

A Comprehensive Survey of Large Language Models and Multimodal Large Language Models in Medicine

Hanguang Xiao^{a,*}, Feizhong Zhou^{a,1}, Xingyue Liu^a, Tianqi Liu^a, Zhipeng Li^a, Xin Liu^a and Xiaoxuan Huang^a

^a*School of Artificial Intelligence, Chongqing University of Technology, Chongqing 401120, China*

ARTICLE INFO

Keywords:

Large language model
Multimodal large language model
Medicine
Healthcare
Clinical application

ABSTRACT

Since the release of ChatGPT and GPT-4, large language models (LLMs) and multimodal large language models (MLLMs) have attracted widespread attention for their exceptional capabilities in understanding, reasoning, and generation, introducing transformative paradigms for integrating artificial intelligence into medicine. This survey provides a comprehensive overview of the development, principles, application scenarios, challenges, and future directions of LLMs and MLLMs in medicine. Specifically, it begins by examining the paradigm shift, tracing the transition from traditional models to LLMs and MLLMs, and highlighting the unique advantages of these LLMs and MLLMs in medical applications. Next, the survey reviews existing medical LLMs and MLLMs, providing detailed guidance on their construction and evaluation in a clear and systematic manner. Subsequently, to underscore the substantial value of LLMs and MLLMs in healthcare, the survey explores five promising applications in the field. Finally, the survey addresses the challenges confronting medical LLMs and MLLMs and proposes practical strategies and future directions for their integration into medicine. In summary, this survey offers a comprehensive analysis of the technical methodologies and practical clinical applications of medical LLMs and MLLMs, with the goal of bridging the gap between these advanced technologies and clinical practice, thereby fostering the evolution of the next generation of intelligent healthcare systems.

1. Introduction

The introduction of the Transformer [1] has revolutionized the fields of Natural Language Processing (NLP) and Computer Vision (CV). The Transformer's robust parallel computing capabilities and self-attention mechanism facilitate the integration of vast training datasets, forming the foundation for LLM and MLLM development. A variety of Transformer-based LLMs and MLLMs have emerged to date, with this survey primarily focusing on the vision-language modality. Notable examples include the PaLM series [2, 3], GPT series [4, 5], and LLaMA series [6, 7, 8] among LLMs, and Gemini [9], GPT-4 [10], and LLaVA [11] among MLLMs. Their exceptional capabilities in understanding, reasoning, and generation have enabled them to achieve state-of-the-art performance across various downstream tasks, including text generation, machine translation, and visual question answering (VQA). LLMs and MLLMs exhibit increasingly robust generalization abilities, significantly impacting the medical domain and accelerating the convergence of artificial intelligence and medicine [12]. Notably, Google's Med-PaLM 2 [13] scored 86.5 on the United States Medical Licensing Examination (USMLE) [14], achieving expert-level performance [15] and further highlighting the immense potential of LLMs in medicine. Additionally, emerging medical LLMs and MLLMs, such as ChatDoctor [16], ChatCAD [17], and LLaVA-Med [18], represent novel opportunities enabled by artificial intelligence in the medical field. These models offer promising solutions for medical report generation [19, 20], clinical diagnosis [17, 21], mental health services [22, 23], and various other clinical applications.

Despite the significant academic breakthroughs of LLMs and MLLMs in the medical field, hospitals still face notable challenges in training their own medical LLMs and MLLMs and deploying them in practical clinical applications. First, training medical LLMs and MLLMs requires a substantial amount of medical data, which is costly to obtain, demands annotation by medical experts, and raises significant concerns about data privacy [24]. Second, the parameters and computational demands of LLMs and MLLMs are substantial, necessitating significant computational resources for their training and deployment [25, 26], substantially increasing the adoption barrier for

*Corresponding author

simenxiao1211@163.com (Hanguang Xiao)

ORCID(s): 0000-0002-4359-7455 (Hanguang Xiao)

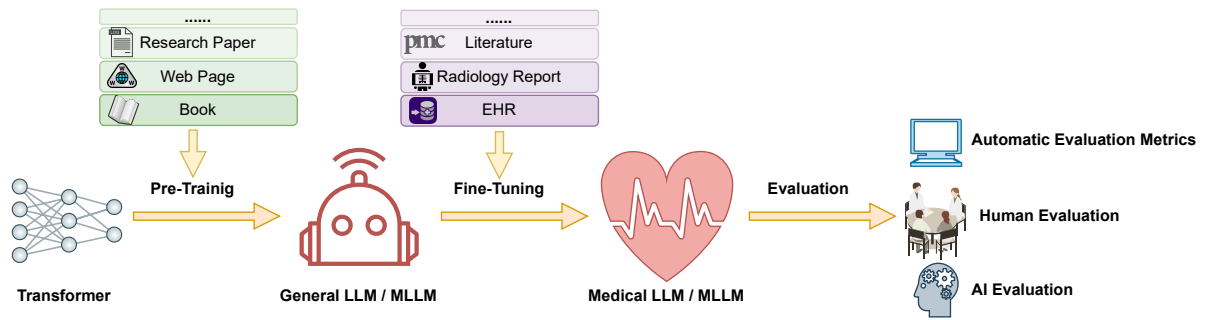


Fig. 1. The process of constructing and evaluating medical LLMs and MLLMs.

hospitals. Third, given the unique characteristics of the medical field, it is essential not only to evaluate the performance of LLMs and MLLMs on benchmarks but also to assess their instruction-following ability [27, 5, 11], safety, and ethical considerations, necessitating additional training and evaluation strategies to enhance and measure the models' performance across multiple dimensions. Furthermore, the development of LLMs and MLLMs in the medical field is still in its infancy, with many of their potential application scenarios remaining undefined. Moreover, they face a range of challenges, including hallucinations [28, 29, 30] and a lack of up-to-date information [12], which significantly impede their practical use in clinical settings.

To address the aforementioned challenges, this survey begins by tracing the evolution of LLMs and MLLMs through the lens of paradigm shifts. Subsequently, it reviews existing medical LLMs and MLLMs, summarizing their structural characteristics. The survey then gathers datasets suitable for training medical LLMs and MLLMs and elaborates on methods for training and evaluating these models, as shown in Fig. 1. Furthermore, to highlight the significant potential impact of LLMs and MLLMs in medicine, this survey summarizes their applications in clinical practice and analyzes current limitations and potential solutions. Finally, the survey explores the future directions of medical LLMs and MLLMs, offering forward-looking and insightful perspectives.

Medicine is a multimodal field [31, 32], making the study of medical MLLMs particularly important, as they can integrate and analyze information from various modalities to enhance clinical decision support, disease diagnosis, and treatment planning. However, the articles relevant to this survey mainly focus on medical LLMs and lack a detailed examination of medical MLLMs [15, 33, 34]. Additionally, most articles focus on the applications and impacts of LLMs in medicine but lack detailed discussions of technical aspects [12, 35, 26, 36, 37], such as datasets, model architectures, and construction methods. In contrast, this survey not only examines the background and principles of LLMs and MLLMs but also explores their applications and impacts in medicine, offering a clear logical structure and substantial depth and breadth. In summary, the contributions of this survey are as follows:

- This survey offers a thorough overview of medical LLMs and MLLMs, starting with an examination of their developmental background and architectural frameworks. Building on this foundation, it catalogs existing medical LLMs and MLLMs while providing a detailed analysis of their structural variations and key components.
- This survey systematically elucidates the complete process of medical LLMs and MLLMs, from training to evaluation, covering fine-tuning methods, evaluation strategies, and relevant medical datasets. Additionally, it highlights how to select appropriate datasets, fine-tuning methods, and evaluation strategies to assist researchers in the rapid development of medical LLMs and MLLMs.
- This survey summarizes the applications, challenges, and potential solutions of medical LLMs and MLLMs in clinical practice, while providing a forward-looking analysis of future developmental trajectories. It seeks to offer visionary perspectives that inspire advancements in the field, benefiting medical professionals and researchers alike.

This survey aims to advance the development of LLMs and MLLMs for clinical medicine applications, thereby promoting deeper integration between artificial intelligence and healthcare. The structure of this survey is outlined in Fig. 2: Section 2 reviews the development background of LLMs and MLLMs. Section 3 describes the architectures of current LLMs and MLLMs and highlights the structural differences among them. Section 4 covers the datasets for

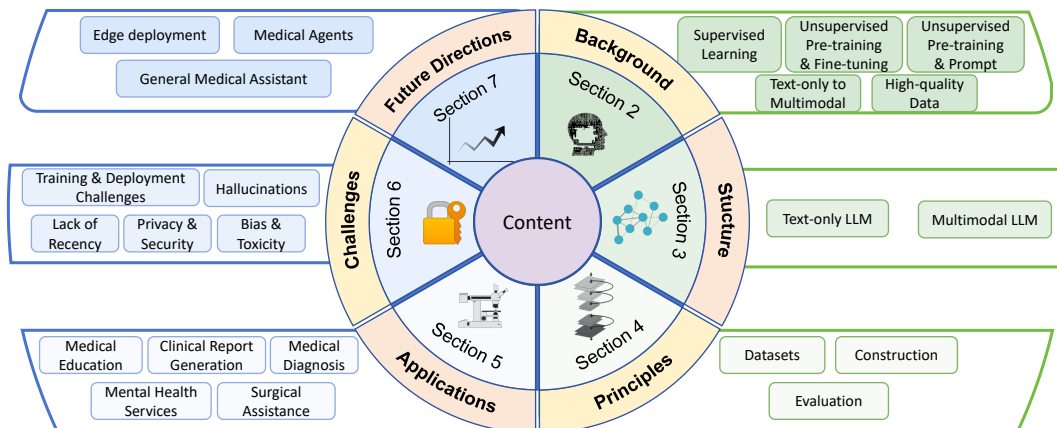


Fig. 2. The overall structure of the survey. Section 2 to Section 4 are biased toward principles of medical LLMs and MLLMs; Section 5 to Section 7 are biased toward practical clinical applications.

training medical LLMs and MLLMs, along with an overview of training and evaluation methods. Section 5 examines current potential applications of medical LLMs and MLLMs. Section 6 discusses the challenges and limitations of LLMs and MLLMs in clinical settings and proposes potential solutions. Section 7 offers a forward-looking perspective on the future of medical LLMs and MLLMs. Finally, Section 8 concludes the survey. In summary, readers interested in the foundational knowledge and principles of medical LLMs and MLLMs should refer to Section 2 to Section 4, while those focused on applications, challenges, and future directions should consult Section 5 and Section 7.

2. Background of LLMs and MLLMs

This section focuses on paradigm shifts, categorizing the development of NLP into four distinct stages, as illustrated in Fig. 3: (1) Supervised Learning; (2) Unsupervised Pre-training and Fine-tuning; (3) Unsupervised Pre-training and Prompt; (4) Text-only to Multimodal. Recent research [38] underscores the importance of high-quality datasets for LLMs and MLLMs. Accordingly, we introduce (5) High-quality Data, which examines the shift from reliance on large-scale data to an emphasis on high-quality data.

2.1. Supervised Learning

Supervised learning is a fundamental paradigm in machine learning that focuses on minimizing a loss function. The objective can be expressed as follows:

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x_i; \theta), y_i) + \lambda \Omega(\theta) \quad (1)$$

where the first term represents the empirical risk, and the second term represents the regularization term. In supervised learning, a model is trained to map input variables x to output variables y by minimizing the discrepancy between $f(x; \theta)$ and y , where θ denotes model parameters, x may consist of manually extracted features or raw text, and y represents supervision signals such as category labels, text, or other forms.

Before the advent of pre-training methods, the supervised learning paradigm dominated the NLP field. Early NLP relied heavily on feature engineering [39], requiring researchers to extract and select features from datasets to perform tasks like text classification [40] and machine translation [41]. The advent of deep learning [42] enabled end-to-end model training, shifting research focus from feature engineering to model architecture design, with CNN and LSTM models emerging as prominent approaches. The supervised learning era in NLP marked a shift from feature selection to model architecture design, signifying a transition from feature engineering to structure engineering.

2.2. Unsupervised Pre-training and Fine-tuning

Supervised learning depends on annotated datasets, which establish clear standards for model optimization. However, acquiring sufficient annotated data is challenging for certain tasks, particularly in medical domains, due

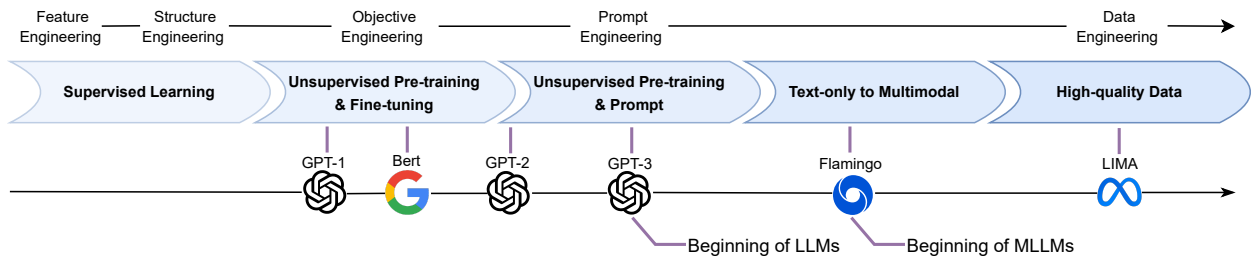


Fig. 3. Evolution of LLMs and MLLMs. The Evolution of LLMs and MLLMs. The upper section illustrates the research focuses and paradigm shifts across the evolution of these models, while the lower section highlights key milestones achieved at each stage.

to the scarcity of specialized annotators and the complexity of the annotation process [24]. The introduction of Transformer [1] revolutionized the NLP learning paradigm, rendering supervised learning increasingly marginalized [39].

Transformer-based models like GPT [43] and BERT [44] achieved state-of-the-art results through unsupervised pre-training on large-scale unlabeled text, followed by supervised fine-tuning (SFT) using task-specific objective functions. The emergence of GPT and BERT introduced a new NLP paradigm: unsupervised pre-training & fine-tuning. This paradigm also revolutionized language model development in the medical field, giving rise to prominent models like PubMedBERT [45] and BioBERT [46]. Compared to earlier models, these approaches offer several advantages: (1) Pre-training data can be drawn from any unannotated text corpus, such as biomedical literature, mitigating challenges related to limited annotated data in medical domains; (2) Training on large-scale unlabeled data enables the acquisition of general and abstract language representations, improving generalization; (3) Fine-tuning for downstream tasks requires only task-specific objective functions, eliminating extensive architectural modifications and facilitating a transition from structural to objective engineering.

2.3. Unsupervised Pre-training and Prompt

While GPT and BERT achieved state-of-the-art results in tasks like machine translation, sentiment analysis, and question-answering (QA), they still required task-specific fine-tuning for each downstream task. To develop a general language model capable of handling diverse tasks without task-specific fine-tuning, Radford et al. [47] curated a dataset of over 8 million documents, totaling 40 GB of text data, encompassing examples from multiple domains and tasks, and used it to train GPT-2. GPT-2 set state-of-the-art benchmarks on 7 out of 8 language modeling tasks without requiring task-specific fine-tuning.

To further enhance the generalization capabilities of language models, Brown et al. [4] scaled the model to 175 billion parameters and significantly expanded the training dataset. This resulted in GPT-3, which demonstrated a qualitative leap in performance, showcasing remarkable few-shot capabilities without requiring fine-tuning. GPT-3 could handle unfamiliar tasks based solely on provided examples, often achieving performance on par with fine-tuned state-of-the-art models. As a result, GPT-3 is widely regarded as the beginning of LLMs [35]. The proposal of GPT-3 further revolutionized NLP by shifting the paradigm from unsupervised pre-training & fine-tuning to unsupervised pre-training & prompt [39]. Such models can handle most tasks effectively using user prompts and contextual examples. For example, Flan-PaLM [13] achieved state-of-the-art performance on MedQA, MedMCQA, PubMedQA, and MMLU clinical benchmarks by employing advanced prompting strategies. However, these models are highly sensitive to user-provided prompts, which directly influence the quality of their responses. This sensitivity has spurred researchers to explore prompting strategies in depth [48], initiating a shift from objective engineering to prompt engineering.

