>
LLaMA 2 <sup>140</sup>	7-70B	Decoder-only	Guo et al. <sup>27</sup> , KREIMEYER et al. <sup>150</sup> , CaLM <sup>26</sup> , MMRAG <sup>151</sup> , Yu et al. <sup>152</sup> , RAMIE <sup>107</sup> , EyeGPT <sup>153</sup> , KG-RAG <sup>100</sup> , JMLR <sup>3</sup> , BiomedRAG <sup>4</sup> , Self-BioRAG <sup>106</sup> , CaLM <sup>26</sup> , Du et al. <sup>154</sup> , SurgeryLLM <sup>28</sup> , Kreimeyer et al. <sup>150</sup>
LLaMA 3 <sup>38</sup>	8-70B	Decoder-only	MMRAG <sup>151</sup> , Fatharanihttps et al. <sup>155</sup> , RAMIE <sup>107</sup> , Woo et al. <sup>156</sup> , RAG-HPO <sup>72</sup> , i-MedRAG <sup>157</sup> , Hewitt et al. <sup>144</sup> , PodGPT <sup>34</sup> , SurgeryLLM <sup>28</sup>
Phi-3 Mini <sup>158</sup>	3.8B	Decoder-only	Fatharanihttps et al. <sup>155</sup>
Mistral 7B <sup>141</sup>	7B	Decoder-only	Boulos et al. <sup>159</sup> , RAMIE <sup>107</sup> , Woo et al. <sup>156</sup> , RAGPR <sup>80</sup> , LITURAt <sup>62</sup> , Kreimeyer et al. <sup>150</sup>
QWen <sup>160</sup>	32B	Decoder-only	NEKO <sup>124</sup> , Thetbanthad <sup>127</sup>
Baichuan <sup>161</sup>	7-13B	Decoder-only	ChatZOC <sup>149</sup>
Falcon <sup>162</sup>	7-180B	Decoder-only	AskFDALabel <sup>42</sup> , CaLM <sup>26</sup>
Zephyr <sup>163</sup>	7B	Decoder-only	Moser et al. <sup>164</sup>
AceGPT <sup>165</sup>	7-13B	Decoder-only	InfectA-Chat <sup>94</sup>
Gemma <sup>166</sup>	1-27B	Decoder-only	PodGPT <sup>34</sup>
Deepseek R1 <sup>41</sup>	671B	Decoder-only	Feng et al. <sup>167</sup>
Commercial LLMs			
ChatGPT-3.5/4o <sup>35</sup>	Proprietary	Decoder-only	Woo et al. <sup>156</sup> , RAGPR <sup>80</sup> , ChatZOC <sup>149</sup> , KG-RAG <sup>100</sup> , Clinfo.ai <sup>14</sup> , Yu et al. <sup>152</sup> , ChatENT <sup>89</sup> , CaLM <sup>26</sup> , Rau et al. <sup>168</sup> , Endo-chat <sup>88</sup> , BiomedRAG <sup>4</sup> , Lu <sup>18</sup> , Tan <sup>169</sup> , GNQA <sup>77</sup> , Puts <sup>170</sup> , DRAGON-AI <sup>83</sup> , Choi <sup>171</sup> , RISE <sup>85</sup> , Du et al. <sup>154</sup> , i-MedRAG <sup>157</sup> , Hewitt et al. <sup>144</sup> , RECTIFIER <sup>91</sup> , Kainer et al. <sup>95</sup> , Gong et al. <sup>49</sup> , VAIA <sup>53</sup> , RareDxGPT <sup>16</sup> , Selcuk et al. <sup>94</sup> , FAVOR-GPT <sup>87</sup> , Markey et al. <sup>172</sup> ,
Claude-3.5 <sup>36</sup>	Proprietary	Decoder-only	Woo et al. <sup>156</sup> , RISE <sup>85</sup> , Hewitt et al. <sup>144</sup>

Table 6 continued from previous page

Model	Parameters	Architecture	Studies
Gemini <sup>37</sup>	Proprietary	Multimodal	Upadhyaya <i>et al.</i> <sup>32</sup> , Tozuka <i>et al.</i> <sup>128</sup> , MEREDITH <sup>15</sup>
<b>Biomedical Specialized LLMs</b>			
MedAlpaca <sup>173</sup>	7-13B	LLaMA-based	RAMIE <sup>107</sup>
PMC-LLaMA <sup>174</sup>	7B	LLaMA-based	RAMIE <sup>107</sup> , BiomedRAG <sup>4</sup> , Ozaki <i>et al.</i> <sup>145</sup>
BioMistral <sup>175</sup>	7B	Mistral-based	RAMIE <sup>107</sup> , pRAGe <sup>176</sup>
MEDITRON <sup>177</sup>	7-70B	LLaMA-based	Ozaki <i>et al.</i> <sup>145</sup> , Alkaeed <i>et al.</i> <sup>146</sup>
Meerkat <sup>178</sup>	7B	Mistral-based	Sohn <i>et al.</i> <sup>179</sup>

### Domain Adaptation

Domain adaptation in biomedical RAG systems constitutes a methodological framework that can be strategically implemented across various RAG pipeline stages to optimize performance for medical applications<sup>180</sup>. This approach facilitates enhanced performance on biomedical tasks while simultaneously preserving the fundamental capabilities inherent to foundation models. Based on our systematic analysis of the literature, we propose a taxonomy of domain adaptation techniques categorized according to their implementation within specific RAG pipeline components.

**Retriever Stage Adaptation.** Specialized retrieval components constitute the initial layer of domain adaptation, wherein techniques such as contrastive learning (implemented in Contriever<sup>57</sup> and MedCPT<sup>33</sup>) enable the development of more semantically rich biomedical embeddings. These adaptations augment the system's capacity to identify and retrieve pertinent medical documents through enhanced comprehension of domain-specific terminology and conceptual relationships. Additional retriever adaptation methodologies include continued pre-training on biomedical corpora (as demonstrated in BMRETRIEVER<sup>55</sup>) and domain-specific architectural modifications (exemplified by LADER<sup>97</sup> with PubMedBERT initialization). These approaches substantially enhance retrieval precision for specialized medical literature, effectively mitigating the lexical and semantic disparities between general and biomedical language domains.

**Reranker Stage Adaptation.** The reranking phase presents a subsequent opportunity for domain specialization, employing methodologies such as the development of dedicated biomedical rerankers (e.g., MedCPT's CrossEnc reranker<sup>33</sup>) and the implementation of domain-specific scoring mechanisms (e.g., BiomedRAG's chunk scorer<sup>4</sup>). Notably, the aforementioned **Model-based rerankers** and **LLM-based rerankers** (Sec ) can be conceptualized as domain adaptation strategies operationalized within the reranking stage. These specialized rerankers, trained with domain-specific datasets and optimization objectives tailored to biomedical document ranking, significantly enhance the system's capacity to prioritize documents based on domain-relevant criteria. Furthermore, they facilitate the identification of subtle domain-specific signals that general models typically fail to detect, thereby improving retrieval precision for specialized medical queries.

**LLM Stage Adaptation.** At the generation stage, adaptation of the language model itself through methodologies such as fine-tuning (MEDGENIE<sup>6</sup>, SeRTS<sup>56</sup>), parameter-efficient approaches like Low-Rank Adaptation (LoRA) (implemented in AskFDALabel<sup>42</sup>), and continued pre-training (RALL<sup>27</sup>, ChatENT<sup>89</sup>) enables more accurate synthesis and explication of biomedical information. JLMR<sup>3</sup> joint train LLM and retrieval model, which enhances biomedical RAG's ability to retrieve clinical guidelines and leverage medical knowledge to reason and answer questions. These adaptations enhance the model's comprehension of domain knowledge, facilitate the generation of appropriate medical terminology, and help improve factual accuracy in response generation.

This multi-level adaptation framework creates biomedical RAG systems that leverage the general capabilities of foundation models while incorporating the specialized knowledge representation and reasoning patterns essential for healthcare and life sciences applications. Through the strategic application of adaptation techniques at each stage of the RAG pipeline, these systems achieve

enhanced precision, improved contextual understanding, and more reliable generation of biomedical information, ultimately advancing the state-of-the-art in biomedical information retrieval and synthesis.

DATASETS

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

Knowledge Source Datasets

The effectiveness of biomedical RAG systems considerably depends on the quality, comprehensiveness, and diversity of their knowledge sources<sup>1,181</sup>. These knowledge repositories serve as the informational backbone from which these systems retrieve and synthesize medical content. High-quality knowledge sources can improve the generated responses, making them not only relevant but also reflect the current biomedical consensus and best practices<sup>89,106,182</sup>.

Biomedical knowledge spans multiple dimensions—from basic science research to clinical guidelines, from pharmaceutical data to standardized medical terminology. To address this complexity, effective biomedical RAG systems integrate diverse knowledge sources that complement each other in scope, specialization, and format. This integration enables systems to comprehensively address the multifaceted nature of medical queries<sup>55,101</sup>.

Table 7 shows the public knowledge source utilized in the surveyed studies. Among these sources, PubMed, UMLS<sup>183</sup>, PMC, and Medical textbooks<sup>184</sup> demonstrate predominant adoption, primarily attributable to their comprehensive coverage of biomedical concepts and extensive document collections. Additionally, we notice that several studies<sup>80,154,171</sup> have implemented biomedical RAG systems based on proprietary clinical data that is unpublished due to privacy concerns, while another subset of studies<sup>50,76,92,165,185,186</sup> focus on a specific medical specialty, wherein specialized medical guidelines about particular clinical specialties are leveraged for constructing a targeted knowledge source.

Table 7: Comprehensive Medical Knowledge Sources

Source	Type	Description	Studies
MIMIC-IV <sup>187</sup>	EHR	Decade of hospital admissions between 2008 and 2019 with data from about 300K patients and 430K admissions	Moser et al. <sup>164</sup>
PubMed	Literature	Repository of over 36 million biomedical research citations and abstracts, containing about 4.5B words	Clinifo.ai <sup>14</sup> , MEDRAG <sup>61</sup> , Self-BioRAG <sup>106</sup> , JMLR <sup>3</sup> , Ascle <sup>30</sup> , RISE <sup>85</sup> , MEREDITH <sup>15</sup> , LITURAt <sup>62</sup> , DR.KNOWS <sup>17</sup> , Aftab <sup>101</sup> , CliniqIR <sup>52</sup> , PM-Search <sup>24</sup> , Kreimeyer <sup>150</sup> , Kim <sup>53</sup>
PMC	Literature	8 million full-text biomedical and life sciences articles with free access, about 13.5B words	Self-BioRAG <sup>106</sup> , Soong et al. <sup>188</sup> , PodGPT <sup>34</sup>
TripClick <sup>189</sup>	Literature	Biomedical literature search and retrieval dataset	LADER <sup>97</sup>
GNQA <sup>77</sup>	Literature	3000 peer reviewed publications on on aging, dementia, Alzheimer’s and diabetes	GNQA <sup>77</sup>
UMLS <sup>183</sup>	Knowledge Base	2 million entities for 900K biomedical concepts	BiomedRAG <sup>4</sup> , CLEAR <sup>99</sup> , Guo <sup>27</sup> , CliniqIR <sup>52</sup>

Continued on next page

**Table 7 continued from previous page**

Source	Type	Description	Studies
ICD-10 <sup>190</sup>	Knowledge Base	International statistical classification of diseases and related health problems	Puts et al. <sup>170</sup>
HPO <sup>191</sup>	Knowledge Base	Human phenotype ontology	Albayrak et al. <sup>192</sup> , RAG-HPO <sup>72</sup> , Kainer <sup>95</sup>
MeSH	Taxonomy	Hierarchical organization of biomedical and health-related topics	BMRETRIEVER <sup>55</sup>
AlzKB <sup>125</sup>	Knowledge Graph	Knowledge graph for Alzheimer’s disease research	ESCARGOT <sup>125</sup> , KRAGEN <sup>121</sup>
DrugBank <sup>193</sup>	Database	Comprehensive database integrating detailed drug data with drug target information	BMRETRIEVER <sup>55</sup> , Kim <sup>53</sup>
SPOKE <sup>194</sup>	Knowledge Graph	Integrating more than 40 publicly available biomedical knowledge sources of separate domains	KG-RAG <sup>100</sup>
COD <sup>149</sup>	Knowledge Base	Comprehensive ophthalmic dataset	ChatZOC <sup>149</sup>
HEALIE KG <sup>29</sup>	Knowledge Graph	Drawing from various medical ontologies, resources, and insights from domain experts	KAKALOU et al. <sup>29</sup>
ADRD Care-giving <sup>26</sup>	Corpus	Collection of resources supporting Alzheimer’s disease caregivers	CaLM <sup>26</sup>
StatPearls <sup>23</sup>	Clinical Guide	Collection of 9,330 clinical decision support articles available through NCBI Bookshelf	MEDRAG <sup>61</sup> , i-MedRAG <sup>157</sup>
Medical Text-books <sup>184</sup>	Educational	18 core medical textbooks commonly used in USMLE preparation	MEDRAG <sup>61</sup> , Self-BioRAG <sup>106</sup> , JMLR <sup>3</sup> , i-MedRAG <sup>157</sup>
Ophthalmology text-books <sup>153</sup>	Educational	14 specialized ophthalmology textbooks	EyeGPT <sup>153</sup>
CheXpert <sup>195</sup>	Image-Report	224,316 chest radiographs with associated reports	CLEAR <sup>99</sup>

### Medical Task Datasets

Medical task datasets represent essential benchmarking instruments for evaluating RAG systems within healthcare applications<sup>2,10</sup>. These datasets simulate real clinical information challenges, providing structured frameworks to assess how systems process, interpret, and generate medical content.

These evaluation datasets can be systematically categorized according to distinct information-handling processes in the biomedical domain: biomedical information extraction, entity recognition, question answering (QA), biomedical multiple-choice examination, dialogue, and generation tasks. Table 8 presents a comprehensive overview of the evaluation datasets employed across the surveyed biomedical RAG systems. Apart from works that make use of public biomedical datasets, several research groups<sup>81,85,90</sup> in favor of human-curated questions, designing proprietary data for evaluation, and inviting healthcare professionals (e.g., physicians or nurses) to conduct assessments.

**Table 8:** Comprehensive Biomedical Tasks Datasets

Dataset	Year	Task	Studies
ChemProt <sup>196</sup>	2010	Information Extraction	BiomedRAG <sup>4</sup> , ChemProt
DDI <sup>197</sup>	2013	Information Extraction	BiomedRAG <sup>4</sup> , ChemProt
NCBI <sup>198</sup>	2014	Entity Recognition	GeneGPT <sup>199</sup> , NetMe <sup>120</sup> , RT <sup>74</sup>
BioASQ <sup>200</sup>	2015	Information Retrieval	BMRETRIEVER <sup>55</sup> , SeRTS <sup>56</sup>
MIMIC-III <sup>20</sup>	2016	Text Summarization	DR.KNOWS <sup>17</sup> , CliniqIR <sup>52</sup>
BC5CDR <sup>201</sup>	2016	Entity Recognition	RT <sup>74</sup>
DC3 <sup>202</sup>	2019	Diagnostic classification	CliniqIR <sup>52</sup>
MIMIC-CXR <sup>203</sup>	2019	Text Summarization	Ascle <sup>30</sup>
MedQuAD <sup>204</sup>	2019	QA	Ascle <sup>30</sup>
PubMedQA <sup>205</sup>	2019	Multiple-choice	BMRETRIEVER <sup>55</sup> , Liévin et al. <sup>206</sup> , PodGPT <sup>34</sup>
MMLU <sup>207</sup>	2021	Multiple-choice	Self-BioRAG <sup>106</sup> , JMLR <sup>3</sup> , i-MedRAG <sup>157</sup> , InfectA-Chat <sup>94</sup> , PodGPT <sup>34</sup>
MedQA <sup>184</sup>	2021	Multiple-choice	Self-BioRAG <sup>106</sup> , JMLR <sup>3</sup> , Ascle <sup>30</sup> , Liévin et al. <sup>206</sup> , i-MedRAG <sup>157</sup> , PodGPT <sup>34</sup>
HealthSearchQA <sup>208</sup>	2022	QA	
MedMCQA <sup>209</sup>	2022	Multiple-choice	Self-BioRAG <sup>106</sup> , JMLR <sup>3</sup> , Ascle <sup>30</sup> , Liévin et al. <sup>206</sup> , EyeGPT <sup>153</sup> , PodGPT <sup>34</sup>
ChatDoctor <sup>210</sup>	2023	Dialogue	EyeGPT <sup>153</sup>
MedAlpaca <sup>173</sup>	2023	Dialogue	EyeGPT <sup>153</sup>
GeneTuring <sup>211</sup>	2023	Genomics QA	GeneGPT <sup>199</sup>
GIT <sup>212</sup>	2023	Information Extraction	BiomedRAG <sup>4</sup>
MIRAGE <sup>61</sup>	2024	Multiple-choice & QA	MEDRAG <sup>61</sup>
MedExpQA <sup>51</sup>	2024	Multilingual Medical QA	MedExpQA <sup>51</sup> , PodGPT <sup>34</sup>
PubMedRS-200 <sup>14</sup>	2024	Medical QA	Clinifo.ai <sup>14</sup>
MedicineQA <sup>182</sup>	2024	Multi-round dialogue	RagPULSE <sup>182</sup>
CELLS <sup>27</sup>	2024	Lay language generation	Guo <sup>27</sup>

## APPLICATIONS

The application of RAG in the biomedical domain represents a significant advancement in healthcare informatics, combining the strengths of LLMs with domain-specific knowledge retrieval. This integration helps mitigate challenges inherent to healthcare applications—factual accuracy, domain specificity, knowledge recency, and explainability—that conventional generative AI approaches struggle to overcome<sup>213</sup>. In this section, we analyze key healthcare domains where biomedical RAG systems have demonstrated substantial utility and are organized according to their functional contributions to clinical practice, precision medicine, and healthcare knowledge dissemination.

### Clinical Decision Support Systems

#### Medical Question Answering

Medical question-answering systems provide clinicians and patients with access to accurate healthcare information at the point of need<sup>100</sup>. This application is essential because it enables rapid retrieval of evidence-based information, supporting clinical decision-making while reducing the cognitive load on healthcare providers and improving information access for patients with varying levels



of health literacy.

Recent advancements have produced several sophisticated medical QA systems with distinctive approaches to knowledge retrieval. Clinfo.ai<sup>14</sup> provides evidence-based healthcare QA by retrieving information from trusted sources like PubMed using the Entrez API, ensuring responses reflect current medical evidence. RALL<sup>27</sup> employs Dense Passage Retriever to generate lay language explanations of medical concepts, addressing the critical need for accessible medical information. Self-BioRAG<sup>106</sup> implements the MedCPT retriever with self-reflection techniques to improve accuracy for complex clinical questions with multiple components, while WeiseEule<sup>101</sup> integrates multiple retrievers (MedCPT, BioBERT, BioGPT) with Pinecone DB to handle diverse medical information needs across specialties. These medical QA systems collectively demonstrate RAG's potential for healthcare information access, providing contextually relevant, evidence-based responses to complex medical queries for both clinical decision support and patient education.

### **Diagnostic and Treatment Decision Support**

Diagnostic and treatment decision support systems assist clinicians in navigating complex clinical pathways through evidence retrieval and synthesis<sup>214</sup>. They can help healthcare providers identify relevant evidence, recognize patterns, and identify appropriate interventions, ultimately improving diagnostic accuracy and treatment selection while reducing cognitive biases and diagnostic errors.

Several innovative implementations showcase RAG's potential in clinical decision support across medical specialties. Expert-Guided LLMs<sup>15</sup> support specialized decision-making in precision oncology by retrieving relevant genomic information and matching it with treatment guidelines and clinical trial data. DR.KNOWS<sup>17</sup> facilitates diagnostic prediction by connecting EHRs with medical knowledge repositories through knowledge graph-enhanced retrieval that models complex relationships between symptoms, diseases, and treatments. CliniqIR<sup>52</sup> combines MedCPT and BM25 retrieval to match patient symptoms with potential diagnoses, while SurgeryLLM<sup>28</sup> delivers surgical decision support by retrieving relevant procedures, techniques, and complication data. RECTIFIER<sup>91</sup> employs OpenAI Text-Embedding with Faiss for clinical trial screening, helping match heart failure patients with appropriate clinical trials by analyzing eligibility criteria against patient data. By connecting patient-specific data with evidence-based medical knowledge, these RAG systems significantly enhance clinical decision-making while reducing the cognitive burden on healthcare providers, potentially improving both efficiency and quality of care.

### **Rare Disease Identification and Management**

Rare disease identification and management represents a particularly challenging healthcare domain due to limited literature, dispersed expertise, and complex symptom patterns<sup>215</sup>. RAG systems are especially valuable in this context because they can aggregate and synthesize knowledge from diverse sources, helping clinicians recognize unusual patterns and connect them with appropriate diagnostic and management strategies for conditions they may encounter only rarely in clinical practice.

Recent research has produced several specialized rare disease applications leveraging different knowledge retrieval approaches. Zelin et al.<sup>16</sup> developed a RAG system specifically focused on rare disease diagnosis, connecting uncommon symptom patterns with rare disease knowledge drawn from specialized databases. RAG-HPO<sup>72</sup> utilizes FastEmbed for automated deep phenotyping of rare genetic disorders, linking phenotypic observations with potential genetic causes through human phenotype ontology mapping. Yang et al.<sup>216</sup> created RDguru, an intelligent conversational agent that employs a multi-source fusion diagnostic model combining GPT-4, PheLR, and phenotype matching strategies to significantly improve diagnostic recall for rare diseases. Some authors<sup>217</sup> have also evaluated the capabilities of LLMs and RAG for phenotype-driven gene prioritization in rare genetic disease diagnosis. These applications demonstrate how RAG systems can effectively bridge knowledge gaps in rare disease management, connecting dispersed information to support timely diagnosis and appropriate treatment strategies for conditions where clinical expertise is often limited.

## Clinical Report Generation

Automatic and semi-automatic report generation can help healthcare providers by significantly reducing the time they spend on documentation, allowing them to dedicate more time to patient care<sup>218</sup>. A critical area of report generation is radiology, which involves image analysis, document consultation, and data evaluation<sup>219</sup>. This process faces multiple challenges including ensuring factual accuracy, maintaining clinical relevance, minimizing hallucinations, and providing sufficient interpretability for clinical adoption<sup>220</sup>.

RAG frameworks have been proposed to address these challenges by grounding language generation in relevant retrieved content, resulting in more accurate and trustworthy outputs<sup>7</sup>. Recent works have demonstrated the effectiveness of RAG in this domain. For example, multimodally-aligned embeddings with general-domain LLMs such as GPT-4<sup>221</sup>, ensuring relevant content retrieval and mitigating hallucinations while allowing customization based on clinical intent. To enhance interpretability, some authors<sup>222</sup> proposed a multi-agent RAG framework utilizing concept bottleneck models that guide report generation with clinically meaningful concept vectors, bridging the gap between performance and explainability. RAG pipelines<sup>223</sup> have also been used to improve image-text matching techniques, incorporating semantic segmentation and triple contrastive loss, leading to more accurate and fluent reports with reduced errors. Differently, LaB-RAG<sup>224</sup> employs label-boosted retrieval by converting medical images into descriptive text labels, which are then used for text-based retrieval and report generation, bypassing the need for direct image processing by the LLM. Lastly, some authors<sup>225</sup> have proposed fact-aware multimodal retrieval strategies to address factual inconsistencies in generated reports by leveraging RadGraph-based factual mining and multimodal retrievers, ensuring high-quality reference retrieval and improved factual completeness.

## Precision Medicine Applications

### Genomic Medicine

Genomic medicine involves analyzing genetic data to inform clinical decision-making. The rapid knowledge evolution and inherent complexity in variant interpretation within genomic medicine impose significant cognitive burdens on clinicians, particularly given the exponential expansion of genomic information and the challenges in ascertaining variant pathogenicity<sup>226</sup>.

Several innovative RAG applications address these challenges in genomic medicine through specialized knowledge retrieval. FAVOR-GPT<sup>87</sup> utilizes OpenAI Text-Embedding with Weaviate DB to facilitate genome variant annotations, connecting genetic variants with clinical significance and supporting evidence from multiple databases. Lu et al.<sup>18</sup> developed a genomics analysis system using Azure AI Search with keyword matching for more accurate genetic variant interpretation within clinical contexts. RAG-HPO<sup>72</sup> connects clinical observations with genetic variants in rare disorders through sophisticated knowledge mapping, while JGCLLM<sup>81</sup> implements GLuCoSE-base-ja vectors to provide genetic counseling support that translates complex genomic findings into clinically meaningful information. By connecting genomic findings with clinical knowledge through sophisticated retrieval mechanisms, these RAG systems significantly improve the clinical application of genetic information, such as disease inference<sup>18</sup> and phenotypic and genotypic analysis<sup>72</sup>, while supporting more personalized treatment decisions in the emerging era of genomic medicine.

### Personalized Health Management

Personalized health management, a tailored healthcare approach, involves generating individualized health information and recommendations based on patient-specific characteristics, preferences, and health status<sup>227</sup>. This application mitigates a important healthcare challenge: the need for more effective and patient-centered care that optimizes treatment adherence, self-management, and health outcomes, in contrast to traditional standardized approaches that often fail to accommodate specific patient needs, literacy levels, and medical conditions.

Recent research has produced several personalized health management systems employing different knowledge retrieval strategies. HEALIE<sup>29</sup> leverages knowledge graphs to generate personalized medical content incorporating both clinical knowledge and socio-economic factors relevant to individual patients. RISE<sup>85</sup> implements OpenAI Text-Embedding with FAISS to provide person-

alized responses to diabetes-related inquiries, helping patients navigate complex self-management requirements. A Dual RAG System<sup>50</sup> combines dense retrieval and BM25 with language-specific tokenizers to support diabetes management across multiple languages, expanding access to personalized health information for diverse populations. RAP<sup>92</sup> utilizes Amazon Titan Text Embeddings to provide personalized nutrition-related question answering, supporting dietary management for various health conditions. These personalized health management systems demonstrate RAG's ability to bridge the gap between complex medical information and individual patient needs through contextually relevant knowledge retrieval, ultimately supporting better self-management and improved health outcomes through personalized guidance.

## Healthcare Education

Medical education and training systems support the development of healthcare professionals through evidence-based learning resources and simulated clinical experiences<sup>228</sup>. These systems help address workforce needs by enabling more efficient knowledge acquisition, supporting the development of clinical reasoning skills, and allowing for personalized learning experiences that adapt to individual learner needs and knowledge gaps.