2.4. Text-only to Multimodal

Inspired by GPT-3, researchers have intensified efforts in developing LLMs, resulting in prominent works such as GLM-130B [49], PaLM [2, 3], and LLaMA [6, 7, 8]. However, these LLMs remain text-focused, and despite progress in multimodal research, they often require fine-tuning for new tasks [50, 51] or lack text generation capabilities [52, 53], which restricts their application scope. Inspired by few-shot learners like GPT-3, Alayrac et al. [54] curated a large-scale multimodal dataset from the web, comprising primarily text-image pairs, to train an MLLM named Flamingo. Flamingo adapts seamlessly to visual tasks via few-shot learning, eliminating the need for task-specific fine-tuning.

Flamingo's robust multimodal in-context learning (ICL) and few-shot capabilities mark it as the GPT-3 moment for the multimodal domain, positioning it as the beginning of MLLMs. These MLLMs employ LLMs as cognitive engines, preserving their inherent capabilities while integrating powerful visual functionalities [55]. This advancement paves the way for the development of generalist medical AI systems. For example, Med-PaLM M [21] achieved performance levels comparable to or exceeding the state-of-the-art expert models across 14 different medical tasks, showcasing the potential of MLLMs as general medical assistants.

2.5. High-quality Data

A key factor behind the success of LLMs and MLLMs is their use of large-scale training data, which enables the acquisition of universal representations transferable to diverse language understanding and generation tasks [38]. However, much of this training data comes from web sources like WebText [47] and Common Crawl, and it is inevitable that there are some toxicities and biases in these large amounts of web data, which are also carried over to LLMs and MLLMs [25]. To mitigate the adverse effects of large-scale datasets and improve model performance, researchers often employ high-quality datasets for fine-tuning.

For example, Li et al. [18] utilized GPT-4 to create open-ended instruction-following data derived from biomedical image-caption datasets, subsequently training LLaVA-Med on this data. LLaVA-Med exhibited remarkable multimodal conversational capabilities, adeptly addressing biomedical image queries based on user instructions. To further refine multimodal medical datasets, Xie et al. [56] introduced a dataset enriched with multi-grained annotations. These annotations encompass global context, modality details, and localized descriptions of medical images. Training LLaVA-Med on this dataset led to a nearly 10% average performance boost across three biomedical VQA datasets, emphasizing the critical role of high-quality datasets. Notably, LIMA [38] fine-tuned using only 1,000 meticulously curated prompts and responses with standard supervised loss, outperformed both Alpaca[57] and Bard in human and GPT-4 preference scores. Ablation experiments on LIMA demonstrated that enhancing data quality provides greater benefits than merely increasing data quantity [38]. Consequently, it can be predicted that data engineering is emerging as a pivotal area of research focus.

This section outlines the development trajectory of LLMs, emphasizing the shift from supervised learning to unsupervised pre-training & fine-tuning, and ultimately to unsupervised pre-training & prompting. The success of LLMs has catalyzed rapid growth in the multimodal domain, leading to the development of MLLMs built on LLM foundations. Notably, GPT-3 and Flamingo, with their robust few-shot capabilities, mark the beginning of LLMs and MLLMs, respectively. Recent studies highlighting the role of high-quality datasets in LLMs and MLLMs suggest that data engineering will emerge as a key research focus. In summary, the evolution of LLMs and MLLMs reflects a progressive shift from initial feature engineering to structure engineering, objective engineering, and now, to prompt and data engineering.

3. Structure of LLMs and MLLMs

Existing LLMs are universally based on the Transformer architecture, which adopts an encoder-decoder framework. Accordingly, these LLMs have evolved into three structural variants based on the Transformer architecture: (1) Encoder-only, represented by models such as BERT [44]; (2) Decoder-only, represented by models such as the GPT series [4, 5]; (3) Encoder-Decoder, represented by models like T5 [58]. Current MLLMs extend LLMs by integrating a vision encoder to process visual information and a modality alignment module [55, 59] to bridge the gap between vision and text modalities. This section provides a comprehensive overview of existing medical LLMs and MLLMs, focusing on their respective model architectures. Section 3.1 reviews medical LLMs categorized by the three aforementioned structures. Section 3.2 discusses the common vision encoders, LLM backbones, and modality alignment techniques employed in medical MLLMs. For clarity, Table 1 and Table 2 detail and categorize existing medical LLMs and MLLMs.

3.1. Structure of LLMs

3.1.1. Encoder-only

Encoder-only language models (LMs) consist of multiple encoder layers within the Transformer architecture, with BERT being the earliest and most representative example. Inspired by BERT, additional encoder-only LMs such as DeBERTa DeBERTa [85], ALBERT [86], and RoBERTa [87] have been developed. Encoder-only LMs commonly

Table 1

Detailed information on existing medical LLMs categorized by architecture type.

Category	Model Name	Base Model	Para.(B)	Training Data Source	Construction Method	Evaluation Method	Date
Decoder-Only	Med-PaLM [60]	PaLM	540	MultiMedQA	IFT	AEM, Human	2022/12
	ChatDoctor [16]	LLaMA	7	Alpaca-52k, HealthCareMagic-100k	IFT	AI	2023/03
	Baize-Healthcare [61]	LLaMA	7	Quora, MedQuAD	SFT	AI	2023/04
	BenTsao [62]	LLaMA	7	CMcKG	SFT	Human	2023/04
	MedAlpaca [63]	LLaMA	7 / 13	Medical Meadow	IFT	AEM	2023/04
	PMC-LLaMA [64]	LLaMA	7 / 13	MedC-K, MedC-I	CPT, IFT	AEM	2023/04
	Med-PaLM 2 [13]	PaLM 2	340	MultiMedQA	IFT	AEM, Human	2023/05
	Clinical Camel [65]	LLaMA 2	13 / 70	ShareGPT, PubMed, MedQA	SFT	AEM	2023/05
	HuatuoGPT [66]	BLOOMZ	7	Hybrid Data	SFT, RLAI	AEM, Human, AI	2023/05
	GatorTronGPT [67]	GPT-3	5 / 20	Clinical Text from UF Health, Pile	PT	AEM	2023/06
	ClinicalGPT [68]	BLOOM	7	Three MedQA, MD-EHR, MedDialog	SFT, RLHF	AEM	2023/06
	Zhongjing [69]	Ziya-LLaMA	13	CMtMedQA, ChatMed, CMcKG	CPT, SFT, RLHF	Human, AI	2023/08
	Radiology-Llama2 [70]	LLaMA 2	7	MIMIC-CXR, OpenI	IFT	AEM, Human	2023/08
	MedChatZH [71]	Baichuan	7	Books, med-mix-2M	CPT, IFT	AEM, AI	2023/09
	CPLLM [72]	LLaMA 2	13	eICU-CRD, MIMIC-IV	IFT	AEM	2023/09
	ChatCounselor [23]	Vicuna	7	Psych8k	IFT	AI	2023/09
	Qilin-Med [73]	Baichuan	7	ChiMed	CPT, SFT, DPO	AEM	2023/10
	AlpaCare [74]	LLaMA	7 / 13	MedInstruct-52k	IFT	AI	2023/10
	TCM-GPT [75]	BLOOM	7	TCM-Corpus-1B, TCM-EXAM, TCM-EHR	CPT, SFT	AEM	2023/11
	HuatuoGPT-II [76]	Baichuan 2	7 / 13	Web Corpus, Books, Literature, Encyclopedia	IFT	AEM, Human, AI	2023/11
	MEDITRON [77]	LLaMA 2	7 / 70	GAP-Replay, MedMCQA, PubMedQA, MedQA	CPT, SFT	AEM	2023/11
	AMIE [21]	PaLM 2	340	MedQA, MultiMedBench, MIMIC-III, RealWorld Dialogue	IFT	Human, AI	2024/01
	BioMistral [78]	Mistral	7	PubMed Central	CPT, SFT	AEM	2024/02
	Me-LLaMA [79]	LLaMA 2	13 / 70	Pile, MIMIC-III, MIMIC-IV, MIMIC-CXR, RedPajama	CPT, IFT	AEM	2024/02
	Apollo [80]	Qwen	7	ApolloCorpora	CPT, SFT	AEM	2024/03
BioMedLM [81]	Transformer	2.7	PubMed Center, Pile	PT, SFT	AEM	2024/03	
PediatricsGPT [82]	Baichuan 2	7 / 13	PedCorpus	CPT, SFT, DPO	AEM, Human, AI	2024/05	
Encoder-Decoder	DoctorGLM [83]	ChatGLM	6	ChatDoctor, HealthcareMagic, MedDialog, CMD.	IFT	Human	2023/04
	BianQue [84]	ChatGLM	6	BianQueCorpus	IFT	AEM	2023/10
	SoulChat [22]	ChatGLM	6	SoulChatCorpus	IFT	AEM, Human	2023/11

¹ Encoder-only models are not included as they typically belong to the PLM, not LLM.² "CPT" means continuous pre-training, "IFT" means instruction fine-tuning, "SFT" means supervised fine-tuning, "RLHF" means reinforcement learning from human feedback, "RLAI" means reinforcement learning from AI feedback, "DPO" means direct preference optimization.³ "AEM" means automatic evaluation metrics.

utilize the masked language modeling (MLM) task during pre-training, where random tokens in sentences are masked, and the model is trained to predict these tokens accurately. This pre-training approach equips encoder-only LMs with exceptional natural language understanding capabilities, allowing them to effectively encode and comprehend medical knowledge, thereby improving performance in various medical tasks. Consequently, researchers have focused on developing dedicated encoder-only LMs tailored specifically for the medical domain [46, 88, 45]. For example, BioBERT [46], pre-trained on biomedical corpora, achieved state-of-the-art results in tasks such as biomedical named entity recognition, relation extraction, and QA. MentalBERT [88] was trained on datasets of mental health disorders (e.g., depression, anxiety, suicidal ideation) sourced from social platforms like Reddit and Twitter, facilitating its application in mental health research.

Despite the presence of numerous encoder-only LMs in the medical domain, these models are better classified as pre-trained language models (PLMs) [33, 34] rather than LLMs. This distinction arises because they require fine-tuning for downstream tasks, lacking the robust ICL and few-shot capabilities of models like GPT-3. Therefore, these PLMs will not be further addressed in the following sections.

3.1.2. Decoder-only

Decoder-only models are the dominant architecture for LLMs, comprising multiple decoder layers within the Transformer. The first decoder-only language model was GPT, and GPT-3 later marked a beginning of LLMs, paving the way for numerous other notable decoder-only models [2, 3, 5, 6, 7]. Decoder-only LLMs primarily use next-token prediction (NTP) as their pre-training objective. In this process, the model learns to predict the next token in a sequence based on all preceding tokens. This training paradigm equips decoder-only LLMs with remarkable generative capabilities, enabling them to convert discriminative tasks into generative ones. This unification of task formats enhances both their generalization and adaptability across application scenarios. For example, Med-PaLM M [21] excels in a range of tasks, such as text-based QA, VQA, image classification, radiology report generation, and summarization, achieving or exceeding state-of-the-art performance.

Compared to encoder-only LMs, these decoder-only LLMs utilize NTP as the pre-training task, which enhances their proficiency in text generation. Additionally, studies [89, 90] have demonstrated that decoder-only LLMs exhibit the best few-shot and zero-shot performance across diverse downstream tasks, which is one of the reasons why decoder-only has become the predominant framework for LLMs today.

3.1.3. Encoder-Decoder

Encoder-decoder LLMs leverage the Transformer architecture, combining stacks of encoders and decoders. The encoder processes input sequence and outputs representations with contextual information, which the decoder uses for text generation. Prominent examples of encoder-decoder LLMs T5 [58] and GLM [91]. Similar to encoder-only and decoder-only architectures, encoder-decoder LLMs have also been adapted for medical applications. For example, Chen et al. [22] fine-tuned ChatGLM using the empathetic dialogue dataset SoulChatCorpus. The resulting model exhibited robust empathetic capabilities, assisting users in articulating their thoughts and offering suitable suggestions during psychological counseling.

While encoder-decoder LLMs integrate the strengths of encoder-only and decoder-only architectures, balancing text understanding and generation, Wang et al. [89] showed that decoder-only LLMs excel in zero-shot scenarios without fine-tuning. In contrast, encoder-decoder LLMs require multitask fine-tuning with annotated data to reach optimal performance. Since the prevailing LLM training paradigm relies on unsupervised learning on large-scale corpora, decoder-only architectures, with their superior zero-shot performance, are better suited to exploit unlabeled data. As a result, decoder-only architectures remain the predominant choice for LLMs.

3.2. Structure of MLLMs

As shown in Fig. 4, this section provides a detailed discussion of three critical modules in MLLMs: the Vision Encoder, the LLM Backbone, and the Modality Alignment Module. The method of leveraging expert models to construct MLLMs is treated as a form of prompt augmentation method [92] and is discussed alongside other modality alignment modules.

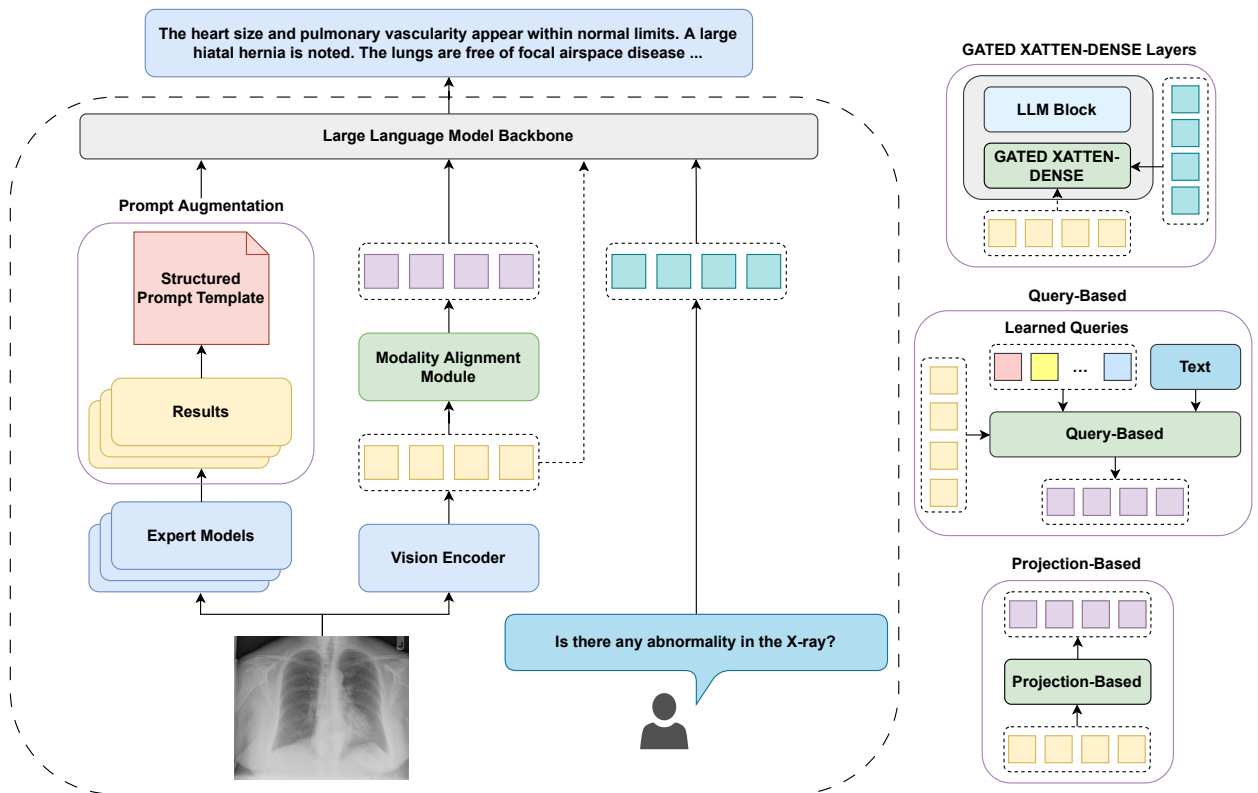


Fig. 4. The core modules and pipeline of MLLMs. On the far right are three types of modality alignment modules. The approach of utilizing expert models to construct MLLMs is regarded as a type of prompt augmentation method, classified under modality alignment modules for further elaboration.