ChatENT<sup>89</sup> leverages Retrieval-Augmented Language Modeling (RALM) to support medical education in Otolaryngology-Head and Neck Surgery (OHNS). EyeGPT<sup>153</sup> transforms a general LLM into an ophthalmology-specialized assistant through prompt engineering, fine-tuning with eye-specific terminology, and RAG from ophthalmology textbooks, supporting both medical education and patient inquiries. Sevgi et al.<sup>186</sup> developed targeted applications in ophthalmology, including "EyeTeacher" for generating educational questions and "EyeAssistant" for clinical support, both enhanced through the integration of clinical guidelines and peer-reviewed documents for retrieval. In nursing education, Zhao et al.<sup>229</sup> leveraged RAG and semantic search technology to answer questions from the Chinese Nursing Licensing Exam, showing the potential for specialized knowledge retrieval in professional certification. By connecting learning objectives with relevant medical knowledge through sophisticated retrieval mechanisms, these RAG-based education systems serve as valuable educational aid tools for healthcare practice skill training.

## Clinical Research

Clinical research application tools are designed to support clinical studies by enhancing efficiency, leveraging knowledge retrieval and synthesis to accelerate the development of medical knowledge and therapeutic innovations.<sup>230</sup>

Recent literature highlights several innovative RAG applications in clinical research leveraging different retrieval approaches. PodGPT<sup>34</sup> implements the BAAI embedding model for scientific literature retrieval informing research design and methodology development. LITURAt<sup>62</sup> utilizes the Entrez API for PubMed retrieval, supporting scientific data analysis in clinical research through relevant literature identification. These applications showcase RAG's ability to optimize and improve the clinical research process through intelligent knowledge retrieval, highlighting its potential as a valuable tool for both research and education, and ultimately accelerating the translation of research findings into clinical practice.

## CHALLENGES AND FUTURE DIRECTION

### Trustworthiness

Recent research indicates that RAG systems can be both maliciously and unintentionally manipulated to produce unreliable decisions and potentially harm humans<sup>1,231</sup>. Such vulnerabilities are particularly unacceptable in the biomedical domain due to the potential for severe consequences.

Several dimensions of trustworthiness in biomedical RAG systems raise critical concerns. **Robustness** refers to the system's ability to exclude malicious content introduced by attackers or filter inadvertent misinformation contained in web corpora that might contaminate the knowledge source<sup>181</sup>. Although certain re-ranking methodologies described in Section can partially mitigate these issues, significant challenges remain in biomedical RAG applications. **Explainability** requires a comprehensive understanding of user intent to generate reliable responses aligned with ground-truth documents. A notable challenge emerges from diverse linguistic styles in input information, creating misalignment between user queries and knowledge sources. For instance, in medical QA systems, users may formulate queries in conversational language while knowledge sources are typically based on academic

publications<sup>139</sup>. This linguistic discrepancy poses substantial challenges for biomedical information retrieval, consequently diminishing the performance of biomedical RAG systems. **Fairness** dictates that trustworthy RAG systems must avoid discriminatory outputs during the generation process. Due to the uneven distribution of demographic attributes in training corpora<sup>2</sup>, which may be present in both the pre-training data of LLMs and the knowledge sources utilized in RAG systems, biomedical RAG applications are susceptible to intrinsic biases. Addressing the fairness concerns also represents an important direction for future research.

### Multilinguality

Among our surveyed publications, English remains the dominant language in all Biomedical RAG systems. This predominance is primarily attributed to the dominance of English in public biomedical knowledge repositories and medical task datasets. The performance of biomedical RAG systems in other languages, particularly those with limited resources, remains inadequately explored<sup>232,233</sup>. This fact presents two significant challenges in the development of multi-lingual biomedical RAG systems.

First, training LLMs specialized in low-resource languages is inherently difficult, which diminishes the generative capabilities of Biomedical RAG in other linguistic contexts. Second, the scarcity of available knowledge sources in non-dominant languages significantly constrains the augmentation potential afforded by knowledge retrieval mechanisms. Although some researchers have attempted to address this limitation through synthetic data generated by LLMs or translated from resource-rich languages<sup>234</sup>, consistent and sustained efforts are still required to overcome this challenge in multilingual biomedical information retrieval and generation.

### Multimodality

Another challenge lies in extending retrieval capabilities beyond textual data to encompass multi-modalities highly prevalent in the biomedical domain, including images (e.g. X-ray<sup>203,235,236</sup>), time series (e.g. ECG<sup>152</sup>), sound recordings (e.g. lung sounds recorded<sup>237</sup>), and video<sup>238</sup>. Among the studies we identified, the overwhelming majority focus exclusively on text-based retrieval, with only a limited subset addressing image-based retrieval, while other modalities remain largely unexplored. Multi-modal data presents several intractable challenges for retrieval systems. For instance, in medical imaging, although we can utilize a visual encoder model as a retriever, high intra-class variance significantly complicates image retrieval processes<sup>239</sup>. Specifically, identical medical conditions can manifest with remarkable heterogeneity across visual representations; for instance, images of histologically identical tumors may exhibit substantially different morphological characteristics when visualized across different patients or at varying stages of disease progression. Attention also should be paid to developing cross-modal RAG in the Biomedical domain and improving its accuracy and trustworthiness in the future.

### Restricted Computing Resources

The application of RAG in the biomedical domain faces significant limitations due to computing resource constraints, which restrict its utility in practical scenarios. Within the RAG framework, LLMs play a crucial role in contextual understanding and human-like response generation based on retrieved knowledge<sup>1</sup>. However, these models typically require high-performance computing clusters to operate effectively. For instance, the minimum deployment unit of DeepSeek V3's decoding stage consists of 40 nodes with 320 GPUs<sup>240</sup>.

Empirical research consistently demonstrates that LLMs with larger parameter counts generally exhibit superior capacity for understanding and generation tasks<sup>241,242</sup>. Nevertheless, in real-world healthcare environments such as hospitals, maintaining high-performance computing infrastructure is often impractical and cost-prohibitive. Several promising approaches have emerged to address these constraints. One solution involves developing smaller, distilled LLMs derived from their larger counterparts<sup>243</sup>. Models such as QWQ<sup>160</sup> and Mistral Small 3.1<sup>40</sup> demonstrate competitive performance despite having significantly reduced parameter counts. Applying Portable LLMs<sup>244</sup> as generation models is also a promising method. Another viable approach incorporates RAG systems that combine local knowledge bases with external LLM API services. This hybrid configuration presents an acceptable compromise when computing resources are limited, although it introduces potential privacy concerns, as local data must be transmitted to commercial providers for response generation.

## Privacy

Privacy concerns in biomedical RAG systems exist throughout both retrieval and generation processes. During knowledge retrieval, reference documents frequently contain patients' clinical histories, which may include personally identifiable information that individuals would prefer to keep confidential<sup>245</sup>. Such information handling might not comply with stringent regulatory frameworks such as the General Data Protection Regulation<sup>246</sup> and the Health Insurance Portability and Accountability Act<sup>247</sup>. Furthermore, the generation phase presents additional challenges, particularly when utilizing commercial LLMs for biomedical RAG development. This approach necessitates the transmission of sensitive data across the internet, significantly increasing the risk of data breaches and unauthorized access. Consequently, developing biomedical RAG systems with robust privacy-preserving mechanisms remains a critical direction for future research endeavors.

## CONCLUSION

In conclusion, this comprehensive survey has provided a detailed overview of the application of RAG in the biomedical domain. It has examined the basic framework of biomedical RAG and discussed its retriever, reranker, and generation modules, as well as explored KG-enhanced biomedical RAG and multi-modal biomedical RAG. The survey has also explored clinical applications where biomedical RAG can be applied and reviewed recent studies in this field. Additionally, the paper has identified key challenges, including trustworthiness, multilingual capabilities, multimodal integration, limited computational resources, and privacy concerns. Overall, this work provides researchers and practitioners with a thorough understanding of the current state of biomedical RAG systems, ultimately contributing to the development of more accurate, accessible, and impactful healthcare AI systems that can meaningfully enhance patient care through evidence-based information access and support.

## REFERENCES

- [1] W. Fan, Y. Ding, L. Ning, S. Wang, H. Li, D. Yin, T.-S. Chua, and Q. Li, "A survey on rag meeting llms: Towards retrieval-augmented large language models," in *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '24)*. New York, NY, USA: Association for Computing Machinery, 2024, pp. 6491–6501. [Online]. Available: <https://doi.org/10.1145/3637528.3671470>
- [2] K. He, R. Mao, Q. Lin, Y. Ruan, X. Lan, M. Feng, and E. Cambria, "A survey of large language models for healthcare: from data, technology, and applications to accountability and ethics," *arXiv preprint arXiv:2310.05694*, 2023, v3, last revised 27 Jan 2025. [Online]. Available: <https://arxiv.org/abs/2310.05694>
- [3] J. Wang, Z. Yang, Z. Yao et al., "JMLR: Joint Medical LLM and Retrieval training for enhancing reasoning and professional question answering capability," *arXiv preprint arXiv:2402.17887*, 2024. [Online]. Available: <https://arxiv.org/abs/2402.17887>
- [4] M. Li, H. Kilicoglu, H. Xu, and R. Zhang, "Biomedrag: A retrieval augmented large language model for biomedicine," *Journal of Biomedical Informatics*, vol. 162, p. 104769, Feb 2025, epub 2025 Jan 13.
- [5] H. Xiao, F. Zhou, X. Liu, T. Liu, Z. Li, X. Liu, and X. Huang, "A comprehensive survey of large language models and multimodal large language models in medicine," *Information Fusion*, vol. 117, p. 102888, 2025. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1566253524006663>
- [6] G. Frisoni, A. Cocchieri, A. Presepi, G. Moro, and Z. Meng, "To generate or to retrieve? on the effectiveness of artificial contexts for medical open-domain question answering," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Bangkok, Thailand: Association for Computational Linguistics, 2024, pp. 9878–9919.
- [7] G. Xiong, Q. Jin, Z. Lu, and A. Zhang, "Benchmarking retrieval-augmented generation for medicine," in *Findings of the Association for Computational Linguistics ACL 2024*, 2024, pp. 6233–6251.
- [8] Y. Li, C. Li, Z. Wang, D. Sui, and J. Yan, "A privacy-preserving framework for medical chatbot based on llm with retrieval augmented generation." Berlin, Heidelberg: Springer-Verlag, 2024, p. 15–28. [Online]. Available: [https://doi.org/10.1007/978-981-97-9437-9\\_2](https://doi.org/10.1007/978-981-97-9437-9_2)
- [9] S. Santra, P. Kukreja, K. Saxena, S. Gandhi, and O. V. Singh, "Navigating regulatory and policy challenges for ai enabled combination devices," *Frontiers in Medical Technology*, vol. 6, p. 1473350, 2024. [Online]. Available: <https://www.frontiersin.org/articles/10.3389/fmedt.2024.1473350/full>
- [10] L. Liu, X. Yang, J. Lei, Y. Shen, J. Wang, P. Wei, Z. Chu, Z. Qin, and K. Ren, "A survey on medical large language models: Technology, application, trustworthiness, and future directions," *arXiv preprint arXiv:2406.03712*, 2024. [Online]. Available: <https://doi.org/10.48550/arXiv.2406.03712>
- [11] Y. Gao, Y. Xiong, X. Gao et al., "Retrieval-augmented generation for large language models: A survey," *arXiv preprint*, vol. arXiv:2312.10997, no. 2, 2023. [Online]. Available: <https://arxiv.org/abs/2312.10997>
- [12] W. X. Zhao, K. Zhou, J. Li et al., "A survey of large language models," *arXiv preprint*, vol. arXiv:2303.18223, no. 1(2), 2023. [Online]. Available: <https://arxiv.org/abs/2303.18223>
- [13] D. Stammbach, B. Zhang, and E. Ash, "The choice of textual knowledge base in automated claim checking," *J. Data and Information Quality*, vol. 15, no. 1, Mar. 2023. [Online]. Available: <https://doi.org/10.1145/3561389>
- [14] A. Lozano, S. L. Fleming, C. C. Chiang, and N. Shah, "Clinfo.ai: An Open-Source Retrieval-Augmented Large Language Model System for Answering Medical Questions using Scientific Literature," *Pacific Symposium on Biocomputing*, vol. 29, pp. 8–23, 2024.
- [15] J. Lammert, T. Dreyer, S. Mathes, L. Kuligin, K. J. Borm, U. A. Schatz, M. Kiechle, A. M. Lörsch, J. Jung, S. Lange, N. Pfarr, A. Durner, K. Schwamborn, C. Winter, D. Ferber, J. N. Kather, C. Mogler, A. L. Illert, and M. Tschochoei, "Expert-guided large language

models for clinical decision support in precision oncology," *JCO Precision Oncology*, vol. 8, p. e2400478, October 2024, epub 2024 Oct 30.