3.2.1. Vision Encoder

MLLMs extend LLMs by incorporating a vision encoder, thereby equipping them with visual capabilities. Specifically, the vision encoder V encodes visual input I_x into visual features Z_x , illustrated below:

$$Z_x = V(I_x) \quad (2)$$

Various options exist for the vision encoder V , such as ResNet [93], the ViT[94], and CLIP-ViT[52]. Here, we provide several insights regarding different vision encoders. (1) Transformer-based vision encoders typically outperform CNN-based counterparts due to their superior scalability. Pre-training on large-scale datasets enables these encoders to extract more effective visual representations. Furthermore, as most LLMs are based on the Transformer architecture, Transformer-based vision encoders integrate more seamlessly with these models. (2) Chen et al. [95] observed that contrastive-learning-based vision encoders often outperform classification-based ones, particularly in localization and visual-text understanding tasks. This advantage likely arises from the pre-training tasks of contrastive-learning-based encoders, which inherently involve vision-language alignment, enabling the extraction of visual features more aligned with semantic spaces for MLLMs. (3) In the medical domain, vision encoders pre-trained on medical imaging datasets outperform those trained on natural scene datasets. Models such as LLaVA-Med [18], MAIRA-1 [96], and PathChat [97] provide evidence supporting this claim.

In summary, while ResNet remains a strong choice among CNNs, Transformer-based ViT models are increasingly preferred by researchers. Additionally, contrastive-learning-based ViT models, such as CLIP-ViT and EVA-CLIP ViT [98], often outperform classification-pretrained ViT models when used as vision encoders for MLLMs. Consequently, ViT models trained using contrastive learning represent the mainstream choice for vision encoders. Furthermore, vision encoders pre-trained on medical imaging datasets significantly enhance the performance of medical MLLMs.

3.2.2. LLM Backbone

As the cognitive engine of MLLMs, the LLM backbone is the most critical module among the three primary components and contains the largest number of parameters. It equips MLLMs with capabilities such as text interaction, ICL, and reasoning. The operational principle of the LLM backbone in MLLMs is illustrated below:

$$R = L(H_x, T_x) \quad (3)$$

Where L represents the LLM backbone, R denotes the response output of the LLM, T_x indicates the embedded tokens of the text input, and H_x are visual representations that LLM can understand. The specific meaning of H_x is explained in Equation (4).

Although powerful LLMs like ChatGPT have not yet been open-sourced, numerous high-quality open-source LLMs are available for researchers. Among these, the LLaMA series[6, 7, 8] developed by Meta, stands out as one of the most popular open-source LLMs and is frequently used as the backbone for MLLMs. Additionally, fine-tuned versions of LLaMA, such as Vicuna-13B [99], achieve performance comparable to 90% of ChatGPT and Bard. Notably, different models demonstrate varying levels of performance across languages. For example, Mistral [100] excels in French, Qwen [101] is optimized for Chinese, and GPT-4 offers robust support for a broader range of languages. Consequently, researchers can choose LLM backbones based on specific linguistic requirements.

3.2.3. Modality Alignment

Although integrating a vision encoder into LLMs enables them to process visual inputs, LLMs trained exclusively on text datasets cannot interpret the output features Z_x produced by the vision encoder. Consequently, modality alignment is required to transform Z_x into a format that LLMs can understand, as shown in Equation (4):

$$H_x = f(Z_x) \quad (4)$$

Where f represents the modality alignment method, and H_x refers to visual representations that LLMs can understand. Modality alignment plays a critical role in enabling MLLMs to interpret visual information and significantly enhances their multimodal capabilities. The subsequent sections introduce four established modality alignment methods: GATED XATTN-DENSE Layers, Query-Based method, Projection-Based method, and Prompt Augmentation.

Table 2
Detailed information on existing medical MLLMs.

Modality Alignment Method	Model Name	Vision Encoder	LLM Backbone	Data Source	Evaluation Method	Date
GATED XATTN-DENSE Layers	Med-Flamingo [102]	CLIP-ViT	LLaMA	MTB, PMC-OA	AEM, Human	2023/07
Query-Based	MedBLIP [103]	EVA-CLIP-ViT	BioMedLM	ADNI, NACC, OASIS	AEM	2023/05
	XrayGLM [104]	ViT-G	ChatGLM	MIMIC-CXR, OpenI	/	2023/05
	PCLMed [105]	EVA-CLIP-ViT	ChatGLM	ImageCLEF 2023 caption prediction	AEM	2023/06
	RadFM [106]	3D ViT	MedLLaMA-13B	MedMD, RadMD	AEM, Human	2023/08
	CheXagent [107]	EVA-CLIP-ViT	Mistral	CheXinstruct	AEM, Human	2024/01
Projection-Based	CLIP-ViT w/ GPT2 [108]	CLIP-ViT	GPT2-XL	Slake, PathVQA, OVQA	AEM	2023/05
	MedVInt [109]	PMC-CLIP-ViT	PMC-LLaMA	PMC-VQA	AEM	2023/05
	PathAsst [110]	PathCLIP-ViT	Vicuna	PathCap, PathInstruct	/	2023/05
	LLaVA-Med [18]	CLIP-ViT	LLaMA	PMC-15M, VQA-RAD, SLAKE, PathVQA	AEM, AI	2023/06
	XrayGPT [111]	MedCLIP-ViT	Vicuna	MIMIC-CXR, OpenI	AEM, AI	2023/06
	Med-PaLM M [21]	ViT-e, ViT-22B	PaLM	MultiMedBench	AEM, Human	2023/07
	R2GenGPT [20]	Swin Transformer	LLaMA 2	IU-Xray, MIMIC-CXR	AEM	2023/09
	Qilin-Med-VL [112]	ViT	LLaMA-2-Chinese	ChiMed-VL	/	2023/10
	MAIRA-1 [96]	RAD-DINO	Vicuna	MIMIC-CXR	AEM, Human	2023/11
	PeFoM-Med [113]	EVA-CLIP-ViT	LLaMA 2	ROCO, VQA-RAD	AEM, Human	2024/01
	M3D-LaMed [114]	3D ViT	LLaMA-2	M3D-Data	AEM, AI	2024/03
	MoE-TinyMed [115]	CLIP-ViT	Phi-2	LLaVA-Med, VQA-RAD, SLAKE, PathVQA	AEM	2024/04
	MAIRA-2 [116]	Rad-DINO	Vicuna	MIMIC-CXR, PadChest, USMix	AEM	2024/06
	PathChat [97]	UNI	LLaMA 2	PubMed, WSIs	AEM, Human	2024/06
	HuatuoGPT-Vision [117]	CLIP-ViT	Yi-1.5	PubMedVision, HuatuoGPT-II	AEM	2024/06
	miniGPT-Med [118]	EVA-CLIP-ViT	LLaMA 2	MIMIC, NLST, SLAKE, RSNA, RadVQA	AEM	2024/07
	SkinGPT-4 [119]	ViT	LLaMA 2	SKINCON, Dermnet	Human	2024/07
	LLaVA-Med++ [56]	CLIP-ViT	LLaMA	MedTrinity-25M, VQA-RAD, SLAKE, PathVQA	AEM	2024/08
	SigPhi-Med [120]	SigLIP	Phi-2	LLaVA-Med, VQA-RAD, SLAKE, PathVQA	AEM	2024/10
	Prompt Augmentation	ChatCAD [17]	Expert models	ChatGPT	MIMIC-CXR, CheXpert	AEM
Visual-Med-Alpaca [92]		Expert models	Med-Alpaca	ROCO, BigBio	/	2023/04
ChatCAD+ [121]		Expert models	ChatGPT	CheXpert, MIMIC-CXR	AEM	2023/05
OphGLM [122]		Expert models	ChatGLM	Web data, MedDialog	AEM	2023/06

GATED XATTN-DENSE Layers, introduced in Flamingo [54], incorporate dense cross-attention layers into a frozen pre-trained LLM. These cross-attention layers receive information from the vision encoder’s output, which is typically processed through a Perceiver Resampler [123] to reduce the computational complexity of vision-text cross-attention. Using additional cross-attention layers, the LLM generates text responses based on visual representations. Subsequent works such as Med-Flamingo [102], which are based on Flamingo, also utilize these cross-attention layers for modality alignment.

Query-Based method, often considered a multimodal perceiver [124], extracts information from visual representations through a set of learnable query vectors. For example, the Q-Former introduced in BLIP-2 [125] extracts visual features from a frozen vision encoder, enabling LLMs to generate text responses aligned with visual information. This query-based approach can be effectively extended to 3D spaces, as demonstrated by Chen et al. [103] with MedBLIP, which adapts the querying mechanism for 3D medical imaging. Although these methods can represent images with only a few queries, thus reducing training costs, they risk losing critical visual information. Moreover, Yao et al. [126] demonstrated that the Q-Former functions merely as an inefficient visual token compressor. For cost-effective token reduction, adaptive average pooling outperforms it.

Projection-Based method can be regarded as a type of multimodal converter [124]. It is simpler than the query-based method, as it maps visual representations from the vision encoder’s output to the word embedding space using a simple projection layer, allowing LLMs to interpret images. For example, LLaVA-Med, Qilin-Med-VL, and XrayGPT [111] employ a simple linear layer to map visual representations, and MedVIntTE [109] and LLaVA-1.5 [127] rely on MLPs for this purpose. The mapped visual representations, combined with textual representations, act as inputs to the LLM backbone. Compared to query-based methods, projection-based methods preserve more visual information because they do not reduce the number of vision tokens, though this advantage comes at the cost of increased training complexity.

Prompt Augmentation typically processes images using expert models, integrates their results with text prompt templates to serve as input for LLMs, thereby linking visual information with text. For example, OphGLM [122] extracts information from fundus images through classification and segmentation models, embeds this information into structured text templates, and forms diagnostic reports that act as inputs to LLMs. Similarly, in ChatCAD [17], X-ray images are processed by trained computer-aided diagnosis (CAD) models to generate outputs, which are subsequently transformed into natural language via prompt templates and used as inputs to LLMs. Compared to query-based and projection-based methods, prompt augmentation leverages expert models without requiring additional modality alignment training. However, this method heavily depends on the quality of prompt templates and the performance of expert models.

Despite their differences, all four approaches share a text-centered design philosophy, leveraging text as a modality space to convert visual information into textual space, thereby enabling LLMs to interpret visual input.

4. Principles of Medical LLMs and MLLMs

To assist researchers and medical professionals in developing their own medical LLMs and MLLMs, this section summarizes the medical datasets available for training, explains the methods for fine-tuning medical LLMs and MLLMs, and finally discusses three approaches for evaluating the performance of medical LLMs and MLLMs.

4.1. Training Datasets

The currently available datasets are categorized into six main types: electronic health records (EHRs), scientific literature, QA, dialogue, medical image-text pairs, and instruction-following data. Table 3 provides detailed information about these datasets.

Electronic Health Records: EHRs include personal health records, such as demographic information, summaries of major diseases and health issues, and primary healthcare records. The Medical Information Mart for Intensive Care III (MIMIC-III) [128] is one of the largest and most widely used publicly available EHR datasets, containing approximately 2 million de-identified notes across 13 specialties, including cardiology, respiratory medicine, and radiology. The MIMIC-III dataset is a valuable resource for developing medical LLMs, as evidenced by models such as AMIE [152] and GatorTron [153], which were trained using this dataset. Besides MIMIC-III, other widely used EHR datasets include the Clinical Practice Research Datalink (CPRD) [130] and the updated version of MIMIC-III, MIMIC-IV [129].

Scientific Literature: Scientific literature, which provides accurate and authoritative medical knowledge, serves as a key source for medical datasets. PubMed, the most widely used repository for biomedical and life science literature, provides access to key resources, including MEDLINE, PubMed Central (PMC), and the NCBI Bookshelf. It indexes citations from over 34 million biomedical research articles. PubMed abstracts include approximately 4.5 billion words, making it an excellent resource for medical training datasets. In addition to PubMed, PubMed Central (PMC) is a widely used repository offering free full-text access, with its articles collectively containing around 13.5 billion words. PubMed and PMC offer high-quality medical literature, often used as sources for other datasets. For instance, PMC-OA [150], PMC-VAQ [109], and PMC-15M [151] are three biomedical multimodal datasets extracted from PMC, significantly facilitating the development of medical LLMs [65, 64, 77] and MLLMs [18, 102].

Question-Answer: QA datasets are categorized into two types: discriminative QA [14, 133] and generative QA [109]. Discriminative QA datasets primarily include multiple-choice questions, whereas generative QA focuses on open-ended questions. Typical QA datasets include PubMedQA [132], MedQA [14], PMC-VQA [109], and MultiMedQA [60]. MultiMedQA, in particular, is a comprehensive medical QA dataset encompassing seven sub-datasets that assess the authenticity, helpfulness, accuracy, and potential harm of LLMs' responses. Beyond text-based QA datasets, the medical domain also includes VQA datasets. For example, the classic VQA-RAD [141] and SLAKE [142] are prominent VQA datasets comprising radiology images, questions, and answers, which span various organs and anatomical regions. Such medical VQA datasets have greatly contributed to the advancement of medical MLLMs.

Dialogue: While datasets like EHRs, scientific literature, and QA pairs enrich LLMs and MLLMs with medical knowledge, relying exclusively on these may result in insufficient long-dialogue interaction capabilities, limiting their practical clinical application. As a result, researchers are focused on developing high-quality dialogue datasets to enhance the models' capabilities in multi-turn conversations. For example, Li et al. [16] curated HealthCareMagic-100k by collecting around 100K authentic doctor-patient dialogues from the HealthCareMagic platform, followed by extensive filtering. To bypass the labor-intensive process of gathering authentic dialogue datasets—requiring extensive filtering and deduplication—Li et al. [16] simulated real dialogue scenarios using ChatGPT to create a synthetic dataset named GenMedGPT-5k.

Image-Text Pairs: In addition to the aforementioned VQA datasets, image-text pairs, including image-caption and image-report datasets, are essential for training medical MLLMs. These datasets offer crucial context and annotations, enriching the ability of MLLMs to interpret visual information in medical scenarios. For example, PMC-OA [150] comprises 1.65 million medical image-text pairs sourced from PMC, and it has been utilized to train models like PMC-CLIP [150] and Med-Flamingo [102]. Zhang et al. [109] built upon PMC-OA and utilized ChatGPT to generate

Table 3
Summary of medical datasets for pre-training and fine-tuning.

Datasets	Type	Description	AI Synthesis
<i>Datasets For Medical LLMs</i>			
MIMIC-III [128]	EHR	Approximately 2M de-identified notes.	✗
MIMIC-IV [129]	EHR	About 300K patients, 430K admissions.	✗
CPRD [130]	EHR	Anonymized medical records for over 11.3M patients.	✗
PubMed	Literature	Over 34M citations and abstracts of biomedical literature, about 4.5B words.	✗
PMC	Literature	Provides free full-text access to PubMed, about 13.5B words.	✗
CORD-19 [131]	Literature	More than 140K papers, with more than 72K full text.	✗
PubMedQA [132]	QA	1K labeled, 612K unlabeled and 211.3K manually generated QA.	✗
MedQA (USMLE) [14]	QA	61,097 multiple-choice QA pairs.	✗
MedMCQA [133]	QA	194K multiple-choice QA pairs.	✗
cMedQA2 [134]	QA	100K questions and 200k answers.	✗
MultiMedQA [60]	QA	Includes six existing datasets and one new dataset.	✗
MedQuAD [135]	QA	47,457 question-answer pairs from trusted medical sources.	✗
Medical Meadow [63]	QA	Over 160K QA pairs.	✓
Huatuo-26M [136]	QA	26M QA pairs.	✗
Psych8k [23]	QA	8,187 query-answer pairs.	✓
CMtMedQA [69]	Dialogue	70K multi-round conversation datasets from real doctor-patient conversations.	✓
MedDialog [137]	Dialogue	3.4M Chinese conversations and 0.6 million English conversations.	✗
HealthCareMagic-100k [16]	Dialogue	100K authentic patient-doctor conversations.	✗
GenMedGPT-5k [16]	Dialogue	5K generated conversations between patients and physicians from ChatGPT.	✓
MedC-1 [64]	Instruction-Following Data	202M tokens.	✓
MedInstruct-52k [74]	Instruction-Following Data	52K instruction-response pairs generated by GPT-4.	✓
UMLS [138]	Knowledge Base	2M entities for 900K concepts.	✗
CMeKG [139]	Knowledge Base	Chinese medical knowledge graph.	✗
COMETA [140]	Web Data	Consisting of 20K English biomedical entity mentions.	✗
TCM-Corpus-1B [75]	Web Data	20GB dataset collected from Baidu Baike, Wikipedia and other sources.	✗
ChiMed [73]	Hybrid	Composed of various data such as QA, books, dialogues, etc.	✗
GAP-REPLAY [77]	Hybrid	Includes data from clinical practice guidelines, abstracts, and original articles.	✗
<i>Datasets For Medical MLLMs</i>			
PMC-VQA [109]	QA	Contains 149K images, 227K VQA pairs.	✓
VQA-RAD [141]	QA	315 radiology images and 3515 QA pairs generated by clinicians.	✗
Slake [142]	QA	642 radiology images and over 7000 diverse QA pairs.	✗
PathVQA [143]	QA	4,998 pathology images with 32,799 QA pairs.	✗
MIMIC-CXR [144]	Image-Report	227,835 imaging studies for 65,379 patients.	✗
OpenI [145]	Image-Report	7,470 images and 3,955 reports.	✗
CheXpert [146]	Image-Report	224,316 chest X-rays with reports.	✗
ROCO [147]	Image-Caption	Contains more than 81K radiologic images, each with a corresponding title, keywords.	✗
OpenPath [148]	Image-Caption	208,414 pathology images paired with natural language descriptions.	✗
MedCaT [149]	Image-Caption	160K images with captions and inline references.	✗
PathCap [110]	Image-Caption	142K high quality pathology image-caption pairs.	✓
MedMD [106]	Image-Caption	15.5M 2D scans, 180k 3D scans, with corresponding captions or diagnosis labels.	✓
PMC-OA [150]	Image-Caption	1.6M image-caption pairs.	✗
PMC-15M [151]	Image-Caption	15M figure-caption pairs from over 3M articles.	✗
LLaVA-Med-Alignment [18]	Image-Caption	600K image-caption pairs from PMC-15M.	✓
ChiMed-VL-Alignment [112]	Image-Caption	580,014 images and context information or descriptions.	✓
PubMedVision-Alignment [117]	Image-Caption	647,031 image-caption pairs.	✓
MedTrinity-25M [56]	Image-Annotation	Contains 25M samples along with their multigranular annotations.	✓
LLaVA-Med-Instruct [18]	Instruction-Following Data	60K instruction-following data.	✓
ChiMed-VL-Instruction [112]	Instruction-Following Data	469,441 question-answer pairs.	✓
PubMedVision-Alignment [117]	Instruction-Following Data	647,031 instruction-following data.	✓
PathInstruct [110]	Instruction-Following Data	180K instruction-following data.	✓
CheXInstruct [107]	Instruction-Following Data	An instruction-tuning dataset curated from 28 publicly available datasets	✓
ApolloCorpora [80]	Hybrid	2.5B tokens with data in multiple languages.	✗
PedCorpus [82]	Hybrid	The corpus includes pediatric textbooks, clinical guidelines, and knowledge graphs.	✓
M3D-Data [114]	Hybrid	Comprising 120K image-text pairs and 662K instruction-response pairs.	✓

¹ "Hybrid" means that the dataset is a mixture of multiple types of data.

² "AI Synthesis" indicates that generative AI such as chatGPT and GPT-4 were used during the development of the dataset to assist in generating the data.

a diverse set of high-quality QA pairs. After filtering, they created PMC-VQA, which provides 227K VQA pairs. PMC-15M [151], derived from PMC articles, includes 15 million figure-caption pairs, making it two orders of magnitude larger than MIMIC-CXR [144]. Other notable medical image-text datasets, such as ChiMed-VL [112], RadMD [106], and Open-I [145], have also played a significant role in advancing medical MLLMs.

Instruction-Following Data: The effectiveness of medical LLMs and MLLMs in downstream tasks depends not only on the medical knowledge they acquire but also on their capacity to follow user instructions. Enhanced instruction-following ability enable more accurate comprehension and execution of user directives, thereby improving downstream task performance. To strengthen the instruction-following ability of medical LLMs and MLLMs,

training on instruction-following datasets has become widely accepted. These datasets often include instruction-response pairs or image-instruction-response, such as PathInstruct [110] with 180K image-instruction-response pairs, LLaVA-Med-Instruct [18] with 60K image-instruction-response pairs, and ChiMed-VL-Instruction [112] with 469K image-instruction-response pairs. Such datasets effectively enhance models' instruction-following ability, leading to improved zero-shot performance. Notably, instructions can include explicit commands or questions; as a result, some works [18, 112] also considered QA data as instruction-following data.

In addition to the aforementioned data types, various medical knowledge bases [138] and web data [75] also serve as training sources for medical LLMs and MLLMs. Different types of data are typically utilized at various training stages. For medical LLMs, scientific literatures and web data are primarily used during pre-training and continual pre-training phases to incorporate medical knowledge and facilitate domain adaptation. In contrast, QA, dialogue, and instruction-following data are typically employed during the fine-tuning phase to enhance interaction capabilities and instruction-following performance. For medical MLLMs, image-caption pairs are widely used during pre-training to align visual features with text representations, while instruction-following data is commonly applied for fine-tuning. Furthermore, the extensive content of EHRs and scientific literature often makes them foundational sources for other data types. For example, PMC-15M [151] is derived from 3 million PMC articles, yielding 15 million image-caption pairs. Finally, studies have demonstrated that fine-tuning models with substantial high-quality synthetic data generated by ChatGPT can significantly enhance downstream task performance [154]. As a result, AI-assisted data generation has emerged as a prevalent strategy in the data-scarce medical field. For example, Liu et al. [23] transformed 260 real psychological counseling audio recordings into text and utilized GPT-4 to extract question-answer pairs and generate key summaries, providing supplementary contextual information for constructing the Psych8k dataset.

4.2. Fine-Tuning Methods

The extensive parameters in LLMs and MLLMs make training medical LLMs and MLLMs from scratch computationally intensive. Consequently, the prevalent method for constructing medical LLMs and MLLMs involves fine-tuning general foundation models using medical datasets. This section outlines six fine-tuning methods, as illustrated in Fig. 5, to aid researchers in developing medical LLMs and MLLMs. In addition, this section explores the characteristics of these fine-tuning methods, providing practical guidance for selecting the appropriate fine-tuning strategies.

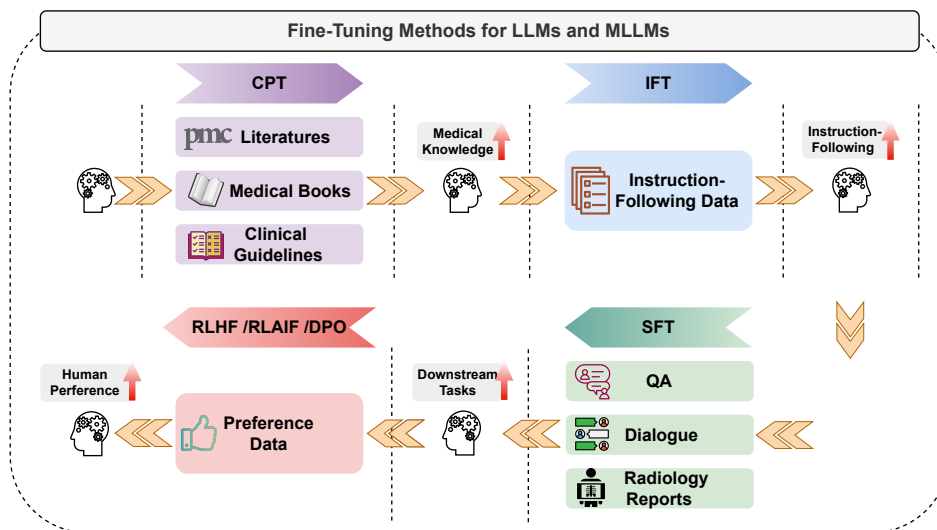


Fig. 5. Overview of six fine-tuning methods. In our analysis of the related work on medical LLMs and MLLMs, we found that Continuous Pre-Training (CPT) is commonly used to inject medical knowledge into LLMs and MLLMs; Instruction Fine-Tuning (IFT) enhances the models' ability to follow instructions and their zero-shot performance; Supervised Fine-Tuning is frequently employed to improve model performance on specific tasks; and Reinforcement Learning from Human Feedback (RLHF), Reinforcement Learning from AI Feedback (RLAIF), and Direct Preference Optimization (DPO) are used to align model behavior with human preferences.

Continual Pre-Training: CPT [34, 155] refers to the further pre-training of a general pre-trained model on large-scale medical datasets. Given that general LLMs and MLLMs often lack domain-specific medical knowledge, CPT aims to incorporate such knowledge into these models. For LLMs, datasets comprising extensive medical texts, including scientific literature, books, and clinical guidelines, are typically used in CPT to ensure LLMs acquire sufficient medical knowledge. For example, MEDITRON-70B [77], which is based on LLaMA 2, utilizes a medical mixed dataset composed of clinical guidelines, PubMed papers, and abstracts for CPT. For MLLMs, medical image-caption pairs are frequently utilized as CPT datasets. This method aligns visual features with the semantic space by training MLLMs to predict captions for medical images. Prominent models, including LLaVA-Med [18], Qilin-Med-VL [112], and HuatuoGPT-Vision [117], utilize large-scale medical image-caption datasets for CPT.

Instruction Fine-Tuning: While CPT using large-scale biomedical corpora integrates medical knowledge into LLMs and MLLMs, clinical performance also relies on their instruction-following capability. A lack of instruction-following ability in LLMs and MLLMs can result in unpredictable behavior, underscoring the necessity of fine-tuning these models with instruction-following datasets to improve their compliance with diverse human instructions [18, 156]. Instruction fine-tuning (IFT) enables models to accurately understand and execute human directives, thereby significantly enhancing their zero-shot performance. For example, Singhal et al. [60] demonstrated substantial improvements in Flan-PaLM's performance on MedQA, MedMCQA, and PubMedQA after fine-tuning it with instruction data.

Supervised Fine-Tuning: To further improve the performance of LLMs and MLLMs in downstream tasks, SFT is typically performed on datasets tailored to those tasks. For instance, Chen et al. [77] fine-tuned MEDITRON on MedQA, PubMedQA, and MedMCQA datasets to enhance its performance in medical QA tasks. Similarly, Li et al. [18] improved LLaVA-Med's performance in medical VQA tasks by fine-tuning it on PathVQA, SLAKE, and VQA-RAD datasets. Hyland et al. [96] leveraged large-scale image-report pairs to enhance MAIRA-1's performance in radiology report generation tasks. While prior studies often do not strictly distinguish between SFT and IFT, we argue that IFT focuses on utilizing instruction datasets to enhance models' instruction-following ability and zero-shot performance. Conversely, SFT targets high-quality task-specific datasets to fine-tune LLMs and MLLMs, thereby boosting their performance in specialized downstream tasks.

Reinforcement Learning from Human Feedback: Reinforcement learning from human feedback (RLHF) [5, 157] is a technique for aligning model behavior more closely with human preferences and directives. Compared to the preceding three fine-tuning approaches, RLHF is more intricate and involves three distinct stages [5, 157, 158]: collecting human feedback, training the reward model, and policy optimization, as shown in Fig. 6. During the human feedback collection stage, the primary task is gathering comparison data. Typically, an LLM is provided with a prompt, generating multiple outputs that experts annotate and score based on quality [158]. These annotated outputs, along with the prompts, form the comparison data. For example, Yang et al. [69] employed 6 medical graduate students or clinical doctors as labelers to rank the model's outputs based on dimensions such as safety, professionalism, and fluency, forming a comparison dataset. During reward model training, a reward model learns from the comparative data to produce scalar rewards that represent human preferences. During policy optimization, a new prompt is provided as input to the LLM, whose response is evaluated by the reward model, outputting a scalar reward. The LLM is then fine-tuned using Proximal Policy Optimization (PPO) based on these rewards. Notably, the reward model's data quality is typically lower than that used for SFT [33], directly transitioning from pre-training to RLHF may lead to suboptimal fine-tuning results, so RLHF is often conducted after IFT and SFT [7, 157].

Reinforcement Learning from AI Feedback: Reinforcement learning from AI feedback (RLAIF) is a cost-effective alternative to RLHF, where the reward model learns from AI feedback without requiring human annotation [159]. In the medical domain, Zhang et al. [66] sampled multiple responses from the fine-tuned model after IFT and SFT and used ChatGPT to evaluate them across dimensions such as informativeness, coherence, adherence to human preferences, and accuracy. This comparison data was then used to train a reward model. Training reward models through AI feedback eliminates the need for manual data labeling in RLHF, significantly reducing labor costs.

Direct Preference Optimization: While RLHF and RLAIF align models with human preferences and ethical norms, they typically involve fitting a reward model that reflects human preferences and combining reinforcement learning to fine-tune LLMs and MLLMs. However, this process is both complex and often unstable. Direct preference optimization (DPO) [160] is a simpler and more efficient paradigm for aligning models with human preferences, bypassing the need for a reward model by directly optimizing the model using preference data. The core idea of DPO is to use an analytical mapping from the reward function to the optimal policy, transforming the reward function loss into a policy loss and thereby eliminating the need for explicit reward modeling. For example, Qilin-Med [73] uses

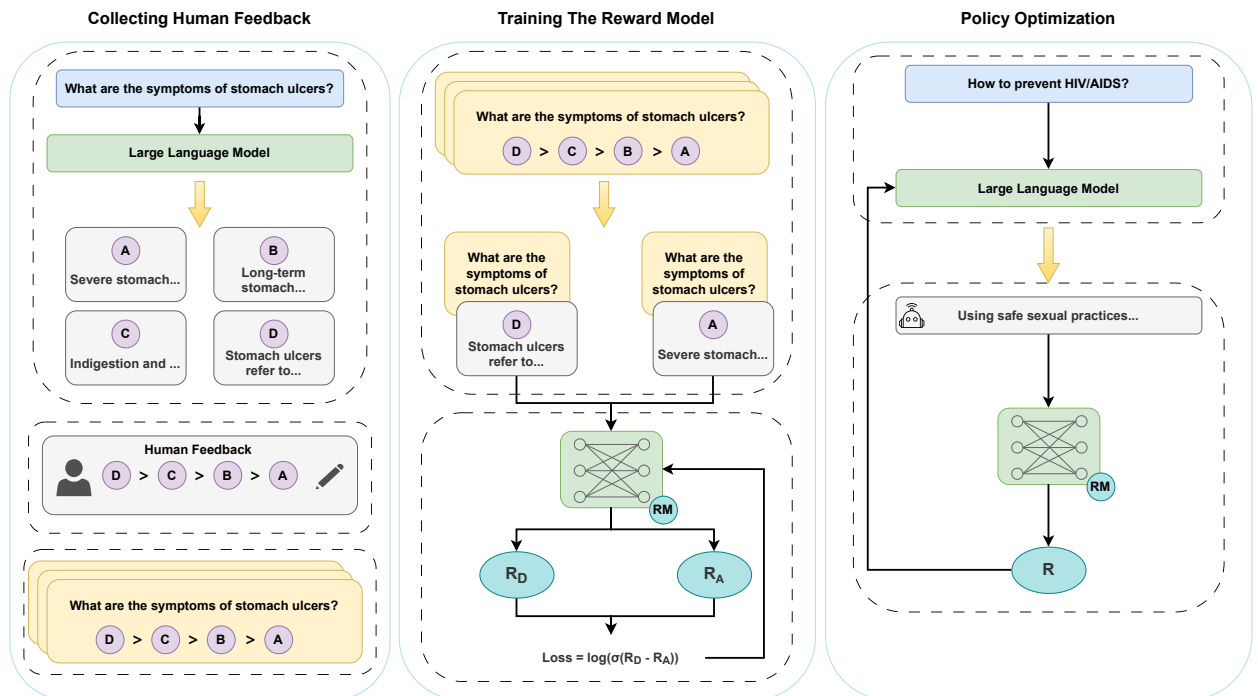


Fig. 6. Pipeline of Reinforcement Learning from Human Feedback. **Left** illustrates the Collect Human Feedback phase: A labeler provides a single prompt to the model, ranks multiple responses, and collect the prompt along with the labeled responses. **Mid** shows the Reward Model Training phase: A prompt and two responses are randomly sampled from the dataset and used to train the reward model. **Right** represents the Policy Optimization phase: A new prompt is provided, and the reward model generates a scalar reward for the model's response, which is then used for policy optimization.

two publicly available preference datasets after SFT to optimize the model through DPO, achieving stable and efficient training while aligning it with human preferences.