- [16] C. Zelin, W. K. Chung, M. Jeanne, G. Zhang, and C. Weng, "Rare disease diagnosis using knowledge guided retrieval augmentation for chatgpt," *Journal of Biomedical Informatics*, vol. 157, p. 104702, September 2024, epub 2024 Jul 29.
- [17] Y. Gao, R. Li, E. Croxford, J. Caskey, B. W. Patterson, M. Churpek, T. Miller, D. Dligach, and M. Afshar, "Leveraging medical knowledge graphs into large language models for diagnosis prediction: Design and application study," *JMIR AI*, vol. 4, p. e58670, Feb 2025.
- [18] S. Lu and E. Cosgun, "Boosting gpt models for genomics analysis: generating trusted genetic variant annotations and interpretations through rag and fine-tuning," *Bioinformatics Advances*, vol. 5, no. 1, p. vbaf019, February 5 2025.
- [19] S. Raminedi, S. Shridevi, and D. Won, "Multi-modal transformer architecture for medical image analysis and automated report generation," *Scientific Reports*, vol. 14, no. 1, p. 19281, August 20 2024.
- [20] A. E. Johnson, T. J. Pollard, L. Shen, L.-w. H. Lehman, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. A. Celi, and R. G. Mark, "Mimic-iii, a freely accessible critical care database," *Scientific Data*, vol. 3, no. 1, pp. 1–9, 2016.
- [21] E. Herrett, A. M. Gallagher, K. Bhaskaran, H. Forbes, R. Mathur, T. Van Staa, and L. Smeeth, "Data resource profile: clinical practice research datalink (cprd)," *International Journal of Epidemiology*, vol. 44, no. 3, pp. 827–836, 2015.
- [22] A. Cohan, S. Feldman, I. Beltagy, D. Downey, and D. S. Weld, "SPECTER: Document-level Representation Learning using Citation-informed Transformers," in *ACL*, 2020.
- [23] K. Bashir and H. Bukumiric, *Basic Metabolic Panel*. Treasure Island (FL): StatPearls Publishing, 2024, available from: <https://www.ncbi.nlm.nih.gov/books/NBK430685/>.
- [24] Q. Jin, C. Tan, M. Chen, M. Yan, N. Zhang, S. Huang, and X. Liu, "State-of-the-art evidence retriever for precision medicine: Algorithm development and validation," *JMIR Medical Informatics*, vol. 10, no. 12, p. e40743, Dec. 2022.
- [25] B. Zhang, N. Naderi, R. Mishra, and D. Teodoro, "Online health search via multidimensional information quality assessment based on deep language models: Algorithm development and validation," *JMIR AI*, vol. 3, p. e42630, 2024.
- [26] B. Parmanto, B. Aryoyudanta, T. W. Soekinto, I. M. A. Setiawan, Y. Wang, H. Hu, A. Saptono, and Y. K. Choi, "A reliable and accessible caregiving language model (calm) to support tools for caregivers: Development and evaluation study," *JMIR Formative Research*, vol. 8, p. e54633, Jul 31 2024.
- [27] Y. Guo, W. Qiu, C. Leroy, S. Wang, and T. Cohen, "Retrieval augmentation of large language models for lay language generation," *Journal of Biomedical Informatics*, vol. 149, p. 104580, Jan 2024, epub 2023 Dec 30.
- [28] C. S. Ong, N. T. Obey, Y. Zheng, A. Cohan, and E. B. Schneider, "Surgeryllm: a retrieval-augmented generation large language model framework for surgical decision support and workflow enhancement," *NPJ Digital Medicine*, vol. 7, no. 1, p. 364, Dec 18 2024.
- [29] C. Kakalou, C. Karamanidou, T. Dalamagas, and M. Koubarakis, "Enhancing patient empowerment and health literacy: Integrating knowledge graphs with language models for personalized health content delivery," in *Studies in Health Technology and Informatics*, vol. 316, Aug 2024, pp. 1018–1022.
- [30] R. Yang, Q. Zeng, K. You, Y. Qiao, L. Huang, C. C. Hsieh, B. Rosand, J. Goldwasser, A. Dave, T. Keenan, Y. Ke, C. Hong, N. Liu, E. Chew, D. Radev, Z. Lu, H. Xu, Q. Chen, and I. Li, "Ascle-a python natural language processing toolkit for medical text generation: Development and evaluation study," *Journal of Medical Internet Research*, vol. 26, p. e60601, Oct 2024.
- [31] L. Sun, J. Zhao, M. Han et al., "Fact-aware multimodal retrieval augmentation for accurate medical radiology report generation," *arXiv preprint*, vol. arXiv:2407.15268, 2024. [Online]. Available: <https://arxiv.org/abs/2407.15268>
- [32] D. P. Upadhyaya, A. G. Shaikh, G. B. Cakir, K. Prantzaos, P. Golnari, F. F. Ghasia, and S. S. Sahoo, "A 360° view for large language models: Early detection of amblyopia in children using multi-view eye movement recordings," *medRxiv [Preprint]*, p. 2024.05.03.24306688, May 2024.
- [33] Q. Jin, W. Kim, Q. Chen, D. C. Comeau, L. Yeganova, W. J. Wilbur, and Z. Lu, "Medcpt: Contrastive pre-trained transformers with large-scale pubmed search logs for zero-shot biomedical information retrieval," *Bioinformatics*, vol. 39, no. 11, p. btad651, 2023.
- [34] S. Jia, S. Bit, E. Searls, M. V. Lauber, L. A. Claus, P. Fan, V. H. Jasodanand, D. Veerapaneni, W. M. Wang, R. Au, and V. B. Kolachalama, "Podgpt: An audio-augmented large language model for research and education," *medRxiv [Preprint]*, November 2024, PMID: 39040167; PMCID: PMC11261953. [Online]. Available: <https://doi.org/10.1101/2024.07.11.24310304>
- [35] OpenAI. (2025) Text generation - openai api. OpenAI. [Online]. Available: <https://platform.openai.com/docs/guides/text-generation>
- [36] Anthropic. (2025) Getting started - anthropic api documentation. Anthropic. [Online]. Available: <https://docs.anthropic.com/en/api/getting-started>
- [37] Google AI. (2025) Gemini api: Text generation guide. Google. [Online]. Available: <https://ai.google.dev/gemini-api/docs/text-generation>
- [38] M. AI, "Llama 3: 70b high-performance open language models," *Meta AI Research*, 2024. [Online]. Available: <https://ai.meta.com/llama/papers-and-research/>
- [39] M. N. Team. (2023) Introducing mpt-7b: A new standard for open-source, commercially usable llms. Accessed: 2023-05-05. [Online]. Available: [www.mosaicml.com/blog/mpt-7b](https://www.mosaicml.com/blog/mpt-7b)
- [40] Mistral AI, "Mistral-small-3.1-24b-instruct-2503," <https://huggingface.co/mistralai/Mistral-Small-3.1-24B-Instruct-2503>, 2024, accessed: 2025-03-16.
- [41] DeepSeek-AI Team, "Deepseek-ai/deepseek-r1," <https://huggingface.co/deepseek-ai/DeepSeek-R1>, 2024, accessed: 2025-03-16.
- [42] Author(s), "A framework enabling llms into regulatory environment for transparency and trustworthiness and its application to drug labeling document," *Journal/Conference Name*, 2024.
- [43] J. Wang, J. X. Huang, X. Tu, J. Wang, A. J. Huang, M. T. R. Laskar, and A. Bhuiyan, "Utilizing bert for information retrieval: Survey, applications, resources, and challenges," *ACM Comput. Surv.*, vol. 56, no. 7, Apr. 2024. [Online]. Available: <https://doi.org/10.1145/3648471>
- [44] L. Tamine and L. Coeuriot, "Semantic information retrieval on medical texts: Research challenges, survey, and open issues," *ACM Comput. Surv.*, vol. 54, no. 7, Sep. 2021. [Online]. Available: <https://doi.org/10.1145/3462476>
- [45] S. Sivarajkumar, H. A. Mohammad, D. Oniani, K. Roberts, W. Hersh, H. Liu, D. He, S. Visweswaran, and Y. Wang, "Clinical information retrieval: A literature review," *Journal of Healthcare Informatics Research*, vol. 8, no. 2, pp. 313–352, 2024.
- [46] S. Robertson and H. Zaragoza, "The probabilistic relevance framework: Bm25 and beyond," *Found. Trends Inf. Retr.*, vol. 3, no. 4, p. 333–389, Apr. 2009. [Online]. Available: <https://doi.org/10.1561/15000000019>



- [47] K. S. Jones, "A statistical interpretation of term specificity and its application in retrieval," *Journal of Documentation*, vol. 28, no. 1, pp. 11–21, 1972.
- [48] S. Das, Y. Ge, Y. Guo, S. Rajwal, J. Hairston, J. Powell, D. Walker, S. Peddireddy, S. Lakamana, S. Bozkurt, M. Reyna, R. Sameni, Y. Xiao, S. Kim, R. Chandler, N. Hernandez, D. Mowery, R. Wightman, J. Love, A. Spadaro, J. Perrone, and A. Sarker, "Two-layer retrieval-augmented generation framework for low-resource medical question answering using reddit data: Proof-of-concept study," *Journal of Medical Internet Research*, vol. 27, p. e66220, Jan. 2025.
- [49] E. J. Gong, C. S. Bang, J. J. Lee, J. Park, E. Kim, S. Kim, M. Kimm, and S. H. Choi, "The potential clinical utility of the customized large language model in gastroenterology: A pilot study," *Bioengineering (Basel)*, vol. 12, no. 1, p. 1, Dec. 2024.
- [50] J. Lee, H. Cha, Y. Hwangbo, and W. Cheon, "Enhancing large language model reliability: Minimizing hallucinations with dual retrieval-augmented generation based on the latest diabetes guidelines," *Journal of Personalized Medicine*, vol. 14, no. 12, p. 1131, Nov. 2024.
- [51] I. Alonso, M. Oronoz, and R. Agerri, "Medexpqa: Multilingual benchmarking of large language models for medical question answering," *Artificial Intelligence in Medicine*, vol. 155, p. 102938, 2024. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0933365724001805>
- [52] T. Abdullahi, L. Mercurio, R. Singh, and C. Eickhoff, "Retrieval-based diagnostic decision support: Mixed methods study," *JMIR Medical Informatics*, vol. 12, p. e50209, Jun. 2024.
- [53] S. Kim and J. Yoon, "Vaiv bio-discovery service using transformer model and retrieval augmented generation," *BMC Bioinformatics*, vol. 25, no. 1, p. 273, Aug. 2024.
- [54] E. Kamalloo, N. Thakur, C. Lassance, X. Ma, J.-H. Yang, and J. Lin, "Resources for brewing beer: Reproducible reference models and statistical analyses," ser. SIGIR '24. New York, NY, USA: Association for Computing Machinery, 2024, p. 1431–1440. [Online]. Available: <https://doi.org/10.1145/3626772.3657862>
- [55] R. Xu, W. Shi, Y. Yu, Y. Zhuang, Y. Zhu, M. D. Wang, J. C. Ho, C. Zhang, and C. Yang, "BMRetriever: Tuning large language models as better biomedical text retrievers," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. Miami, Florida, USA: Association for Computational Linguistics, 2024, pp. 22 234–22 254.
- [56] M. Hu, L. Zong, H. Wang, J. Zhou, J. Li, Y. Gao, K.-F. Wong, Y. Li, and I. King, "Serts: Self-rewarding tree search for biomedical retrieval-augmented generation," in *Findings of the Association for Computational Linguistics: EMNLP 2024*. Miami, Florida, USA: Association for Computational Linguistics, 2024, pp. 1321–1335.
- [57] G. Izacard, M. Caron, L. Hosseini, S. Riedel, P. Bojanowski, A. Joulin, and E. Grave, "Unsupervised dense information retrieval with contrastive learning," 2021. [Online]. Available: <https://arxiv.org/abs/2112.09118>
- [58] Elastic, "Elasticsearch Documentation: Getting Started Guide," <https://www.elastic.co/guide/en/elasticsearch/reference/current/getting-started.html>, 2024, accessed: 2025-03-13.
- [59] T. W. Team, "Whoosh documentation: Release notes," <https://whoosh.readthedocs.io/en/latest/releases/index.html>, 2024, accessed: 2025-03-13.
- [60] N. C. for Biotechnology Information (NCBI), "Entrez programming utilities (e-utilities) api," <https://www.ncbi.nlm.nih.gov/home/develop/api/>, 2024, accessed: 2025-03-13.
- [61] G. Xiong, Q. Jin, Z. Lu, and A. Zhang, "Benchmarking retrieval-augmented generation for medicine," in *Findings of the Association for Computational Linguistics ACL 2024*, L.-W. Ku, A. Martins, and V. Srikumar, Eds. Bangkok, Thailand and virtual meeting: Association for Computational Linguistics, Aug. 2024, pp. 6233–6251. [Online]. Available: <https://aclanthology.org/2024.findings-acl.372>
- [62] D. Peasley, R. Kuplicki, S. Sen, and M. Paulus, "Leveraging large language models and agent-based systems for scientific data analysis: Validation study," *JMIR Mental Health*, vol. 12, p. e68135, Feb. 2025.
- [63] W. X. Zhao, J. Liu, R. Ren, and J.-R. Wen, "Dense text retrieval based on pretrained language models: A survey," *ACM Trans. Inf. Syst.*, vol. 42, no. 4, Feb. 2024. [Online]. Available: <https://doi.org/10.1145/3637870>
- [64] Y. Santander-Cruz, S. Salazar-Colores, W. J. Paredes-Garcia, I. Cruz-Aceves, L. Chacón-Hernández, C. E. Galván-Tejada, L. A. Zanella-Calzada, N. K. Gamboa-Rosales, J. M. Celaya-Padilla, J. M. Gálvez et al., "Semantic feature extraction using sbert for dementia detection," *Brain Sciences*, vol. 12, no. 2, p. 270, 2022.
- [65] N. Thakur, N. Reimers, J. Daxenberger, and I. Gurevych, "Augmented sbert: Data augmentation method for improving bi-encoders for pairwise sentence scoring tasks," *arXiv preprint arXiv:2010.08240*, 2020. [Online]. Available: <https://arxiv.org/abs/2010.08240>
- [66] S. Xiao, Z. Liu, P. Zhang, and N. Muennighoff, "C-pack: Packaged resources to advance general chinese embedding," 2023.
- [67] V. Karpukhin, B. Oguz, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, and W. tau Yih, "Dense passage retrieval for open-domain question answering," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 6769–6781.
- [68] A. V. Solatorio, "Gistembed: Guided in-sample selection of training negatives for text embedding fine-tuning," *arXiv preprint arXiv:2402.16829*, 2024. [Online]. Available: <https://arxiv.org/abs/2402.16829>
- [69] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, "Enriching word vectors with subword information," *Transactions of the Association for Computational Linguistics*, vol. 5, pp. 135–146, 2017.
- [70] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized bert pretraining approach," *arXiv preprint arXiv:1907.11692*, 2019.
- [71] T. Gao, X. Yao, and D. Chen, "SimCSE: Simple contrastive learning of sentence embeddings," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, M.-F. Moens, X. Huang, L. Specia, and S. W.-t. Yih, Eds. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 6894–6910. [Online]. Available: <https://aclanthology.org/2021.emnlp-main.552/>
- [72] B. T. Garcia, L. Westerfield, P. Yelemali, N. Gogate, E. Andres Rivera-Munoz, H. Du, M. Dawood, A. Jolly, J. R. Lupski, and J. E. Posey, "Improving automated deep phenotyping through large language models using retrieval augmented generation," *medRxiv [Preprint]*, p. 2024.12.01.24318253, Dec 2 2024.
- [73] Q. Team, "fastembed: A blazing fast embedding library," <https://github.com/qdrant/fastembed>, 2023, accessed: 2025-04-13.
- [74] M. Li, H. Zhou, H. Yang, and R. Zhang, "Rt: a retrieving and chain-of-thought framework for few-shot medical named entity recognition," *Journal of the American Medical Informatics Association*, vol. 31, no. 9, pp. 1929–1938, Sep 1 2024.
- [75] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, 2019, pp. 4171–4186.