In summary, CPT injects medical knowledge into LLMs and MLLMs. IFT enhances instruction-following capabilities and zero-shot performance. SFT improves task-specific performance in downstream tasks. and RLHF, RLAIF, and DPO further align models with human preferences. In the current work on medical LLMs and MLLMs, CPT, STF, and IFT are the most commonly used fine-tuning methods. In contrast, RLHF, RLAIF, and DPO have been used relatively infrequently in medical LLMs and MLLMs, with no existing medical MLLMs reported to utilize RLHF, RLAIF or DPO.

4.3. Evaluation Methods

With the continuous emergence of abilities in medical LLMs and MLLMs, coupled with growing ethical and safety concerns [161], comprehensively evaluating their performance in medical tasks and ensuring their safety in clinical environments has become a pressing challenge. Therefore, this subsection summarizes three commonly used evaluation methods and examines their respective advantages and limitations, as shown in Fig. 7.

4.3.1. Automatic Evaluation Metrics

Accuracy is commonly used to evaluate the performance of medical LLMs and MLLMs on multiple-choice question benchmarks, such as MedQA [14] and MedMCQA [133]. However, accuracy alone is insufficient for evaluating tasks requiring longer text generation, such as medical report writing and summarization. Consequently, the evaluation of medical LLMs and MLLMs must incorporate additional metrics for a more comprehensive assessment.

Bilingual Evaluation Understudy (BLEU) [162] metric assesses the quality of generated text by measuring the similarity of n -grams (consecutive word sequences of length n) between the generated and reference texts. BLEU is categorized into BLEU-1, BLEU-2, BLEU-3, and BLEU-4 based on the value of n , capturing n -gram similarity at different levels. For instance, BLEU-1 reflects word-level accuracy, while BLEU-4 emphasizes text continuity.

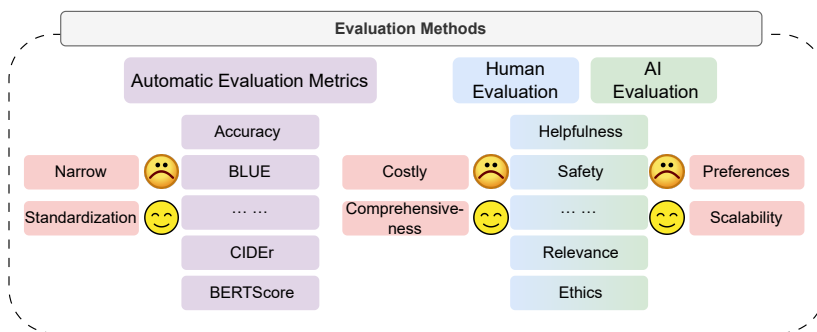


Fig. 7. Overview of three evaluation methods. We summarize three methods for evaluating medical LLMs and MLLMs: Automatic Evaluation Metrics, Human Evaluation, and AI Evaluation, and discuss their respective advantages and disadvantages.

Recall-Oriented Understudy for Gisting Evaluation (ROUGE) [163] includes metrics such as ROUGE-N, ROUGE-L, ROUGE-W, and ROUGE-S. Similar to BLEU, ROUGE-N evaluates n-gram similarity between generated and reference texts; however, it emphasizes recall, whereas BLEU focuses on precision. ROUGE-L assesses textual coherence by calculating the longest common subsequence between the generated and reference texts. ROUGE-W extends ROUGE-L by applying weighted scoring, giving higher importance to continuous and accurate matching subsequences. ROUGE-S extends ROUGE-N by accommodating non-contiguous words in n-grams. Google BLEU (GLEU) [164], a BLEU variant, incorporates lexical overlap and word order, offering a more nuanced evaluation of fluency and naturalness. The Distinct-n [165] metric quantifies text diversity by calculating the ratio of unique n-grams to total n-grams. CIDEr [166], tailored for image caption evaluation, combines n-gram recall and precision, assigning higher weights to rare n-grams to assess key information. BERTScore [167] leverages pre-trained BERT embeddings to compute token-level similarity scores between generated and reference sentences. Compared to n-gram-based metrics, BERTScore more effectively evaluates vocabulary and compositional diversity.

In the medical domain, models such as HuaTuoGPT[66], ClinicalGPT [68] [68], SoulChat [22] and BianQue [84] commonly use these metrics to assess generative performance. While these metrics partially capture the accuracy and fluency of generated text, they fall short in assessing clinical dialogue quality [152] and alignment with human values, necessitating the inclusion of human evaluation.

4.3.2. Human Evaluation

Human evaluation is an essential approach to assess medical LLMs and MLLMs, as it addresses aspects that automatic evaluation metrics may fail to capture. For example, Tu et al. [152] highlighted that metrics such as BLEU and ROUGE fail to reflect the clinical quality of medical consultations. To address this, 23 medical experts from the United States, the United Kingdom, and India were invited to evaluate model-generated responses based on accuracy, appropriateness, and comprehensiveness. Similarly, Yang et al. [69] engaged human experts to assess the safety, accuracy, and ethical implications of model-generated responses.

Clearly, human evaluation can encompass critical aspects such as safety and helpfulness, which are essential for medical LLMs and MLLMs. Despite its ability to evaluate diverse capabilities of medical LLMs and MLLMs, human evaluation remains inherently subjective due to the absence of standardized criteria. Moreover, hiring medical experts is costly, making AI evaluation a practical alternative.

4.3.3. AI Evaluation

Employing advanced AI models, such as ChatGPT and GPT-4, which align with human values, has become the predominant method for evaluating medical LLMs and MLLMs [168]. Wang et al. [168] conducted experiments on five natural language generation datasets, showing that ChatGPT, as an evaluation tool, outperformed traditional metrics in most cases and matched human evaluation. In the medical domain, Li et al. [18] asked GPT-4 to evaluate responses from itself and LLava-Med on criteria such as helpfulness, relevance, accuracy, and level of detail. Liu et al. [23] instructed GPT-4 to evaluate whether LLM responses were acceptable and whether their tone resembled that of human counselors.

Despite its scalability and reduced reliance on human input, AI evaluation has notable limitations. Studies [169, 61] have revealed that GPT-4, as an evaluation tool, tends to favor the first response when multiple answers are presented sequentially. Furthermore, GPT-4 tends to prefer longer responses and those it has generated itself [23]. To mitigate the limitations of the aforementioned methods, integrating multiple evaluation approaches may provide more reliable results. Additionally, training specialized LLMs or MLLMs through reinforcement learning or other methods to align with human judgment criteria could overcome AI evaluation's shortcomings.

5. Applications of LLMs and MLLMs in Medicine

Traditional medical models are tailored for specific tasks, including medical named entity recognition, relation extraction, text classification, and semantic textual similarity. Although these models perform well in their respective tasks, they are limited in their ability to integrate multi-source data for complex clinical applications. In contrast, medical LLMs and MLLMs excel in navigating diverse medical contexts and performing a wide range of tasks, while also gradually surpassing traditional deep learning models in single-task performance, as shown in Fig. 8. Consequently, medical LLMs and MLLMs act as versatile medical assistants with extensive application potential. To aid practitioners in comprehending the developmental trajectory of LLMs and MLLMs in medicine, this section highlights their potential applications in healthcare (Fig. 9) and outlines strategies for employing these models in diverse medical tasks.

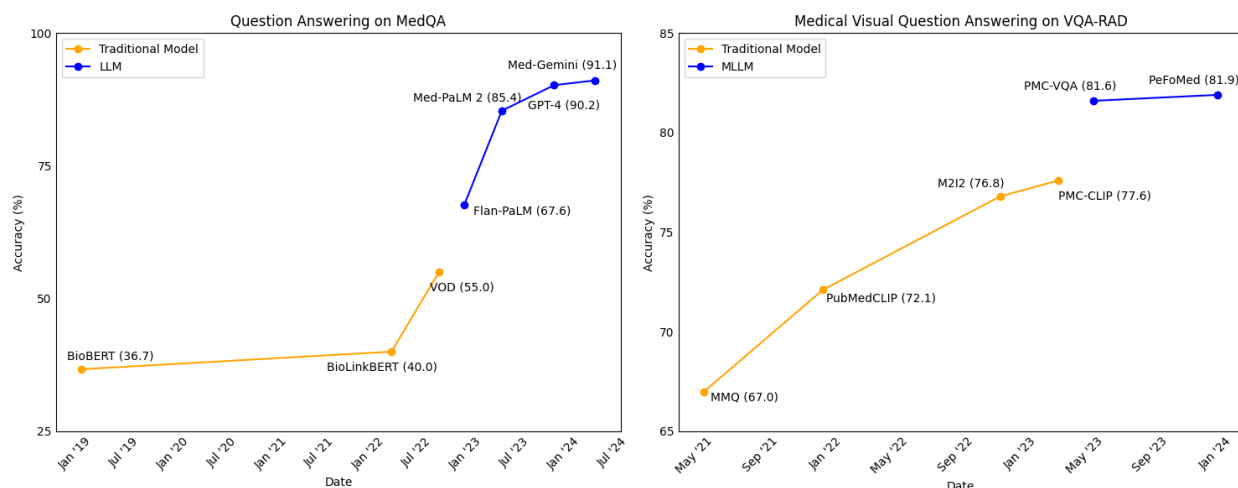


Fig. 8. Performance comparison between traditional deep learning models and medical LLMs / MLLMs on QA and VQA benchmark datasets. Recent advancements in medical LLMs / MLLMs demonstrate their significant superiority over traditional deep learning models on QA and VQA benchmark datasets.

5.1. Medical Diagnosis

AI has been under development in medical diagnosis for several decades [170]. While some breakthroughs have been achieved, its role has largely been confined to assisting with some simple diagnostic tasks, such as medical image segmentation, lesion detection and classification. In recent years, with advancements in LLMs and MLLMs, doctors and patients are now expected to rely on these models for more comprehensive diagnosis. Specifically, these models can process subjective descriptions of disease symptoms or medical images, such as X-rays, and leverage embedded medical knowledge to directly generate diagnostic results. For physicians, medical LLMs and MLLMs shorten diagnostic times and reduces workload. For patients, they offers detailed information about their condition along with recovery suggestions.

Currently, Med-PaLM 2, one of the top-performing medical LLMs, generates responses to consumer medical and adversarial questions, outperforming physician-generated answers across multiple evaluation criteria [13]. This underscores the potential of LLMs as medical diagnostic assistants. Additionally, Yuan et al. [171] showed that multi-turn dialogues with GPT-4 significantly improve its ability to accurately diagnose and recommend effective treatments

for gastrointestinal cancers, achieving a performance level comparable to experienced physicians. In prostate cancer, Zhu et al. [172] developed 22 questions based on patient education guidelines and clinical experience, addressing topics like screening, prevention, treatment options, and postoperative complications. Testing indicated that ChatGPT provided accurate and comprehensive responses, while demonstrating appropriate humanistic care toward patients. Yang et al. [75] further advanced the application of LLMs as diagnostic assistants by training TCM-GPT on traditional Chinese medicine datasets, showing that it outperformed other models in traditional Chinese medicine examination and diagnosis, contributing to the advancement of traditional Chinese medicine. Furthermore, inspired by general MLLMs [10, 11], researchers have developed multimodal medical diagnostic assistants [18, 104, 17, 21, 92, 173, 111], extending diagnostic capabilities from text to medical images. For example, Zhou et al. [119] created SkinGPT-4, a model capable of autonomously analyzing images, identifying skin condition features and categories, conducting in-depth analyses, and offering interactive treatment recommendations.

Medical LLMs and MLLMs, as medical diagnostic assistants, are capable of providing diagnostic recommendations for doctors and offering consultation advice for patients. However, due to inherent limitations of LLMs and MLLMs [29, 161], these models currently serve only as auxiliary tools for physicians, with their diagnostic outputs considered as references rather than definitive results.

5.2. Clinical Report Generation

Clinical reports are standardized documents prepared by doctors for patients. The manual drafting of clinical reports is often tedious and time-consuming, significantly increasing clinicians' workload and reducing overall efficiency. Medical LLMs and MLLMs, equipped with extensive medical knowledge and generative capabilities, serve as efficient tools for clinical report generation.

For example, during medical consultations, doctors typically record key information from patient interactions, which serves as a basis for evaluating conditions or informing other clinical reports. Medical LLMs can function as clinical note-taking tools, automating this task for doctors [65]. Doctors can provide LLMs with doctor-patient interaction records, which the models process to generate detailed medical notes [174]. Doctors can also prompt LLMs to simplify medical notes by removing complex details and generating concise summaries for easier review and analysis [19]. Following a medical diagnosis, doctors often draft diagnostic documents such as radiology reports. Medical MLLMs, capable of processing visual inputs, are particularly effective in assisting with radiology report generation. For example, miniGPT-Med, developed by Alkhaldi et al. [118], achieved state-of-the-art performance in generating medical reports, surpassing prior models in accuracy by 19%. This underscores the feasibility of using medical MLLMs for radiology report generation. During treatment, doctors explain the cause of the disease, the treatment process, and provide detailed clinical information to patients through clinic letters. By utilizing LLMs to generate clinic letters, clinicians can streamline this tedious process, with the resulting letters exhibiting coherence, accuracy, and empathy comparable to those created by humans [175]. After patient recovery, clinicians allocate significant time to drafting discharge summaries, potentially delaying patient discharge. By employing LLMs, clinicians can generate complete discharge summaries in seconds by providing a template and necessary inputs [176]. The quality of these summaries often surpasses those produced by junior doctors [177].

Leveraging advanced LLMs and MLLMs, various clinical reports from patient admission to discharge can be automatically generated. These reports are more comprehensive and accurate than those produced by humans [19, 177], significantly alleviating doctors' workloads and allowing them to dedicate more time to patient care [176]. However, these LLMs and MLLMs are intended to serve exclusively as auxiliary tools for generating clinical reports. They can draft, modify, and summarize reports, but the final versions must be reviewed, edited, and approved by clinicians, who remain accountable for their content [12, 25].

5.3. Medical Education

GPT-4 and Med-PaLM 2 passed the USMLE with scores exceeding 86% [178], while GPT-4V [179] achieved 90.7%, outperforming most medical students on medical image-related questions [180]. This demonstrates that certain LLMs and MLLMs possess the capability to provide educational support for medical students. Moreover, LLMs and MLLMs can role-play in various contexts, enabling the simulation of diverse learning scenarios for users. Consequently, platforms like Khanmigo [181] and Duolingo [182] have integrated tools like GPT-4 to enhance online teaching and learning.

Specifically, medical LLMs and MLLMs can simulate patients in scenarios such as accidents, emergency rooms, or operating rooms, offering simulation training to medical students prior to clinical practice to enhance their professional

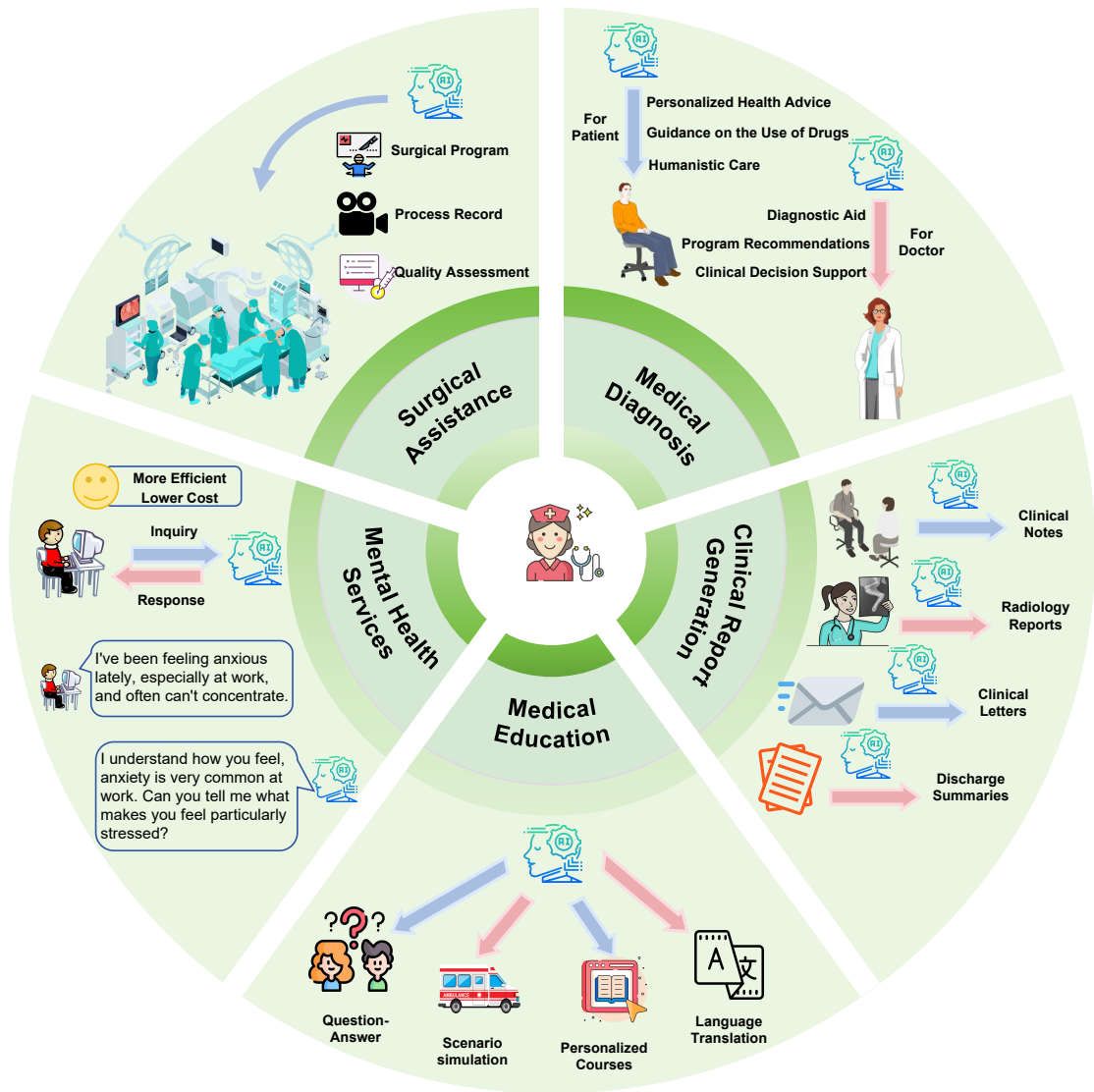


Fig. 9. Overview of potential applications of LLMs and MLLMs in medicine.

and practical skills—tasks unachievable by traditional deep learning [183]. Additionally, medical LLMs and MLLMs can evaluate students' performance in simulated exercises and create personalized learning plans, a process that is typically time-intensive for teachers but can be executed more cost-effectively and efficiently using these models [184]. Finally, given their extensive training on large corpora, medical LLMs and MLLMs excel at translation, offering cross-language capabilities as well as the ability to simplify medical terminology into plain language [185], greatly assisting medical students in reading and writing. In summary, powerful LLMs and MLLMs can enrich medical students' learning experiences by offering comprehensive medical content, personalized curricula, and realistic, diverse scenarios, thereby broadening their horizons in the medical field and laying a strong foundation for clinical practice.

The potential of LLMs and MLLMs in medical education surpasses that of traditional training courses, as educators in these courses cannot always interact with students or offer tailored learning plans. While these models hold great promise in medical education, they should be regarded solely as auxiliary tools and not as replacements for educators, as their inherent biases and hallucinations make it challenging for students to evaluate the accuracy of the generated content [186, 187]. If these models frequently deliver inaccurate content that is hard to identify, they risk misleading

students over time. Thus, LLMs and MLLMs can only serve a supportive role in medical education, requiring students to use them under the guidance and supervision of educators.

5.4. Mental Health Services

Growing societal pressures have led to an increased global demand for mental health services, yet many regions face a severe shortage of mental health specialists due to limited development and resources [188]. Conversation-driven psychological counseling is a central aspect of mental health services, making chatbots powered by LLMs a potential solution for delivering such services in the future [23, 22].

Compared to professional mental health experts, LLM-based mental health chatbots are more accessible and can extend mental health services to remote or underserved areas lacking specialists. Additionally LLM-based chatbots can offer personalized interaction styles tailored to patients' histories and interaction records, including specific emotional patterns, styles, or tones [189]. Furthermore, the high cost of psychological counseling and therapy often deters individuals from seeking mental health services. However, LLM-based chatbots can substantially lower costs [15, 190], making these services more accessible. Finally, studies suggest that individuals are more likely to share negative emotions with chatbots, as certain topics may feel awkward to discuss with humans but easier to disclose to a robot [191]. In terms of convenience, cost, and acceptability, LLM-based mental health chatbots surpass mental health professionals, potentially encouraging more individuals with mental illnesses to seek help [189].

Given that patients with mental illnesses are often more vulnerable and psychologically sensitive, mental health service chatbots must exhibit empathy, trustworthiness, understanding, and comfort in conversations, beyond merely offering advice [22]. While research is advancing the empathetic capabilities of LLMs [22], they still fall short compared to humans. Moreover, despite efforts to align LLMs with human concepts and ethical norms using approaches like SFT and RLHF, they may still produce aggressive or psychologically harmful content [161], which is unacceptable for vulnerable mental health patients. Before deploying LLMs as mental health chatbots in real-world applications, significant work is required to address these challenges, along with stricter control measures for such products.

5.5. Surgical Assistance

Medical robots have undergone rapid development over the past few decades, significantly enhancing surgeons' capabilities [192] and expanding the potential for minimally invasive surgery [193]. Recently, medical robots have entered a new phase with the advent of MLLMs, which not only endow them with visual capabilities but also enhance interactivity and create a more user-friendly environment.

Efforts are currently underway to explore the application of MLLMs in surgical procedures [194]. Integrating MLLMs into surgical robots enables them to perform crucial auxiliary tasks, such as assisting in endoscopic examinations [25]. The robust visual capabilities and specialized knowledge of MLLMs can yield valuable diagnostic conclusions and feasible surgical solutions based on endoscopic images. Furthermore, during surgical procedures, MLLMs can combine video streams to annotate the surgical process, analyze and summarize the steps taken, and record non-compliant operations to assist in the surgeon's post-surgical review and examination.

While medical MLLMs exhibit promising potential for surgical assistance and may play a role in specific medical scenarios, they are not yet suitable for emergency surgeries. This is due to the fact that erroneous information from MLLMs could adversely impact the surgeon's judgment, potentially leading to irreversible consequences. Additionally, current research on MLLMs predominantly focuses on vision-text modalities. We anticipate future investigations will explore other modalities, such as audio and time series, to enable surgical robots to perform more comprehensive and accurate auxiliary tasks while providing more flexible interaction methods.

6. Challenges of LLMs and MLLMs in Medicine

Although LLMs and MLLMs have generated significant interest in the AI community and achieved initial successes in medicine, the unique characteristics of the medical field present numerous challenges and risks for their development and deployment. In this section, we will discuss and analyze the current challenges faced by LLMs and MLLMs in the medical field in detail, as well as propose possible solutions to these challenges.

6.1. Hallucination Phenomenon

Hallucinations refer to the generation of seemingly plausible yet unverified or incorrect information by LLMs and MLLMs [29, 30]. This can result in issues such as misleading radiology reports and the dissemination of incorrect

medical knowledge in medical education [174]. Therefore, addressing the hallucination problem in LLMs and MLLMs is crucial for accelerating the application of these technologies in medicine.

To tackle this challenge, some researchers have proposed a new benchmark dataset specifically for testing hallucinations in medical LLMs and MLLMs [28]. However, these benchmark datasets can only detect hallucination phenomena in models and do not directly mitigate the problem. Other studies indicate that LLMs primarily acquire knowledge during the pre-training phase [38], and the presence of noisy data, such as error messages in the training dataset, can contribute to hallucinations. Thus, the most fundamental approach to reducing hallucinations is to manually or automatically clean unreliable data from the pre-training corpus [30]. However, the pre-training corpus of LLMs and MLLMs typically consists of vast amounts of data, including web-crawled information, which is challenging to clean and necessitates effective selection and filtering strategies. Consequently, fine-tuning LLMs and MLLMs with high-quality medical datasets is an effective strategy for mitigating hallucinations [195, 196]. To further reduce the cost of addressing hallucinations, existing efforts have focused on solutions during the inference stage. For instance, prompting LLMs or MLLMs to verify their own responses has proven effective in alleviating hallucinations [174]. One such method, Chain-of-Verification (CoVe) [197], involves the model drafting an initial response, planning verification questions based on that response, answering these questions to validate the draft, and ultimately generating an optimized answer. Experiments have demonstrated that self-verification methods, such as CoVe, can effectively reduce hallucinations across various tasks. Additionally, retrieval-augmented generation has proven effective in reducing hallucinations [198]. This approach allows the model to retrieve relevant knowledge from external webpages or knowledge bases during the response generation phase [16, 110], significantly addressing the hallucination problem.

Particularly, for MLLMs, the occurrence of hallucinations may be attributed to limited visual perception capabilities. For example, the Q-former uses only 32 learnable vectors to represent an image, which inevitably leads to a loss of visual information. Furthermore, the resolution of vision encoders in current medical MLLMs typically ranges from 224 to 336, hindering their ability to capture the complexities of biomedical images and contributing to hallucinations. To tackle this challenge, integrating more advanced and higher-resolution vision encoders, such as DINOv2 [199] and SAM [200], into medical MLLMs represents a viable solution [201].

6.2. Training and Deployment Challenges

Due to the need for real-time data processing and enhanced privacy protection, local deployment of private LLMs and MLLMs is essential, particularly in medical settings. However, the substantial increase in parameters for LLMs and MLLMs significantly escalates the demand for computational resources, resulting in high operational costs. For example, MEDITRON-70B requires 128 A100 GPUs for training, and the smaller LLaVA-Med 7B necessitates 8 A100 GPUs. Furthermore, even after training and fine-tuning, the large sizes of medical LLMs and MLLMs result in high deployment and inference costs, presenting considerable challenges for most hospitals in locally deploying these models and applying them in practical settings. To facilitate the training and deployment of medical LLMs and MLLMs in hospitals with limited computational resources, this subsection proposes three solutions: optimizing the training process, reducing model parameters, and modifying model architectures.

To achieve efficient training of LLMs and MLLMs while minimizing computational overhead, researchers have proposed a series of parameter-efficient fine-tuning methods [202, 203, 204]. These approaches involve freezing the majority of parameters in LLMs and MLLMs and updating only a small subset to enable effective training. For example, Hu et al. [203] introduced LoRA, which freezes the pre-trained model weights and injects trainable rank decomposition matrices into each layer of the Transformer architecture to facilitate efficient training. However, while PEFT methods facilitate efficient training, they do not address the challenges associated with deployment. To address this issue, reducing model parameters and designing lightweight models are viable solutions [205, 206, 120]. For example, MobileVLM [205] is a customized MLLM for mobile scenarios, which reduces the training and inference budget by downsizing LLaMA and designing an efficient projector. It is capable of running on mobile devices while remaining competitive with other MLLMs across most tasks.

Currently, most LLMs and MLLMs are based on the Transformer architecture, which leads to a quadratic increase in computational complexity with sequence length, resulting in low efficiency for long sequences. To fundamentally address the challenges of training and deploying medical LLMs and MLLMs, selecting model architectures that are more efficient in computation and inference is a viable option [207, 208]. For example, RWKV [207] combines the efficient parallel training of Transformers with effective inference from RNNs, ensuring constant computational and memory complexity during inference while maintaining comparable performance to similarly scaled Transformer models. Furthermore, Mamba [208], based on the State Space Model, outperforms Transformer models in both performance

and inference speed, achieving five times the inference speed of Transformers while remaining comparable in scale. Extending these computationally and inference-efficient model architectures to medical LLMs and MLLMs will help address the current training and deployment challenges faced by these models.

6.3. Lack of Recency

Once medical LLMs and MLLMs are trained, the knowledge they acquire becomes static. However, as medical knowledge is continuously updated, the absence of new concepts will exacerbate the models' inaccuracies and hallucination problems. This is particularly evident when the models encounter new terms that are not present in the training corpus, rendering them unable to comprehend this knowledge [12]. For example, if medical LLMs and MLLMs are trained exclusively on data prior to 2020, they will lack information regarding COVID-19. This limitation may prevent the models from understanding terms such as "COVID-19" or "Long COVID," or they may incorrectly classify COVID-19 as a known viral variant (e.g., SARS or MERS), leading to the provision of misleading advice. Consequently, the lack of recent knowledge will significantly hinder the real-world application of medical LLMs and MLLMs.

To address the lack of recency resulting from the offline learning of medical LLMs and MLLMs, continual parameter updates through fine-tuning methods to synchronize them with current human knowledge is a feasible solution [155]. While fine-tuning can inject new medical concepts and knowledge into the model, it also introduces two challenges: catastrophic forgetting, in which the model forgets previously learned knowledge upon acquiring new information [209]. The second challenge is negative forward transfer, wherein performance on unseen tasks deteriorates when learning new tasks [210]. To address these issues, researchers have proposed model editing [211], one is introducing additional trainable parameters to correct erroneous responses stemming from outdated knowledge while preserving the original parameters of the model [212, 213]. Another approach involves identifying parameters related to specific knowledge and updating them accordingly to integrate relevant new information [214, 215, 216]. In addition to model editing, retrieval-augmented generation can be employed to update the knowledge of medical LLMs and MLLMs by linking the model to an information retrieval component. This allows the model to retrieve relevant content from external knowledge bases as references [16, 110], thereby generating more reliable responses.

6.4. Privacy and Security

Medical LLMs and MLLMs are trained on a large-scale medical corpus that includes data such as EHRs, doctor-patient dialogues, and other information that may involve patient privacy, including names, phone numbers, and email addresses. This information can potentially be directly extracted from medical LLMs or MLLMs using specific prompting methods [217, 218], raising significant privacy and security concerns.

Currently, a common practice to enhance patient privacy protection is data de-identification [16, 23]. This process involves removing or anonymizing sensitive information from datasets, including names, phone numbers, addresses, and medical record numbers. Additionally, sensitive terms in clinician-patient dialogues and medical records are replaced or removed to dissociate these terms from specific patients. Moreover, differential privacy methods can effectively mitigate the risk of privacy breaches by adding noise to obscure individual information in the training data, thus preventing the inference of specific details while still enabling meaningful data analysis [219]. Furthermore, utilizing high-quality synthetic data generated by models such as ChatGPT or GPT-4 during training ensures both the controllability and diversity of the datasets while mitigating the risk of privacy leaks. It is also advisable to monitor and filter the model's outputs. For instance, if the output contains sensitive data such as names or contact information, this information should be removed or modified during post-processing. Finally, we urge developers to adhere to ethical standards, ensuring that models respect privacy rights when handling patient data and provide patients with the right to access and delete their information.

6.5. Bias and Toxicity

Large-scale corpora, particularly data sourced from the internet, inevitably contain various biased viewpoints, which LLMs and MLLMs may learn [26, 220], including biases related to race [221], gender [222], and politics [223]. Additionally, language models may generate toxic responses, such as aggressive and hurtful remarks, with certain groups being more likely to be targeted due to these biases [161]. These biases and toxicities extend to LLMs and MLLMs, posing potential implications and threats to patients, and may have serious consequences for individuals with mental illness.