- [76] R. Xu, Y. Hong, F. Zhang, and H. Xu, "Evaluation of the integration of retrieval-augmented generation in large language model for breast cancer nursing care responses," *Scientific Reports*, vol. 14, no. 1, p. 30794, Dec. 2024.
- [77] S. Darnell, R. Overall, A. Guarracino, V. Colonna, F. Villani, E. Garrison, A. Isaac, P. Muli, F. Muriithi, A. Kabui, M. Kilyungi, F. Lisso, A. Kibet, B. Muhia, H. Nijveen, S. Yousefi, D. Ashbrook, P. Huang, G. Suh, M. Umar, C. Batten, H. Chen, S. Sen, R. Williams, and P. Prins, "Creating a biomedical knowledge base by addressing gpt inaccurate responses and benchmarking context," *bioRxiv*, 2024, preprint. [Online]. Available: <https://doi.org/10.1101/2024.10.16.618663>
- [78] Y. A. Malkov and D. A. Yashunin, "Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs," *CoRR*, vol. abs/1603.09320, 2016. [Online]. Available: <https://arxiv.org/abs/1603.09320>
- [79] E. Klang, I. Tessler, D. U. Apakama, E. Abbott, B. S. Glicksberg, M. Arnold, A. Moses, A. Sakhuja, A. W. Charney, D. L. Reich, J. McGreevy, N. Gavin, B. Carr, R. Freeman, and G. N. Nadkarni, "Assessing retrieval-augmented large language model performance in emergency department icd-10-cm coding compared to human coders," *medRxiv [Preprint]*, Oct 17 2024.
- [80] Y. Zheng, Y. Yan, S. Chen, Y. Cai, K. Ren, Y. Liu, J. Zhuang, and M. Zhao, "Integrating retrieval-augmented generation for enhanced personalized physician recommendations in web-based medical services: model development study," *Frontiers in Public Health*, vol. 13, p. 1501408, Jan 29 2025.
- [81] T. Fukushima, M. Manabe, S. Yada, S. Wakamiya, A. Yoshida, Y. Urakawa, A. Maeda, S. Kan, M. Takahashi, and E. Aramaki, "Evaluating and enhancing japanese large language models for genetic counseling support: Comparative study of domain adaptation and the development of an expert-evaluated dataset," *JMIR Medical Informatics*, vol. 13, p. e65047, Jan 2025.
- [82] A. Fukuchi, Y. Hoshino, and Y. Watanabe, "Glucose (general luke-based contrastive sentence embedding)," <https://huggingface.co/pkshatech/GLUCoSE-base-ja>, 2023, hugging Face.
- [83] S. Toro, A. V. Anagnostopoulos, S. M. Bello, K. Blumberg, R. Cameron, L. Carmody, A. D. Diehl, D. M. Dooley, W. D. Duncan, P. Fey, P. Gaudet, N. L. Harris, M. P. Joachimiak, L. Kiani, T. Lubiana, M. C. Munoz-Torres, S. O'Neil, D. Osumi-Sutherland, A. Puig-Barbe, J. T. Reese, L. Reiser, S. M. Robb, T. Ruemping, J. Seager, E. Sid, R. Stefancsik, M. Weber, V. Wood, M. A. Haendel, and C. J. Mungall, "Dynamic retrieval augmented generation of ontologies using artificial intelligence (dragon-ai)," *Journal of Biomedical Semantics*, vol. 15, no. 1, p. 19, Oct 17 2024.
- [84] OpenAI. (2023) Vector embeddings. [Online]. Available: <https://platform.openai.com/docs/guides/embeddings>
- [85] D. Wang, J. Liang, J. Ye, J. Li, J. Li, Q. Zhang, Q. Hu, C. Pan, D. Wang, Z. Liu, W. Shi, D. Shi, F. Li, B. Qu, and Y. Zheng, "Enhancement of the performance of large language models in diabetes education through retrieval-augmented generation: Comparative study," *Journal of Medical Internet Research*, vol. 26, p. e58041, Nov 8 2024.
- [86] J. Lee, H. Cha, Y. Hwangbo, and W. Cheon, "Enhancing large language model reliability: Minimizing hallucinations with dual retrieval-augmented generation based on the latest diabetes guidelines," *Journal of Personalized Medicine*, vol. 14, no. 12, p. 1131, Nov 30 2024.
- [87] T. C. Li, H. Zhou, V. Verma, X. Tang, Y. Shao, E. Van Buren, Z. Weng, M. Gerstein, B. Neale, S. R. Sunyaev, and X. Lin, "Favor-gpt: a generative natural language interface to whole genome variant functional annotations," *Bioinformatics Advances*, vol. 4, no. 1, p. vbae143, Sep 28 2024.
- [88] Z. Fu, S. Fu, Y. Huang, W. He, Z. Zhong, Y. Guo, and Y. Lin, "Application of large language model combined with retrieval enhanced generation technology in digestive endoscopic nursing," *Frontiers in Medicine (Lausanne)*, vol. 11, p. 1500258, Nov 6 2024.
- [89] C. Long, D. Subburam, K. Lowe, A. Dos Santos, J. Zhang, S. Hwang, N. Saduka, Y. Horev, T. Su, D. W. J. Côté, and E. D. Wright, "Chatent: Augmented large language model for expert knowledge retrieval in otolaryngology-head and neck surgery," *Otolaryngology-Head and Neck Surgery*, vol. 171, no. 4, pp. 1042–1051, Oct 2024.
- [90] D. Steybe, P. Poxleitner, S. Aljohani, B. B. Herlofson, O. Nicolatou-Galitis, V. Patel, S. Fedele, T. G. Kwon, V. Fusco, S. E. C. Pichardo, K. T. Obermeier, S. Otto, A. Rau, and M. F. Russe, "Evaluation of a context-aware chatbot using retrieval-augmented generation for answering clinical questions on medication-related osteonecrosis of the jaw," *J Craniomaxillofac Surg*, pp. S1010–5182(24)00341–X, Jan 10 2025.
- [91] O. Unlu, J. Shin, C. J. Mailly, M. F. Oates, M. R. Tucci, M. Varugheese, K. Waghlikar, F. Wang, B. M. Scirica, A. J. Blood, and S. J. Aronson, "Retrieval augmented generation enabled generative pre-trained transformer 4 (gpt-4) performance for clinical trial screening," *medRxiv [Preprint]*, p. 2024.02.08.24302376, Feb 8 2024.
- [92] I. Azimi, M. Qi, L. Wang, A. M. Rahmani, and Y. Li, "Evaluation of llms accuracy and consistency in the registered dietitian exam through prompt engineering and knowledge retrieval," *Sci Rep*, vol. 15, no. 1, p. 1506, Jan 9 2025.
- [93] Amazon Web Services, "Amazon titan embeddings models - user guide," <https://docs.aws.amazon.com/bedrock/latest/userguide/titan-embedding-models.html>, 2024, accessed: 2025-03-16.
- [94] Y. Selcuk, E. Kim, and I. Ahn, "Infected-chat, an arabic large language model for infectious diseases: Comparative analysis," *JMIR Med Inform*, vol. 13, p. e63881, Feb 10 2025.
- [95] D. Kainer, "The effectiveness of large language models with rag for auto-annotating trait and phenotype descriptions," *Biol Methods Protoc*, vol. 10, no. 1, p. bpaf016, Feb 26 2025.
- [96] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, "BioBERT: a pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, 09 2019. [Online]. Available: <https://doi.org/10.1093/bioinformatics/btz682>
- [97] Q. Jin, A. Shin, and Z. Lu, "Lader: Log-augmented dense retrieval for biomedical literature search," in *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, ser. SIGIR '23. New York, NY, USA: Association for Computing Machinery, 2023, p. 2092–2097. [Online]. Available: <https://doi.org/10.1145/3539618.3592005>
- [98] Y. Gu, R. Tinn, H. Cheng, M. Lucas, N. Usuyama, X. Liu, T. Naumann, J. Gao, and H. Poon, "Domain-specific language model pretraining for biomedical natural language processing," 2020.
- [99] I. Lopez, A. Swaminathan, K. Vedula, S. Narayanan, F. Nateghi Haredasht, S. P. Ma, A. S. Liang, S. Tate, M. Maddali, R. J. Gallo, N. H. Shah, and J. H. Chen, "Clinical entity augmented retrieval for clinical information extraction," *NPJ Digital Medicine*, vol. 8, no. 1, p. 45, Jan 19 2025.
- [100] K. Soman, P. Rose, J. Morris, R. Akbas, B. Smith, B. Peetoom, C. Villouta-Reyes, G. Cerono, Y. Shi, A. Rizk-Jackson, S. Israni, C. Nelson, S. Huang, and S. Baranzini, "Biomedical knowledge graph-optimized prompt generation for large language models," *Bioinformatics*, vol. 40, no. 9, p. btae560, Sep 2024.
- [101] W. Aftab, Z. Apostolou, K. Bouazoune, and T. Straub, "Optimizing biomedical information retrieval with a keyword frequency-driven prompt enhancement strategy," *BMC Bioinformatics*, vol. 25, no. 1, p. 281, Aug 2024.
- [102] R. Luo, L. Sun, Y. Xia, T. Qin, S. Zhang, H. Poon, and T.-Y. Liu, "BioGPT: Generative pre-trained transformer for biomedical text generation and mining," *Briefings in Bioinformatics*, vol. 23, no. 6, p. bbac409, Nov. 2022.