Reducing bias in training data is a fundamental approach to mitigating bias in models. Specifically, careful curation and screening of diverse, balanced, and representative training data ensure that models learn from a broader range

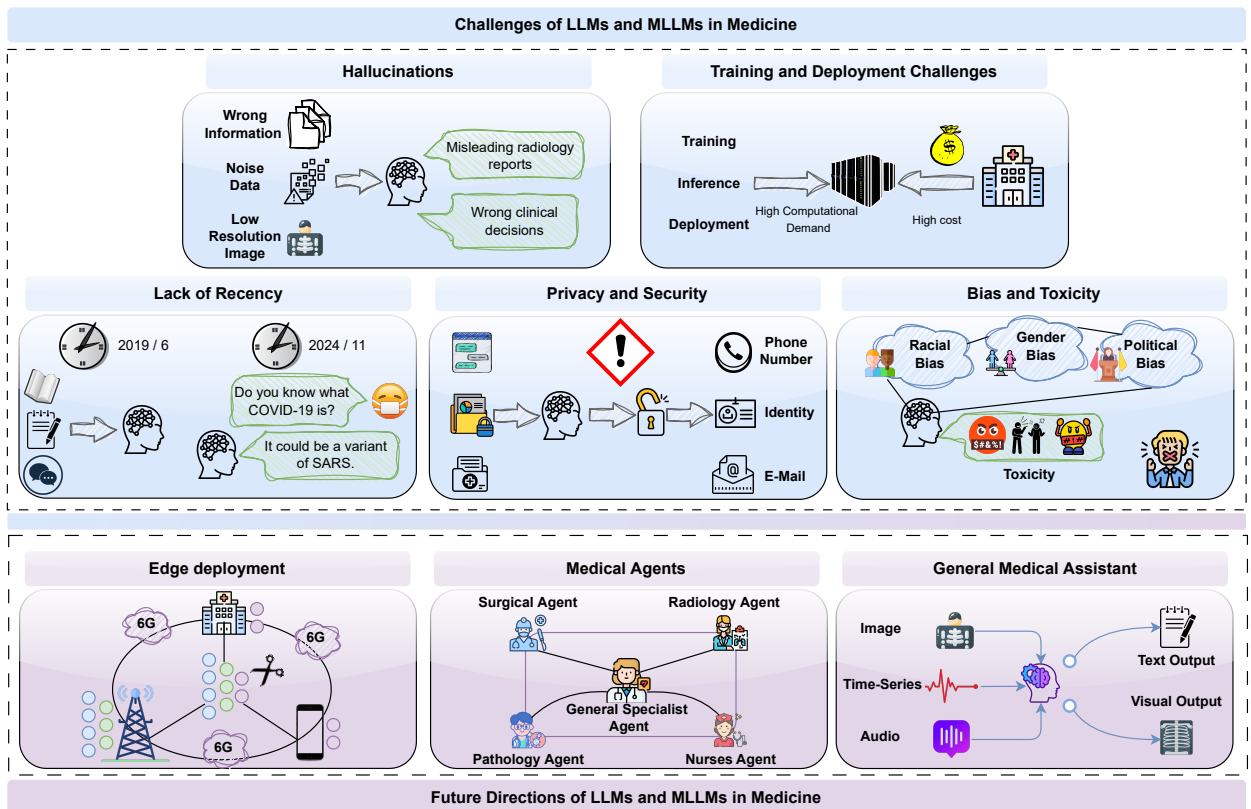


Fig. 10. Overview of challenges and future directions for medical LLMs and MLLMs in clinical settings.

of perspectives and experiences, leading to a more comprehensive understanding and reduced biases across various dimensions [220]. For addressing model toxicity, utilizing empathetic data has been shown to decrease the output of toxic content from models [224]. However, re-screening pre-training datasets and retraining models to minimize biases and toxicities can be costly. Therefore, screening high-quality datasets with anti-bias and anti-toxicity for fine-tuning is a more cost-effective approach. In addition to training, further enhancement is needed in evaluating model bias and toxicity. Designing a comprehensive benchmark for evaluating model bias and toxicity facilitates the detection of these issues, allowing developers to conduct regular reviews of the models [220].

7. Future Directions of LLMs and MLLMs in Medicine

This section examines future development trends in medical LLMs and MLLMs, envisioning their improved integration into clinical settings while offering forward-looking and insightful perspectives.

7.1. Edge Deployment

Due to their substantial computational and storage demands, existing LLMs and MLLMs are predominantly deployed on cloud servers. However, cloud-based deployment faces notable limitations in medical environments. First, in regions with poor network connectivity, communication latency with the cloud hinders the use of medical LLMs and MLLMs in critical clinical scenarios, such as surgical assistance. Second, transmitting data to cloud servers for cloud-based LLMs and MLLMs raises significant concerns about data privacy breaches. To address these challenges, deploying medical LLMs and MLLMs on edge devices is gaining traction as a promising future direction.

Edge deployment offers enhanced patient privacy protection and enables real-time responses [225]. Nevertheless, the limited computational capacity and memory of edge devices present significant challenges for deploying medical LLMs and MLLMs in such environments. Lin et al. [226] proposed a Mobile Edge Computing (MEC) architecture leveraging 6G networks to coordinate tasks between edge devices and servers as a solution to this challenge. For

example, in scenarios demanding high privacy, complete device-side inference can be utilized to ensure full patient privacy protection. In more complex scenarios, this architecture supports dynamic adjustments of model splitting and task allocation. For instance, in surgical assistance, the MEC architecture can deploy lightweight model components on edge devices to enable real-time responses, while offloading computationally intensive tasks to nearby edge servers. This setup supports distributed collaborative inference through 6G communication, reducing the burden on edge devices.

7.2. Medical Agents

Clinical diagnosis is a multidisciplinary process requiring collaboration among experts from fields such as pathology, radiology, and surgery. LLM- and MLLM-based agents can simulate departmental experts, facilitating team collaboration to deliver comprehensive and accurate diagnostic support [15].

For example, when handling a complex medical case, a general practitioner agent can decompose the task into clear sub-tasks based on the patient's specific conditions and requirements, dynamically assigning them to the most suitable expert agents. Each expert agent provides analysis specific to their assigned sub-tasks. A radiology agent analyzes imaging studies to generate a radiology report, a surgical agent develops a surgical plan, and a pathology agent examines tissue samples or blood test results to provide cellular- or molecular-level diagnostic insights. After each agent completes their assigned sub-tasks, a collaborative mechanism consolidates their diagnostic inputs into a comprehensive medical opinion, which is then communicated back to the general practitioner agent [227, 228]. Additionally, during agent collaboration, shared data and insights contribute to improved overall diagnostic recommendations. However, effective collaboration among multiple agents requires low-latency communication, especially in urgent scenarios like surgical assistance, where real-time interaction is essential. Consequently, 6G communication technology can also offer vital technical support for enabling seamless multi-agent collaboration.

7.3. General Medical Assistant

Although medical MLLMs currently support visual and text modalities, enabling tasks such as medical image report generation and visual question answering, clinical environments demand these models to process additional input modalities, including time series and audio data [229] and audio data. These additional input modalities can provide more comprehensive patient information, particularly playing a crucial role in dynamic monitoring and real-time diagnosis.

At present, medical LLMs and MLLMs primarily rely on static data, such as medical records and imaging studies, limiting their capacity to account for patients' dynamic data. Incorporating time series data allows medical MLLMs to analyze health trends and detect early signs of disease progression, leading to more precise alerts and predictions. For example, in an ICU, in an ICU, medical MLLMs can integrate continuous ECG monitoring and blood pressure variations to predict cardiac arrest risks and recommend timely interventions. Similarly, audio data, such as cough recordings, can aid in diagnosing lung diseases or predicting respiratory failure risks. Voice analysis by medical MLLMs can also assess a patient's emotional state, supporting psychiatrists in diagnosing mental health conditions like depression and anxiety [230, 231, 232].

Beyond additional input modalities, medical MLLMs are also expected to support a broader range of output modalities. Currently, most medical MLLMs rely on LLM components as output modules, restricting outputs to text. In clinical practice, medical MLLMs require visual output modules for region- and pixel-level outputs, such as segmenting pathological areas in images according to physician instructions to deliver more detailed diagnostic evidence.

8. Conclusion

In recent years, advancements in LLMs have driven significant breakthroughs in NLP, enabling researchers to make substantial progress toward artificial general intelligence by extending LLMs into the multimodal domain, resulting in the creation of MLLMs. Concurrently, the rapid development and impressive performance of general LLMs and MLLMs have spurred the emergence of numerous medical LLMs and MLLMs. This survey aims to help researchers and medical practitioners understand the technological advancements and developmental status of medical LLMs and MLLMs. It focuses on the paradigm shift of LLMs and MLLMs, highlighting their evolution from feature engineering to structure engineering, objective engineering, and now to prompt engineering and data engineering. The survey summarizes the mainstream architectures of current LLMs and MLLMs, compiles a list of existing medical LLMs and MLLMs, and provides insights into various architectures and model components. Additionally, it presents

a comprehensive guide on existing medical datasets, model fine-tuning approaches, and evaluation methods, aiding researchers and practitioners in building medical LLMs and MLLMs. The survey also examines the broad potential of LLMs and MLLMs in diverse clinical applications, such as medical diagnosis, clinical report generation, medical education, mental health services, and surgical assistance. Despite their remarkable achievements, medical LLMs and MLLMs face several critical challenges and limitations that hinder their practical deployment in clinical settings. Consequently, this survey addresses these challenges, including hallucinations, training and deployment issues, lack of recency, privacy and security, bias and toxicity, while offering potential solutions to support the practical application of future medical LLMs and MLLMs. Lastly, the survey explores future directions for medical LLMs and MLLMs, such as edge deployment, medical agents, and general medical assistant, offering forward-looking and insightful analysis. In conclusion, this survey delivers a comprehensive analysis of medical LLMs and MLLMs, covering their background, principles, applications, challenges, and future directions, with the aim of advancing their development in clinical medicine and fostering AI integration in healthcare.

References

- [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 2017.
- [2] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
- [3] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.
- [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- [5] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [6] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [7] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shrutvi Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [8] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [9] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [10] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [11] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in Neural Information Processing Systems*, 36:34892–34916, 2024.
- [12] Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature Medicine*, 29(8):1930–1940, 2023.
- [13] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, et al. Towards expert-level medical question answering with large language models. *arXiv preprint arXiv:2305.09617*, 2023.
- [14] Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421, 2021.
- [15] Hongjian Zhou, Boyang Gu, Xinyu Zou, Yiru Li, Sam S Chen, Peilin Zhou, Junling Liu, Yining Hua, Chengfeng Mao, Xian Wu, et al. A survey of large language models in medicine: Progress, application, and challenge. *arXiv preprint arXiv:2311.05112*, 2023.
- [16] Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang, and You Zhang. Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge. *Cureus*, 15(6), 2023.
- [17] Sheng Wang, Zihao Zhao, Xi Ouyang, Tianming Liu, Qian Wang, and Dinggang Shen. Interactive computer-aided diagnosis on medical image using large language models. *Communications Engineering*, 3(1):133, 2024.
- [18] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *Advances in Neural Information Processing Systems*, 36:28541–28564, 2024.
- [19] Dave Van Veen, Cara Van Uden, Louis Blankemeier, Jean-Benoit Delbrouck, Asad Aali, Christian Bluethgen, Anuj Pareek, Malgorzata Polacin, Eduardo Pontes Reis, Anna Seehofnerova, et al. Clinical text summarization: Adapting large language models can outperform human experts. *Research Square*, 2023.
- [20] Zhanyu Wang, Lingqiao Liu, Lei Wang, and Luping Zhou. R2gengpt: Radiology report generation with frozen llms. *Meta-Radiology*, 1(3):100033, 2023.

- [21] Tao Tu, Shekoofeh Azizi, Danny Driess, Mike Schaeckermann, Mohamed Amin, Pi-Chuan Chang, Andrew Carroll, Charles Lau, Ryutarō Tanno, Ira Ktena, et al. Towards generalist biomedical ai. *NEJM AI*, 1(3):AIoa2300138, 2024.
- [22] Yirong Chen, Xiaofen Xing, Jingkai Lin, Huimin Zheng, Zhenyu Wang, Qi Liu, and Xiangmin Xu. Soulchat: Improving llms' empathy, listening, and comfort abilities through fine-tuning with multi-turn empathy conversations. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1170–1183, dec 2023.
- [23] June M Liu, Donghao Li, He Cao, Tianhe Ren, Zeyi Liao, and Jiamin Wu. Chatcounselor: A large language models for mental health support. *arXiv preprint arXiv:2309.15461*, 2023.
- [24] Yunkun Zhang, Jin Gao, Zheling Tan, Lingfeng Zhou, Kexin Ding, Mu Zhou, Shaoting Zhang, and Dequan Wang. Data-centric foundation models in computational healthcare: A survey. *arXiv preprint arXiv:2401.02458*, 2024.
- [25] Michael Moor, Oishi Banerjee, Zahra Shakeri Hosseini Abad, Harlan M Krumholz, Jure Leskovec, Eric J Topol, and Pranav Rajpurkar. Foundation models for generalist medical artificial intelligence. *Nature*, 616(7956):259–265, 2023.
- [26] Jianing Qiu, Lin Li, Jiankai Sun, Jiachuan Peng, Peilun Shi, Ruiyang Zhang, Yinzhaohong, Kyle Lam, Frank P-W Lo, Bo Xiao, et al. Large ai models in health informatics: Applications, challenges, and the future. *IEEE Journal of Biomedical and Health Informatics*, 27(12):6074–6087, 2023.
- [27] Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*, 2023.
- [28] Logesh Kumar Umapathi, Ankit Pal, and Malaikannan Sankarasubbu. Med-halt: Medical domain hallucination test for large language models. *arXiv preprint arXiv:2307.15343*, 2023.
- [29] Vipula Rawte, Amit Sheth, and Amitava Das. A survey of hallucination in large foundation models. *arXiv preprint arXiv:2309.05922*, 2023.
- [30] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, 2023.
- [31] Qika Lin, Yifan Zhu, Xin Mei, Ling Huang, Jingying Ma, Kai He, Zhen Peng, Erik Cambria, and Mengling Feng. Has multimodal learning delivered universal intelligence in healthcare? a comprehensive survey. *Information Fusion*, page 102795, 2024.
- [32] Felix Krones, Umar Marikkar, Guy Parsons, Adam Szmul, and Adam Mahdi. Review of multimodal machine learning approaches in healthcare. *Information Fusion*, 114:102690, 2025.
- [33] Kai He, Rui Mao, Qika Lin, Yucheng Ruan, Xiang Lan, Mengling Feng, and Erik Cambria. A survey of large language models for healthcare: from data, technology, and applications to accountability and ethics. *arXiv preprint arXiv:2310.05694*, 2023.
- [34] Benyou Wang, Qianqian Xie, Jiahuan Pei, Zhihong Chen, Prayag Tiwari, Zhao Li, and Jie Fu. Pre-trained language models in biomedical domain: A systematic survey. *ACM Computing Surveys*, 56(3):1–52, 2023.
- [35] Surendrabikram Thapa and Surabhi Adhikari. Chatgpt, bard, and large language models for biomedical research: opportunities and pitfalls. *Annals of Biomedical Engineering*, 51(12):2647–2651, 2023.
- [36] Jesutofunmi A Omiye, Haiwen Gui, Shawheen J Rezaei, James Zou, and Roxana Daneshjou. Large language models in medicine: the potentials and pitfalls: a narrative review. *Annals of Internal Medicine*, 177(2):210–220, 2024.
- [37] Rajesh Bhayana. Chatbots and large language models in radiology: A practical primer for clinical and research applications. *Radiology*, 310(1):e232756, 2024.
- [38] Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36:55006–55021, 2024.
- [39] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023.
- [40] Huan Liu and Lei Yu. Toward integrating feature selection algorithms for classification and clustering. *IEEE Transactions on Knowledge and Data Engineering*, 17(4):491–502, 2005.
- [41] Franz Josef Och, Daniel Gildea, Sanjeev Khudanpur, Anoop Sarkar, Kenji Yamada, Alexander Fraser, Shankar Kumar, Libin Shen, David A Smith, Katherine Eng, et al. A smorgasbord of features for statistical machine translation. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 161–168, 2004.
- [42] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- [43] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- [44] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [45] Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23, 2021.
- [46] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240, 2020.
- [47] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [48] Bertalan Meskó. Prompt engineering as an important emerging skill for medical professionals: tutorial. *Journal of Medical Internet Research*, 25:e50638, 2023.
- [49] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414*, 2022.
- [50] Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, and Yejin Choi. Merlot: Multimodal neural script knowledge models. *Advances in Neural Information Processing Systems*, 34:23634–23651, 2021.
- [51] Lisa Anne Hendricks, John Mellor, Rosalia Schneider, Jean-Baptiste Alayrac, and Aida Nematzadeh. Decoupling the role of data, attention, and losses in multimodal transformers. *Transactions of the Association for Computational Linguistics*, 9:570–585, 2021.