- [103] Y. Anand, Z. Nussbaum, B. Duderstadt, and B. Schmidt, "Gpt4all: Training an assistant-style chatbot with large scale data distillation from gpt-3.5-turbo," <https://github.com/nomic-ai/gpt4all>, 2023, accessed: 2025-03-16.
- [104] S. Myers, T. A. Miller, Y. Gao, M. M. Churpek, A. Mayampurath, D. Dligach, and M. Afshar, "Lessons learned on information retrieval in electronic health records: a comparison of embedding models and pooling strategies," *Journal of the American Medical Informatics Association*, vol. 32, no. 2, pp. 357–364, Feb 2025.
- [105] P. Behnamghader, V. Adlakha, M. Mosbach, D. Bahdanau, N. Chapados, and S. Reddy, "LLM2Vec: Large language models are secretly powerful text encoders," in *First Conference on Language Modeling*, 2024. [Online]. Available: <https://openreview.net/forum?id=IW1PR7vEBf>
- [106] M. Jeong, J. Sohn, M. Sung, and J. Kang, "Improving medical reasoning through retrieval and self-reflection with retrieval-augmented large language models," *Bioinformatics*, vol. 40, no. Suppl 1, pp. i119–i129, Jun. 2024.
- [107] Z. Zhan, S. Zhou, M. Li, and R. Zhang, "Ramie: retrieval-augmented multi-task information extraction with large language models on dietary supplements," *Journal of the American Medical Informatics Association*, vol. 32, no. 3, pp. 545–554, Mar. 2025.
- [108] Anthropic, "Embeddings api," <https://docs.anthropic.com/en/docs/build-with-claude/embeddings>, 2024, accessed: 2025-02-19.
- [109] Google Cloud, "Get text embeddings with vertex ai," <https://cloud.google.com/vertex-ai/generative-ai/docs/embeddings/get-text-embeddings>, 2024, accessed: 2025-03-16.
- [110] D. Lyu, X. Wang, Y. Chen, and F. Wang, "Language model and its interpretability in biomedicine: A scoping review," *iScience*, vol. 27, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:267973118>
- [111] C. Wu, W. Lin, X. Zhang, Y. Zhang, Y. Wang, and W. Xie, "Pmc-llama: Towards building open-source language models for medicine," *arXiv preprint arXiv:2304.14454*, 2023.
- [112] Z. Liu, C. Li, S. Xiao, Y. Shao, and D. Lian, "Llama2Vec: Unsupervised adaptation of large language models for dense retrieval," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, L.-W. Ku, A. Martins, and V. Srikumar, Eds. Bangkok, Thailand: Association for Computational Linguistics, Aug. 2024, pp. 3490–3500. [Online]. Available: <https://aclanthology.org/2024.acl-long.191/>
- [113] K. Luo, Z. Liu, S. Xiao, T. Zhou, Y. Chen, J. Zhao, and K. Liu, "Landmark embedding: A chunking-free embedding method for retrieval augmented long-context large language models," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, L.-W. Ku, A. Martins, and V. Srikumar, Eds. Bangkok, Thailand: Association for Computational Linguistics, Aug. 2024, pp. 3268–3281. [Online]. Available: <https://aclanthology.org/2024.acl-long.180/>
- [114] Y. Han, C. Liu, and P. Wang, "A comprehensive survey on vector database: Storage and retrieval technique, challenge," 2023, arXiv:2310.10844 [cs.DB]. [Online]. Available: <https://arxiv.org/abs/2310.10844>
- [115] Chroma. (2023) Getting started with chroma. [Online]. Available: <https://docs.trychroma.com/docs/overview/getting-started>
- [116] Facebook Research. (2023) Faiss: A library for efficient similarity search. [Online]. Available: <https://github.com/facebookresearch/faiss>
- [117] Pinecone, "Pinecone overview," <https://docs.pinecone.io/guides/getting-started/overview>, 2024, accessed: 2024-03-22.
- [118] Weaviate, "Weaviate developers documentation," <https://weaviate.io/developers/weaviate>, 2024, accessed: 2025-03-16.
- [119] B. Peng, Y. Zhu, Y. Liu, X. Bo, H. Shi, C. Hong, Y. Zhang, and S. Tang, "Graph retrieval-augmented generation: A survey," 2024, arXiv:2408.08092 [cs.AI]. [Online]. Available: <https://arxiv.org/abs/2408.08092>
- [120] A. Di Maria, L. Bellomo, F. Billeci, A. Cardillo, S. Alaimo, P. Ferragina, A. Ferro, and A. Pulvirenti, "Netme 2.0: a web-based platform for extracting and modeling knowledge from biomedical literature as a labeled graph," *Bioinformatics*, vol. 40, no. 5, p. btac194, May 2024.
- [121] N. Matsumoto, J. Moran, H. Choi, M. E. Hernandez, M. Venkatesan, P. Wang, and J. H. Moore, "Kragen: a knowledge graph-enhanced rag framework for biomedical problem solving using large language models," *Bioinformatics*, vol. 40, no. 6, p. btac353, June 2024. [Online]. Available: <https://doi.org/10.1093/bioinformatics/btac353>
- [122] D. Li, S. Yang, Z. Tan, J. Y. Baik, S. Yun, J. Lee, A. Chacko, B. Hou, D. Duong-Tran, Y. Ding, H. Liu, L. Shen, and T. Chen, "DALK: Dynamic Co-Augmentation of LLMs and KG to answer Alzheimer's Disease Questions with Scientific Literature," in *Findings of the Association for Computational Linguistics: EMNLP 2024*. Miami, Florida, USA: Association for Computational Linguistics, 2024, pp. 2187–2205.
- [123] S. Gilbert, J. N. Kather, and A. Hogan, "Augmented non-hallucinating large language models as medical information curators," *NPJ Digital Medicine*, vol. 7, no. 1, p. 100, Apr 2024.
- [124] Z. Xiao, H. B. Pakrasi, Y. Chen, and Y. J. Tang, "Network for knowledge organization (neko): An ai knowledge mining workflow for synthetic biology research," *Metabolic Engineering*, vol. 87, pp. 60–67, 2025. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1096717624001484>
- [125] N. Matsumoto, H. Choi, J. Moran, M. E. Hernandez, M. Venkatesan, X. Li, J. H. Chang, P. Wang, and J. H. Moore, "Escargot: an ai agent leveraging large language models, dynamic graph of thoughts, and biomedical knowledge graphs for enhanced reasoning," *Bioinformatics*, vol. 41, no. 2, p. btac031, Feb 2025.
- [126] M. M. Abootorabi, A. Zobeiri, M. Dehghani, M. Mohammadkhani, B. Mohammadi, O. Ghahroodi, M. S. Baghshah, and E. Asgari, "Ask in any modality: A comprehensive survey on multimodal retrieval-augmented generation," 2025. [Online]. Available: <https://arxiv.org/abs/2502.08826>
- [127] P. Thetbanthad, B. Sathanarugsawait, and P. Praneetpolgrang, "Application of generative artificial intelligence models for accurate prescription label identification and information retrieval for the elderly in northern east of thailand," *J Imaging*, vol. 11, no. 1, p. 11, Jan 2025.
- [128] R. Tozuka, H. Johno, A. Amakawa, J. Sato, M. Muto, S. Seki, A. Komaba, and H. Onishi, "Application of notebooklm, a large language model with retrieval-augmented generation, for lung cancer staging," *Japanese Journal of Radiology*, Nov 2024, epub ahead of print.
- [129] M. Omar, V. Ullanat, M. Loda, L. Marchionni, and R. Umerton, "Chatgpt for digital pathology research," *Lancet Digital Health*, vol. 6, no. 8, pp. e595–e600, Aug 2024, epub 2024 Jul 9.
- [130] M. Ranjit, G. Ganapathy, R. Manuel, and T. Ganu, "Retrieval augmented chest x-ray report generation using openai gpt models," in *Proceedings of the 8th Machine Learning for Healthcare Conference*, ser. Proceedings of Machine Learning Research, K. Deshpande, M. Fiterau, S. Joshi, Z. Lipton, R. Ranganath, I. Urteaga, and S. Yeung, Eds., vol. 219. PMLR, 11–12 Aug 2023, pp. 650–666. [Online]. Available: <https://proceedings.mlr.press/v219/ranjit23a.html>
- [131] P. Xia, K. Zhu, H. Li, T. Wang, W. Shi, S. Wang, L. Zhang, J. Zou, and H. Yao, "Mmed-rag: Versatile multimodal rag system for medical vision language models," *arXiv preprint arXiv:2410.13085*, 2024.

- [132] P. Xia, K. Zhu, H. Li, H. Zhu, Y. Li, G. Li, L. Zhang, and H. Yao, "Rule: Reliable multimodal rag for factuality in medical vision language models," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2024, pp. 1081–1093.
- [133] Y. Yang et al., "Spatio-temporal and retrieval-augmented modelling for chest x-ray report generation," *IEEE Transactions on Medical Imaging*, 2025.
- [134] D. Hu, Z. Jiang, J. Shi, F. Xie, K. Wu, K. Tang, M. Cao, J. Huai, and Y. Zheng, "Pathology report generation from whole slide images with knowledge retrieval and multi-level regional feature selection," *Comput Methods Programs Biomed*, vol. 263, p. 108677, May 2025, epub 2025 Feb 27.
- [135] X. Chen, C.-J. Hsieh, and B. Gong, "When vision transformers outperform resnets without pretraining or strong data augmentations," *arXiv preprint arXiv:2106.01548*, 2021.
- [136] P. Adejumo, P. Thangaraj, S. V. Shankar, L. S. Dhingra, A. Aminorroaya, and R. Khera, "Retrieval-augmented generation for extracting CHA2DS2-VASc risk factors from unstructured clinical notes in patients with atrial fibrillation," *medRxiv* [Preprint], September 2024, pMID: 39371151; PMCID: PMC11451809. [Online]. Available: <https://doi.org/10.1101/2024.09.19.24313992>
- [137] J. Chen, S. Xiao, P. Zhang, K. Luo, D. Lian, and Z. Liu, "M3-embedding: Multi-linguality, multi-functionality, multi-granularity text embeddings through self-knowledge distillation," in *Findings of the Association for Computational Linguistics: ACL 2024*. Bangkok, Thailand: Association for Computational Linguistics, 2024, pp. 2318–2335.
- [138] D. B. Craig and S. Drăghici, "LmRaC: a functionally extensible tool for LLM interrogation of user experimental results," *Bioinformatics*, vol. 40, no. 12, p. btac679, November 2024, pMID: 39546375; PMCID: PMC11645126.
- [139] J. Lee, L. H. Pham, and Ö. Uzuner, "Enhancing consumer health question reformulation: Chain-of-thought prompting integrating focus, type, and user knowledge level," in *Proceedings of the First Workshop on Patient-Oriented Language Processing (CL4Health) @ LREC-COLING 2024*, Torino, Italia, pp. 220–228.
- [140] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, D. Bikel, L. Blecher, C. C. Ferrer, M. Chen, G. Cucurull, D. Esiobu, J. Fernandes, J. Fu, W. Fu, B. Fuller, C. Gao, V. Goswami, N. Goyal, A. Hartshorn, S. Hosseini, R. Hou, H. Inan, M. Kardas, V. Kerkez, M. Khabsa, I. Kloumann, A. Korenev, P. S. Koura, M.-A. Lachaux, T. Lavril, J. Lee, D. Liskovich, Y. Lu, Y. Mao, X. Martinet, T. Mihaylov, P. Mishra, I. Molybog, Y. Nie, A. Poulton, J. Reizenstein, R. Rungta, K. Saladi, A. Schelten, R. Silva, E. M. Smith, R. Subramanian, X. E. Tan, B. Tang, R. Taylor, A. Williams, J. X. Kuan, P. Xu, Z. Yan, I. Zarov, Y. Zhang, A. Fan, M. Kambadur, S. Narang, A. Rodriguez, R. Stojnic, S. Edunov, and T. Scialom, "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint*, 2023. [Online]. Available: <https://arxiv.org/abs/2307.09288>
- [141] Mistral AI, "Mistral-7b-instruct-v0.3," <https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>, 2024, accessed: 2025-03-16.
- [142] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *Journal of Machine Learning Research*, vol. 21, no. 140, pp. 1–67, 2020. [Online]. Available: <http://jmlr.org/papers/v21/20-074.html>
- [143] Google, "Gemma 3 27b instruct - hugging face," <https://huggingface.co/google/gemma-3-27b-it>, 2024, accessed: 2025-04-13.
- [144] K. Hewitt, I. Wiest, Z. Carrero, L. Bejan, T. Millner, S. Brandner, and J. Kather, "Large language models as a diagnostic support tool in neuropathology," *Journal of Pathology: Clinical Research*, vol. 10, no. 6, p. e70009, November 2024.
- [145] S. Ozaki, Y. Kato, S. Feng, M. Tomita, K. Hayashi, R. Obara, M. Oyamada, K. Hayashi, H. Kamigaito, and T. Watanabe, "Understanding the impact of confidence in retrieval augmented generation: A case study in the medical domain," <https://arxiv.org/abs/2412.17895>, 2024, arXiv:2412.17895 [cs.CL].
- [146] M. Alkaeed, S. Abioye, A. Qayyum, Y. M. Mekki, I. Berrou, M. Abdallah, A. Al-Fuqaha, M. Bilal, and J. Qadir, "Open foundation models in healthcare: Challenges, paradoxes, and opportunities with genai driven personalized prescription," <https://arxiv.org/abs/2502.04356>, 2025, arXiv:2502.04356 [cs.CL].
- [147] T. GLM, A. Zeng, B. Xu, B. Wang, C. Zhang, D. Yin, D. Rojas, G. Feng, H. Zhao, H. Lai, H. Yu, H. Wang, J. Sun, J. Zhang, J. Cheng, J. Gui, J. Tang, J. Zhang, J. Li, L. Zhao, L. Wu, L. Zhong, M. Liu, M. Huang, P. Zhang, Q. Zheng, R. Lu, S. Duan, S. Zhang, S. Cao, S. Yang, W. L. Tam, W. Zhao, X. Liu, X. Xia, X. Zhang, X. Gu, X. Lv, X. Liu, X. Liu, X. Yang, X. Song, X. Zhang, Y. An, Y. Xu, Y. Niu, Y. Yang, Y. Li, Y. Bai, Y. Dong, Z. Qi, Z. Wang, Z. Yang, Z. Du, Z. Hou, and Z. Wang, "Chatglm: A family of large language models from glm-130b to glm-4 all tools," 2024.
- [148] L. Zhao, Y. Wang, X. Wang, R. Lin, Z. Gang, and B. Xu, "COPD-chatglm: A chronic obstructive pulmonary disease diagnostic model," in *2024 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 2024, pp. 2965–2970.
- [149] M. Luo, J. Pang, S. Bi, Y. Lai, J. Zhao, Y. Shang, T. Cui, Y. Yang, Z. Lin, L. Zhao, X. Wu, D. Lin, J. Chen, and H. Lin, "Development and evaluation of a retrieval-augmented large language model framework for ophthalmology," *JAMA Ophthalmology*, vol. 142, no. 9, pp. 798–805, 2024.
- [150] K. Kreimeyer, J. V. Canzoniero, M. Fatteh, V. Anagnostou, and T. Botsis, "Using retrieval-augmented generation to capture molecularly-driven treatment relationships for precision oncology," *Studies in Health Technology and Informatics*, vol. 316, pp. 983–987, 2024.
- [151] Z. Zhan, J. Wang, S. Zhou, J. Deng, and R. Zhang, "Mmrag: Multi-mode retrieval-augmented generation with large language models for biomedical in-context learning," 2025. [Online]. Available: <https://arxiv.org/abs/2502.15954>
- [152] H. Yu, P. Guo, and A. Sano, "Zero-shot ecg diagnosis with large language models and retrieval-augmented generation," in *Proceedings of the 3rd Machine Learning for Health Symposium*, ser. Proceedings of Machine Learning Research, S. Hegselmann, A. Parziale, D. Shanmugam, S. Tang, M. N. Asiedu, S. Chang, T. Hartvigsen, and H. Singh, Eds., vol. 225. PMLR, 10 Dec 2023, pp. 650–663. [Online]. Available: <https://proceedings.mlr.press/v225/yu23b.html>
- [153] X. Chen, Z. Zhao, W. Zhang, P. Xu, Y. Wu, M. Xu, L. Gao, Y. Li, X. Shang, D. Shi, and M. He, "Eyegpt for patient inquiries and medical education: Development and validation of an ophthalmology large language model," *Journal of Medical Internet Research*, vol. 26, p. e60063, 2024.
- [154] X. Du, J. Novoa-Laurentiev, J. M. Plasaek, Y. W. Chuang, L. Wang, G. Marshall, S. K. Mueller, F. Chang, S. Datta, H. Paek, B. Lin, Q. Wei, X. Wang, J. Wang, H. Ding, F. J. Manion, J. Du, D. W. Bates, and L. Zhou, "Enhancing early detection of cognitive decline in the elderly: A comparative study utilizing large language models in clinical notes," *medRxiv*, p. 2024.04.03.24305298, 2024, update in: EBioMedicine. 2024 Nov;109:105401. doi:10.1016/j.ebiom.2024.105401. PMID: 38633810; PMCID: PMC11023645. [Online]. Available: <https://doi.org/10.1101/2024.04.03.24305298>
- [155] A. Fatharani and A. Alsayegh, "Pharmacogenomics meets generative ai: transforming clinical trial design with large language models," *Journal of Pharmacology and Pharmacotherapeutics*, p. 0976500X251321885, 2025.

- [156] J. Woo, A. Yang, R. Olsen, S. Hasan, D. Nawabi, B. Nwachukwu, R. r. Williams, and P. Ramkumar, "Custom large language models improve accuracy: Comparing retrieval augmented generation and artificial intelligence agents to noncustom models for evidence-based medicine," *Arthroscopy*, vol. 41, no. 3, pp. 565–573.e6, Mar 2025, epub 2024 Nov 7.
- [157] G. Xiong, Q. Jin, X. Wang, M. Zhang, Z. Lu, and A. Zhang, "Improving retrieval-augmented generation in medicine with iterative follow-up questions," *Pacific Symposium on Biocomputing*, vol. 30, pp. 199–214, 2025.
- [158] Microsoft, "Phi-3-mini-4k-instruct," <https://huggingface.co/microsoft/Phi-3-mini-4k-instruct>, 2024, accessed: 2025-03-16.
- [159] M. N. Kamel Boulos and R. Dellavalle, "Nvidia's "chat with rtx" custom large language model and personalized ai chatbot augments the value of electronic dermatology reference material," *JMIR Dermatology*, vol. 7, p. e58396, Jul 24 2024.
- [160] Qwen Team, "Qwen/qwq-32b," <https://huggingface.co/Qwen/QwQ-32B>, 2024, accessed: 2025-03-16.
- [161] B. Wang, H. Zhao, H. Zhou, L. Song, M. Xu, W. Cheng, X. Zeng, Y. Zhang, Y. Huo, Z. Wang, Z. Zhao, D. Pan, F. Kou, F. Li, F. Chen, G. Dong, H. Liu, H. Zhang, J. He, J. Yang, K. Wu, K. Wu, L. Su, L. Niu, L. Sun, M. Wang, P. Fan, Q. Shen, R. Xin, S. Dang, S. Zhou, W. Chen, W. Luo, X. Chen, X. Men, X. Lin, X. Dong, Y. Zhang, Y. Duan, Y. Zhou, Z. Ma, and Z. Wu, "Baichuan-m1: Pushing the medical capability of large language models," <https://arxiv.org/abs/2502.09298>, 2025, arXiv:2502.09298 [cs.CL], v2.
- [162] E. Almazrouei, H. Alobeidli, A. Alshamsi, A. Cappelli, R. Cojocaru, M. Debbah, É. Goffinet, D. Hessel, J. Launay, Q. Malartic, D. Mazzotta, B. Noun, B. Pannier, and G. Penedo, "The falcon series of open language models," *arXiv preprint arXiv:2311.16867*, 2023. [Online]. Available: <https://arxiv.org/abs/2311.16867>
- [163] L. Tunstall, E. Beeching, N. Lambert, N. Rajani, K. Rasul, Y. Belkada, S. Huang, L. von Werra, C. Fourrier, N. Habib, N. Sarrazin, O. Sansevieri, A. M. Rush, and T. Wolf, "Zephyr: Direct distillation of lm alignment," 2023.
- [164] D. Moser, M. Bender, and M. Sariyar, "Generating synthetic healthcare dialogues in emergency medicine using large language models," *Studies in Health Technology and Informatics*, vol. 321, pp. 235–239, Nov 22 2024.
- [165] H. Huang, F. Yu, J. Zhu, X. Sun, H. Cheng, S. Dingjie, Z. Chen, M. Alharthi, B. An, J. He, Z. Liu, J. Chen, J. Li, B. Wang, L. Zhang, R. Sun, X. Wan, H. Li, and J. Xu, "AceGPT, Localizing Large Language Models in Arabic," in *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. Mexico City, Mexico: Association for Computational Linguistics, 2024, pp. 8139–8163.
- [166] T. G. T. Mesnard, C. Hardin et al., "Gemma: Open models based on gemini research and technology," *arXiv preprint arXiv:2403.08295*, 2024.
- [167] Y. Feng, J. Wang, R. He, L. Zhou, and Y. Li, "A retrieval-augmented knowledge mining method with deep thinking llms for biomedical research and clinical support," *arXiv preprint arXiv:2503.23029*, 2025.
- [168] S. Rau, A. Rau, J. Nattenmüller, A. Fink, F. Bamberg, M. Reiser, and M. F. Russe, "A retrieval-augmented chatbot based on gpt-4 provides appropriate differential diagnosis in gastrointestinal radiology: a proof of concept study," *European Radiology Experimental*, vol. 8, no. 1, p. 60, 2024. [Online]. Available: <https://doi.org/10.1186/s41747-024-00457-x>
- [169] J. An, D. Lim, Q. Le et al., "Chatgpt performance in assessing musculoskeletal mri scan appropriateness based on acr appropriateness criteria," *Scientific Reports*, vol. 15, p. 7140, 2025. [Online]. Available: <https://doi.org/10.1038/s41598-025-88925-1>
- [170] S. Puts, C. M. L. Zegers, A. Dekker, and I. Bermejo, "Developing an icd-10 coding assistant: Pilot study using roberta and gpt-4 for term extraction and description-based code selection," *JMIR Formative Research*, vol. 9, p. e60095, 2025. [Online]. Available: <https://doi.org/10.2196/60095>
- [171] H. Choi, D. Lee, and Y.-k. Kang, "Empowering pet imaging reporting with retrieval-augmented large language models and reading reports database: A pilot single center study," *medRxiv*, 2024. [Online]. Available: <https://www.medrxiv.org/content/early/2024/05/14/2024.05.13.24307312>
- [172] N. Markey, I. El-Mansouri, G. Rensonnet, C. van Langen, and C. Meier, "From rags to riches: Utilizing large language models to write documents for clinical trials," *Clinical Trials*, 2025, epub ahead of print.
- [173] T. Han, L. C. Adams, J.-M. Papaioannou, P. Grundmann, T. Oberhauser, A. Löser, D. Truhn, and K. K. Bressemer, "Medalpaca—an open-source collection of medical conversational ai models and training data," *arXiv preprint arXiv:2304.08247*, 2023.
- [174] C. Wu, W. Lin, X. Zhang et al., "Pmc-llama: Toward building open-source language models for medicine," *Journal of the American Medical Informatics Association*, vol. 31, no. 9, pp. 1833–1843, 2024.
- [175] Y. Labrak, A. Bazoge, E. Morin, P.-A. Gourraud, M. Rouvier, and R. Dufour, "Biomistral: A collection of open-source pretrained large language models for medical domains," 2024.
- [176] I. Buhnila, A. Sinha, and M. Constant, "Retrieve, generate, evaluate: A case study for medical paraphrases generation with small language models," <https://arxiv.org/abs/2407.16565>, 2024, arXiv:2407.16565 [cs.CL].
- [177] Z. Chen, A. Hernández-Cano, A. Romanou, A. Bonnet, K. Matoba, F. Salvi, M. Pagliardini, S. Fan, A. Köpf, A. Mohtashami, A. Sallinen, A. Sakhaeirad, V. Swamy, I. Krawczuk, D. Bayazit, A. Marmet, S. Montariol, M.-A. Hartley, M. Jaggi, and A. Bosselut, "Meditron-70b: Scaling medical pretraining for large language models," 2023.
- [178] H. Kim, H. Hwang, J. Lee, S. Park, D. Kim, T. Lee, C. Yoon, J. Sohn, D. Choi, and J. Kang, "Small language models learn enhanced reasoning skills from medical textbooks," *arXiv preprint arXiv:2404.00376*, 2024.
- [179] J. Sohn, Y. Park, C. Yoon, S. Park, H. Hwang, M. Sung, H. Kim, and J. Kang, "Rationale-guided retrieval augmented generation for medical question answering," <https://arxiv.org/abs/2411.00300>, 2024, arXiv:2411.00300 [cs.CL].
- [180] S. Siriwardhana, R. Weerasekera, E. Wen, T. Kaluarachchi, R. Rana, and S. Nanayakkara, "Improving the domain adaptation of retrieval augmented generation (rag) models for open domain question answering," *Transactions of the Association for Computational Linguistics*, vol. 11, pp. 1–17, 2023.
- [181] Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, M. Wang, and H. Wang, "Retrieval-augmented generation for large language models: A survey," *arXiv preprint arXiv:2312.10997*, 2023. [Online]. Available: <https://arxiv.org/abs/2312.10997>
- [182] Z. Huang, K. Xue, Y. Fan, L. Mu, R. Liu, T. Ruan, S. Zhang, and X. Zhang, "Tool calling: Enhancing medication consultation via retrieval-augmented large language models," 2024, arXiv:2404.17897 [cs.CL], licensed under CC BY 4.0. [Online]. Available: <https://arxiv.org/abs/2404.17897>
- [183] O. Bodenreider, "The unified medical language system (umls): integrating biomedical terminology," *Nucleic Acids Research*, vol. 32, no. suppl\_1, pp. D267–D270, 2004.
- [184] D. Jin, E. Pan, N. Oufattole, W.-H. Weng, H. Fang, and P. Szolovits, "What disease does this patient have? a large-scale open domain question answering dataset from medical exams," *arXiv preprint arXiv:2009.13081*, 2020.
- [185] L. Masanneck, S. G. Meuth, and M. Pawlitzki, "Evaluating base and retrieval augmented llms with document or online support for evidence based neurology," *npj Digital Medicine*, vol. 8, p. 137, 2025.



- [186] M. Sevgi, F. Antaki, and P. A. Keane, "Medical education with large language models in ophthalmology: custom instructions and enhanced retrieval capabilities," *British Journal of Ophthalmology*, vol. 108, no. 10, pp. 1354–1361, 2024, epub ahead of print.
- [187] A. E. Johnson, L. Bulgarelli, L. Shen, A. Gayles, A. Shammout, S. Horng, T. J. Pollard, S. Hao, B. Moody, B. Gow et al., "Mimic-iv, a freely accessible electronic health record dataset," *Scientific Data*, vol. 10, no. 1, p. 1, 2023.
- [188] D. Soong, S. Sridhar, H. Si, J. Wagner, A. Sá, C. Yu, K. Karagoz, M. Guan, S. Kumar, H. Hamadeh, and B. Higgs, "Improving accuracy of gpt-3/4 results on biomedical data using a retrieval-augmented language model," *PLOS Digital Health*, vol. 3, no. 8, p. e0000568, 2024.
- [189] N. Rekabsaz, O. Lesota, M. Schedl, J. Brassey, and C. Eickhoff, "Tripclick: The log files of a large health web search engine," in *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21)*. New York, NY, USA: Association for Computing Machinery, 2021, pp. 2507–2513.
- [190] Centers for Medicare & Medicaid Services, "ICD-10-PCS Official Guidelines for Coding and Reporting 2023," <https://www.cms.gov/files/document/2023-official-icd-10-pcs-coding-guidelines.pdf>, 2023, [Accessed 2024-11-20].
- [191] John Snow Labs, "Entity resolver for human phenotype ontology," [https://nlp.johnsnowlabs.com/2021/05/16/sbiobertresolve\\_HPO\\_en.html](https://nlp.johnsnowlabs.com/2021/05/16/sbiobertresolve_HPO_en.html), 2024, accessed August 30, 2024.
- [192] A. Albayrak, Y. Xiao, P. Mukherjee, S. S. Barnett, C. A. Marcou, and S. N. Hart, "Enhancing human phenotype ontology term extraction through synthetic case reports and embedding-based retrieval: A novel approach for improved biomedical data annotation," *Journal of Pathology Informatics*, vol. 16, p. 100409, Nov 16 2024.
- [193] D. S. Wishart, C. Knox, A. C. Guo, S. Shrivastava, M. Hassanali, P. Stothard, Z. Chang, and J. Woolsey, "Drugbank: A comprehensive resource for in silico drug discovery and exploration," *Nucleic Acids Research*, vol. 34, no. suppl 1, pp. D668–D672, 2006.
- [194] J. H. Morris, K. Soman, R. E. Akbas, X. Zhou, B. Smith, E. C. Meng, C.-C. Huang, G. Ceroni, G. Schenk, A. Rizk-Jackson, A. Harroud, L. Sanders, S. V. Costes, K. Bharat, A. Chakraborty, A. R. Pico, T. Mardrossian, M. Keiser, A. Tang, J. Hardi, Y. Shi, M. Musen, S. Israni, S. Huang, P. W. Rose, C. A. Nelson, and S. E. Baranzini, "The scalable precision medicine open knowledge engine (spoke): a massive knowledge graph of biomedical information," *Bioinformatics*, vol. 39, no. 2, p. btad080, 2023.
- [195] J. Irvin, P. Rajpurkar, M. Ko, Y. Yu, S. Ciurea-Illcus, C. Chute, H. Marklund, B. Haghighi, R. Ball, K. Shpanskaya et al., "Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, 2019, pp. 590–597.
- [196] O. Taboureaux, K. Audouze, A. B. Smit, C. Trevor, S. Whitebread, B.-K. Park, P. Goldsmith, B. Wiestler, T. S. Jensen, and S. Brunak, "Chempot: a disease chemical biology database," *Nucleic Acids Research*, vol. 39, no. suppl\_1, pp. D367–D372, 2010.
- [197] I. Segura-Bedmar, P. Martínez Fernández, and M. Herrero Zazo, "Semeval-2013 task 9: Extraction of drug-drug interactions from biomedical texts (ddiextraction 2013)," in *Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*. Atlanta, Georgia, USA: Association for Computational Linguistics, 2013, pp. 341–350.
- [198] R. I. Doğan, R. Leaman, and Z. Lu, "Ncbi disease corpus: a resource for disease name recognition and concept normalization," *Journal of biomedical informatics*, vol. 47, pp. 1–10, 2014.
- [199] Q. Jin, Y. Yang, Q. Chen, and Z. Lu, "Genegpt: Augmenting large language models with domain tools for improved access to biomedical information," *Bioinformatics*, vol. 40, no. 2, p. btad075, 2024.
- [200] G. Tsatsaronis, G. Balikas, P. Malakasiotis et al., "An overview of the bioasq large-scale biomedical semantic indexing and question answering competition," *BMC Bioinformatics*, vol. 16, pp. 1–28, 2015.
- [201] J. Li, Y. Sun, R. J. Johnson, D. Sciaky, C.-H. Wei, R. Leaman, A. P. Davis, C. J. Mattingly, T. C. Wieggers, and Z. Lu, "Biocreative v cdr task corpus: A resource for chemical disease relation extraction," *Database*, vol. 2016, 2016.
- [202] C. Eickhoff, F. Gmehlin, A. V. Patel, J. Boullier, and H. Fraser, "Dc3 – a diagnostic case challenge collection for clinical decision support," in *Proceedings of the 2019 ACM SIGIR International Conference on Theory of Information Retrieval (ICTIR '19)*. Santa Clara, CA, USA: Association for Computing Machinery, 2019.
- [203] A. E. Johnson, T. J. Pollard, S. J. Berkowitz, N. R. Greenbaum, M. P. Lungren, C.-y. Deng, R. G. Mark, and S. Horng, "Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports," *Scientific Data*, vol. 6, no. 1, p. 317, 2019.
- [204] A. Ben Abacha and D. Demner-Fushman, "A question-entailment approach to question answering," *BMC Bioinform.*, vol. 20, no. 1, pp. 511:1–511:23, 2019. [Online]. Available: <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-019-3119-4>
- [205] Q. Jin, B. Dhingra, Z. Liu, W. Cohen, and X. Lu, "PubMedQA: A dataset for biomedical research question answering," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, November 2019, pp. 2567–2577.
- [206] V. Liévin, C. E. Hother, A. G. Motzfeldt, and O. Winther, "Can large language models reason about medical questions?" *Patterns (N Y)*, vol. 5, no. 3, p. 100943, Mar. 2024.
- [207] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt, "Measuring massive multitask language understanding," *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [208] K. Singhal, S. Azizi, T. Tu et al., "Large language models encode clinical knowledge," *arXiv preprint arXiv:2212.13138*, 2022. [Online]. Available: <https://doi.org/10.48550/arXiv.2212.13138>
- [209] A. Pal, L. K. Umapathi, and M. Sankarasubbu, "MedMCQA: A large-scale multi-subject multi-choice dataset for medical domain question answering," in *Proceedings of the Conference on Health, Inference, and Learning*, vol. 174. PMLR, 2022, pp. 248–260.
- [210] Y. Li, Z. Li, K. Zhang, R. Dan, S. Jiang, and Y. Zhang, "Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge," *Cureus*, vol. 15, no. 6, 2023.
- [211] W. Hou and Z. Ji, "Geneturing tests gpt models in genomics," <https://doi.org/10.1101/2023.03.11.532238>, 2023, bioRxiv preprint.
- [212] M. Li, M. Chen, H. Zhou, and R. Zhang, "Petaylor: Improving large language model by tailored chunk scorer in biomedical triple extraction," *arXiv preprint arXiv:2310.18463*, 2023. [Online]. Available: <https://arxiv.org/abs/2310.18463>
- [213] S. Liu, A. B. McCoy, and A. Wright, "Improving large language model applications in biomedicine with retrieval-augmented generation: a systematic review, meta-analysis, and clinical development guidelines," *Journal of the American Medical Informatics Association*, vol. 32, no. 4, pp. 605–615, April 2025. [Online]. Available: <https://doi.org/10.1093/jamia/ocaf008>
- [214] R. T. Sutton, D. Pincock, D. C. Baumgart, D. C. Sadowski, R. N. Fedorak, and K. I. Kroeker, "An overview of clinical decision support systems: benefits, risks, and strategies for success," *NPJ Digital Medicine*, vol. 3, p. 17, Feb 2020.

- [215] C. Martínez-deMiguel, I. Segura-Bedmar, E. Chacón-Solano, and S. Guerrero-Aspizua, "The raredis corpus: A corpus annotated with rare diseases, their signs and symptoms," *Journal of Biomedical Informatics*, vol. 125, p. 103961, 2022, epub 2021 Dec 5.
- [216] J. Yang, L. Shu, H. Duan, and H. Li, "RDguru: A conversational intelligent agent for rare diseases," *IEEE Journal of Biomedical and Health Informatics*, vol. PP, 2024, epub ahead of print.
- [217] J. Kim, K. Wang, C. Weng, and C. Liu, "Assessing the utility of large language models for phenotype-driven gene prioritization in the diagnosis of rare genetic disease," *American Journal of Human Genetics*, vol. 111, no. 10, pp. 2190–2202, 2024, epub ahead of print.
- [218] P. Sloan, P. Clatworthy, E. Simpson, and M. Mirmehdi, "Automated radiology report generation: A review of recent advances," *IEEE Reviews in Biomedical Engineering*, 2024.
- [219] I. E. Hamamci, S. Er, and B. Menze, "Ct2rep: Automated radiology report generation for 3d medical imaging," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2024, pp. 476–486.
- [220] O. Alfarghaly, R. Khaled, A. Elkorany, M. Helal, and A. Fahmy, "Automated radiology report generation using conditioned transformers," *Informatics in Medicine Unlocked*, vol. 24, p. 100557, 2021.
- [221] M. Ranjit, G. Ganapathy, R. Manuel, and T. Ganu, "Retrieval augmented chest x-ray report generation using openai gpt models," in *Machine Learning for Healthcare Conference*. PMLR, 2023, pp. 650–666.
- [222] H. M. T. Alam, D. Srivastav, M. A. Kadir, and D. Sonntag, "Towards interpretable radiology report generation via concept bottlenecks using a multi-agent rag," *arXiv preprint arXiv:2412.16086*, 2024.
- [223] M. L. Bernardi and M. Cimitile, "Report generation from x-ray imaging by retrieval-augmented generation and improved image-text matching," in *2024 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2024, pp. 1–8.
- [224] S. Song, A. Subramanyam, I. Madejski, and R. L. Grossman, "Lab-rag: Label boosted retrieval augmented generation for radiology report generation," *arXiv preprint arXiv:2411.16523*, 2024.
- [225] L. Sun, J. Zhao, M. Han, and C. Xiong, "Fact-aware multimodal retrieval augmentation for accurate medical radiology report generation," *arXiv preprint arXiv:2407.15268*, 2024.
- [226] H. Burton, T. Cole, and A. M. Lucassen, "Genomic medicine: challenges and opportunities for physicians," *Clinical Medicine (London)*, vol. 12, no. 5, pp. 416–419, October 2012, erratum in: *Clinical Medicine*. 2012 Dec;12(6):504.
- [227] S. Quazi and J. A. Malik, "A systematic review of personalized health applications through human-computer interactions (hci) on cardiovascular health optimization," *Journal of Cardiovascular Development and Disease*, vol. 9, no. 8, p. 273, Aug. 2022.
- [228] K. Lemmetty, T. Kuusela, K. Saranto, and A. Ensio, "Education and training of health information systems—a literature review," *Studies in Health Technology and Informatics*, vol. 122, pp. 176–180, 2006.
- [229] Q. Zhao, H. Wang, R. Wang, and H. Cao, "Deriving insights from enhanced accuracy: Leveraging prompt engineering in custom gpt for assessing chinese nursing licensing exam," *Nurse Education in Practice*, vol. 84, p. 104284, 2025. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S147159532500040X>
- [230] L. Ge, R. Agrawal, M. Singer *et al.*, "Leveraging artificial intelligence to enhance systematic reviews in health research: advanced tools and challenges," *Systematic Reviews*, vol. 13, p. 269, 2024. [Online]. Available: <https://doi.org/10.1186/s13643-024-02682-2>
- [231] G. Deng, Y. Liu, K. Wang, Y. Li, T. Zhang, and Y. Liu, "Pandora: Jailbreak gpts by retrieval augmented generation poisoning," *arXiv preprint arXiv:2402.08416*, 2024.
- [232] A. Névél, H. Dalianis, S. Velupillai *et al.*, "Clinical natural language processing in languages other than english: opportunities and challenges," *Journal of Biomedical Semantics*, vol. 9, no. 12, 2018. [Online]. Available: <https://doi.org/10.1186/s13326-018-0179-8>
- [233] Y. Jin, M. Chandra, G. Verma, Y. Hu, M. De Choudhury, and S. Kumar, "Better to ask in english: Cross-lingual evaluation of large language models for healthcare queries," in *Proceedings of the ACM Web Conference 2024*, ser. WWW '24. New York, NY, USA: Association for Computing Machinery, 2024, p. 2627–2638. [Online]. Available: <https://doi.org/10.1145/3589334.3645643>
- [234] I. Mondshine, T. Paz-Argaman, and R. Tsarfay, "Beyond english: The impact of prompt translation strategies across languages and tasks in multilingual llms," 2025, submitted on 13 Feb 2025. [Online]. Available: <https://arxiv.org/abs/2502.09331>
- [235] O. Pelka, S. Koitka, J. Rückert, F. Nensa, and C. M. Friedrich, "Radiology objects in context (roco): A multimodal image dataset," in *Intravascular Imaging and Computer Assisted Stenting and Large-Scale Annotation of Biomedical Data and Expert Label Synthesis*. Springer International Publishing, 2018, pp. 180–189.
- [236] W. Lin, Z. Zhao, X. Zhang, C. Wu, Y. Zhang, Y. Wang, and W. Xie, "Pmc-clip: Contrastive language-image pre-training using biomedical documents," in *Medical Image Computing and Computer Assisted Intervention*. Springer, 2023, pp. 525–536.
- [237] L. Fraiwan, M. Fraiwan, B. Khassawneh, and A. Ibnian, "A dataset of lung sounds recorded from the chest wall using an electronic stethoscope," *Data in Brief*, vol. 35, p. 106913, 2021.
- [238] D. Gupta, K. Attal, and D. Demner-Fushman, "A dataset for medical instructional video classification and question answering," *Scientific Data*, vol. 10, p. 158, 2023. [Online]. Available: <https://doi.org/10.1038/s41597-023-02036-y>
- [239] K. Zhong, L. Zhang, L. Wang, X. Shu, and Z. Wang, "Class-center-based self-knowledge distillation: A simple method to reduce intra-class variance," *Applied Sciences*, vol. 14, no. 16, p. 7022, 2024. [Online]. Available: <https://doi.org/10.3390/app14167022>
- [240] DeepSeek-AI, "Deepseek-v3 technical report," 2024. [Online]. Available: <https://arxiv.org/abs/2412.19437>
- [241] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 1877–1901.
- [242] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cassirer, H. Jones, K. Tunyasuvunakool, G. Petrides, J.-B. Jones, J. Bradbury *et al.*, "Training compute-optimal large language models," *arXiv preprint arXiv:2203.15556*, 2022.
- [243] X. Xu, M. Li, C. Tao, T. Shen, R. Cheng, J. Li, C. Xu, D. Tao, and T. Zhou, "A survey on knowledge distillation of large language models," *arXiv preprint arXiv:2402.13116*, 2024.
- [244] R. M. S. Khan, P. Li, S. Yun, Z. Wang, S. Nirjon, C.-W. Wong, and T. Chen, "Portllm: Personalizing evolving large language models with training-free and portable model patches," *arXiv preprint arXiv:2410.10870*, 2024. [Online]. Available: <https://arxiv.org/abs/2410.10870>
- [245] S. M. Meystre, F. J. Friedlin, B. R. South, S. Shen, and M. H. Samore, "Automatic de-identification of textual documents in the electronic health record: a review of recent research," *BMC Medical Research Methodology*, vol. 10, no. 1, p. 70, 2010.
- [246] I. Consulting, "General data protection regulation," 2016, [Online; accessed March 26, 2025]. [Online]. Available: <https://gdpr-info.eu/>

[247] A. Act, "Health insurance portability and accountability act of 1996," p. 191, 1996, [Online; accessed March 26, 2025].  
[Online]. Available: <https://www.hhs.gov/hipaa/>