- [52] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR, 18–24 Jul 2021.
- [53] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in Neural Information Processing Systems*, 34:9694–9705, 2021.
- [54] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.
- [55] Duzhen Zhang, Yahan Yu, Jiahua Dong, Chenxing Li, Dan Su, Chenhui Chu, and Dong Yu. MM-LLMs: Recent advances in MultiModal large language models. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 12401–12430. Association for Computational Linguistics, August 2024.
- [56] Yunfei Xie, Ce Zhou, Lang Gao, Juncheng Wu, Xianhang Li, Hong-Yu Zhou, Sheng Liu, Lei Xing, James Zou, Cihang Xie, et al. Medtrinity-25m: A large-scale multimodal dataset with multigranular annotations for medicine. *arXiv preprint arXiv:2408.02900*, 2024.
- [57] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- [58] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- [59] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *National Science Review*, page nwaec403, 11 2024.
- [60] Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180, 2023.
- [61] Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6268–6278. Association for Computational Linguistics, December 2023.
- [62] Haochun Wang, Chi Liu, Nuwa Xi, Zewen Qiang, Sendong Zhao, Bing Qin, and Ting Liu. Huatuo: Tuning llama model with chinese medical knowledge. *arXiv preprint arXiv:2304.06975*, 2023.
- [63] Tianyu Han, Lisa C Adams, Jens-Michalis Papaioannou, Paul Grundmann, Tom Oberhauser, Alexander Löser, Daniel Truhn, and Keno K Bresssem. Medalpaca—an open-source collection of medical conversational ai models and training data. *arXiv preprint arXiv:2304.08247*, 2023.
- [64] Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Weidi Xie, and Yanfeng Wang. Pmc-llama: toward building open-source language models for medicine. *Journal of the American Medical Informatics Association*, 31:1833–1843, 2024.
- [65] Augustin Toma, Patrick R Lawler, Jimmy Ba, Rahul G Krishnan, Barry B Rubin, and Bo Wang. Clinical camel: An open-source expert-level medical language model with dialogue-based knowledge encoding. *arXiv preprint arXiv:2305.12031*, 2023.
- [66] Hongbo Zhang, Junying Chen, Feng Jiang, Fei Yu, Zhihong Chen, Guiming Chen, Jianquan Li, Xiangbo Wu, Zhang Zhiyi, Qingying Xiao, Xiang Wan, Benyou Wang, and Haizhou Li. HuatuoGPT, towards taming language model to be a doctor. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10859–10885. Association for Computational Linguistics, December 2023.
- [67] Cheng Peng, Xi Yang, Aokun Chen, Kaleb E Smith, Nima PourNejatian, Anthony B Costa, Cheryl Martin, Mona G Flores, Ying Zhang, Tanja Magoc, et al. A study of generative large language model for medical research and healthcare. *NPJ Digital Medicine*, 6(1):210, 2023.
- [68] Guangyu Wang, Guoxing Yang, Zongxin Du, Longjun Fan, and Xiaohu Li. Clinicalgpt: large language models finetuned with diverse medical data and comprehensive evaluation. *arXiv preprint arXiv:2306.09968*, 2023.
- [69] Songhua Yang, Hanjie Zhao, Senbin Zhu, Guangyu Zhou, Hongfei Xu, Yuxiang Jia, and Hongying Zan. Zhongjing: Enhancing the chinese medical capabilities of large language model through expert feedback and real-world multi-turn dialogue. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 19368–19376, 2024.
- [70] Zhengliang Liu, Yiwei Li, Peng Shu, Aoxiao Zhong, Longtao Yang, Chao Ju, Zihao Wu, Chong Ma, Jie Luo, Cheng Chen, et al. Radiology-llama2: Best-in-class large language model for radiology. *arXiv preprint arXiv:2309.06419*, 2023.
- [71] Yang Tan, Mingchen Li, Zijie Huang, Huiqun Yu, and Guisheng Fan. Medchatz: a better medical adviser learns from better instructions. *arXiv preprint arXiv:2309.01114*, 2023.
- [72] Ofir Ben Shoham and Nadav Rappoport. Cpllm: Clinical prediction with large language models. *arXiv preprint arXiv:2309.11295*, 2023.
- [73] Qichen Ye, Junling Liu, Dading Chong, Peilin Zhou, Yining Hua, and Andrew Liu. Qilin-med: Multi-stage knowledge injection advanced medical large language model. *arXiv preprint arXiv:2310.09089*, 2023.
- [74] Xinlu Zhang, Chenxin Tian, Xianjun Yang, Lichang Chen, Zekun Li, and Linda Ruth Petzold. Alpacare: Instruction-tuned large language models for medical application. *arXiv preprint arXiv:2310.14558*, 2023.
- [75] Guoxing Yang, Xiaohong Liu, Jianyu Shi, Zan Wang, and Guangyu Wang. Tcm-gpt: Efficient pre-training of large language models for domain adaptation in traditional chinese medicine. *Computer Methods and Programs in Biomedicine Update*, 6:100158, 2024.
- [76] Junying Chen, Xidong Wang, Ke Ji, Anningzhe Gao, Feng Jiang, Shunian Chen, Hongbo Zhang, Song Dingjie, Wenya Xie, Chuyi Kong, Jianquan Li, Xiang Wan, Haizhou Li, and Benyou Wang. HuatuoGPT-II, one-stage training for medical adaptation of LLMs. In *First Conference on Language Modeling*, 2024.
- [77] Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, et al. Meditron-70b: Scaling medical pretraining for large language models. *arXiv preprint arXiv:2311.16079*, 2023.
- [78] Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. Biomistral: A collection of open-source pretrained large language models for medical domains. *arXiv preprint arXiv:2402.10373*, 2024.

- [79] Qianqian Xie, Qingyu Chen, Aokun Chen, Cheng Peng, Yan Hu, Fongci Lin, Xueqing Peng, Jimin Huang, Jeffrey Zhang, Vipina Keloth, et al. Me llama: Foundation large language models for medical applications. *arXiv preprint arXiv:2402.12749*, 2024.
- [80] Xidong Wang, Nuo Chen, Junyin Chen, Yan Hu, Yidong Wang, Xiangbo Wu, Anningzhe Gao, Xiang Wan, Haizhou Li, and Benyou Wang. Apollo: Lightweight multilingual medical llms towards democratizing medical ai to 6b people. *arXiv preprint arXiv:2403.03640*, 2024.
- [81] Elliot Bolton, Abhinav Venigalla, Michihiro Yasunaga, David Hall, Betty Xiong, Tony Lee, Roxana Daneshjou, Jonathan Frankle, Percy Liang, Michael Carbin, et al. Biomedlm: A 2.7 b parameter language model trained on biomedical text. *arXiv preprint arXiv:2403.18421*, 2024.
- [82] Dingkang Yang, Jinjie Wei, Dongling Xiao, Shunli Wang, Tong Wu, Gang Li, Mingcheng Li, Shuaibing Wang, Jiawei Chen, Yue Jiang, Qingyao Xu, Ke Li, Peng Zhai, and Lihua Zhang. PediatricsGPT: Large language models as chinese medical assistants for pediatric applications. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [83] Honglin Xiong, Sheng Wang, Yitao Zhu, Zihao Zhao, Yuxiao Liu, Linlin Huang, Qian Wang, and Dinggang Shen. Doctorglm: Fine-tuning your chinese doctor is not a herculean task. *arXiv preprint arXiv:2304.01097*, 2023.
- [84] Yirong Chen, Zhenyu Wang, Xiaofen Xing, Zhipei Xu, Kai Fang, Junhong Wang, Sihang Li, Jieling Wu, Qi Liu, Xiangmin Xu, et al. Bianque: Balancing the questioning and suggestion ability of health llms with multi-turn health conversations polished by chatgpt. *arXiv preprint arXiv:2310.15896*, 2023.
- [85] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*, 2020.
- [86] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.
- [87] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [88] Shaoxiong Ji, Tianlin Zhang, Luna Ansari, Jie Fu, Prayag Tiwari, and Erik Cambria. MentalBERT: Publicly available pretrained language models for mental healthcare. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 7184–7190. European Language Resources Association, June 2022.
- [89] Thomas Wang, Adam Roberts, Daniel Hesslow, Teven Le Scao, Hyung Won Chung, Iz Beltagy, Julien Launay, and Colin Raffel. What language model architecture and pretraining objective works best for zero-shot generalization? In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162, pages 22964–22984. PMLR, 17–23 Jul 2022.
- [90] Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. Why can GPT learn in-context? language models secretly perform gradient descent as meta-optimizers. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 4005–4019. Association for Computational Linguistics, July 2023.
- [91] Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. GLM: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, pages 320–335. Association for Computational Linguistics, May 2022.
- [92] Chang Shu, Baian Chen, Fangyu Liu, Zihao Fu, Ehsan Shareghi, and Nigel Collier. Visual med-alpaca: A parameter-efficient biomedical llm with visual capabilities. <https://github.com/cambridgeltl/visual-med-alpaca>, 2023.
- [93] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [94] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- [95] Xi Chen, Xiao Wang, Lucas Beyer, Alexander Kolesnikov, Jialin Wu, Paul Voigtlaender, Basil Mustafa, Sebastian Goodman, Ibrahim Alabdulmohsin, Piotr Padlewski, et al. Pali-3 vision language models: Smaller, faster, stronger. *arXiv preprint arXiv:2310.09199*, 2023.
- [96] Shaurya Srivastav, Mercy Ranjit, Fernando Pérez-García, Kenza Bouzid, Shruthi Bannur, Daniel C. Castro, Anton Schwaighofer, Harshita Sharma, Maximilian Ilse, Valentina Salvatelli, Sam Bond-Taylor, Fabian Falck, Anja Thieme, Hannah Richardson, Matthew P. Lungren, Stephanie L. Hyland, and Javier Alvarez-Valle. MAIRA at RRG24: A specialised large multimodal model for radiology report generation. In *Proceedings of the 23rd Workshop on Biomedical Natural Language Processing*, pages 597–602. Association for Computational Linguistics, August 2024.
- [97] Ming Y Lu, Bowen Chen, Drew FK Williamson, Richard J Chen, Melissa Zhao, Aaron K Chow, Kenji Ikemura, Ahrong Kim, Dimitra Pouli, Ankush Patel, et al. A multimodal generative ai copilot for human pathology. *Nature*, 634:466–473, 2024.
- [98] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19358–19369, 2023.
- [99] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023.
- [100] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [101] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- [102] Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Yash Dalmia, Jure Leskovec, Cyril Zalka, Eduardo Pontes Reis, and Pranav Rajpurkar. Med-flamingo: a multimodal medical few-shot learner. In *Proceedings of the 3rd Machine Learning for Health Symposium*, volume 225, pages 353–367. PMLR, 10 Dec 2023.

- [103] Qiuhui Chen, Xinyue Hu, Zirui Wang, and Yi Hong. Medblip: Bootstrapping language-image pre-training from 3d medical images and texts. *arXiv preprint arXiv:2305.10799*, 2023.
- [104] Rongsheng Wang, Yaofei Duan, Junrong Li, Patrick Pang, and Tao Tan. Xrayglm: The first chinese medical multimodal model that chest radiographs summarization. <https://github.com/WangRongsheng/XrayGLM>, 2023.
- [105] Bang Yang, Asif Raza, Yuexian Zou, and Tong Zhang. Customizing general-purpose foundation models for medical report generation. *arXiv preprint arXiv:2306.05642*, 2023.
- [106] Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Towards generalist foundation model for radiology. *arXiv preprint arXiv:2308.02463*, 2023.
- [107] Zhihong Chen, Maya Varma, Jean-Benoit Delbrouck, Magdalini Paschali, Louis Blankemeier, Dave Van Veen, Jeya Maria Jose Valanarasu, Alaa Youssef, Joseph Paul Cohen, Eduardo Pontes Reis, Emily Tsai, Andrew Johnston, Cameron Olsen, Tanishq Mathew Abraham, Sergios Gatidis, Akshay S Chaudhari, and Curtis Langlotz. Chexagent: Towards a foundation model for chest x-ray interpretation. In *AAAI 2024 Spring Symposium on Clinical Foundation Models*, 2024.
- [108] Tom Van Sonsbeek, Mohammad Mahdi Derakhshani, Ivona Najdenkoska, Cees GM Snoek, and Marcel Worring. Open-ended medical visual question answering through prefix tuning of language models. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 726–736. Springer, 2023.
- [109] Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. Pmc-vqa: Visual instruction tuning for medical visual question answering. *arXiv preprint arXiv:2305.10415*, 2023.
- [110] Yuxuan Sun, Chenglu Zhu, Sunyi Zheng, Kai Zhang, Lin Sun, Zhongyi Shui, Yunlong Zhang, Honglin Li, and Lin Yang. Pathasst: A generative foundation ai assistant towards artificial general intelligence of pathology. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 5034–5042, 2024.
- [111] Omkar Chakradhar Thawakar, Abdelrahman M. Shaker, Sahal Shaji Mullappilly, Hisham Cholakkal, Rao Muhammad Anwer, Salman Khan, Jorma Laaksonen, and Fahad Khan. XrayGPT: Chest radiographs summarization using large medical vision-language models. In *Proceedings of the 23rd Workshop on Biomedical Natural Language Processing*, pages 440–448. Association for Computational Linguistics, August 2024.
- [112] Junling Liu, Ziming Wang, Qichen Ye, Dading Chong, Peilin Zhou, and Yining Hua. Qilin-med-vl: Towards chinese large vision-language model for general healthcare. *arXiv preprint arXiv:2310.17956*, 2023.
- [113] Jinlong He, Pengfei Li, Gang Liu, Zixu Zhao, and Shenjun Zhong. Pefomed: Parameter efficient fine-tuning on multimodal large language models for medical visual question answering. *arXiv preprint arXiv:2401.02797*, 2024.
- [114] Fan Bai, Yuxin Du, Tiejun Huang, Max Q-H Meng, and Bo Zhao. M3d: Advancing 3d medical image analysis with multi-modal large language models. *arXiv preprint arXiv:2404.00578*, 2024.
- [115] Songtao Jiang, Tuo Zheng, Yan Zhang, Yeying Jin, and Zuozhu Liu. Moe-tinymed: Mixture of experts for tiny medical large vision-language models. *CoRR*, 2024.
- [116] Shruthi Bannur, Kenza Bouzid, Daniel C Castro, Anton Schwaighofer, Sam Bond-Taylor, Maximilian Ilse, Fernando Pérez-García, Valentina Salvatelli, Harshita Sharma, Felix Meissen, et al. Maira-2: Grounded radiology report generation. *arXiv preprint arXiv:2406.04449*, 2024.
- [117] Junying Chen, Chi Gui, Ruyi Ouyang, Anningzhe Gao, Shunian Chen, Guiming Hardy Chen, Xidong Wang, Zhenyang Cai, Ke Ji, Xiang Wan, and Benyou Wang. Towards injecting medical visual knowledge into multimodal LLMs at scale. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7346–7370. Association for Computational Linguistics, November 2024.
- [118] Asma Alkhalidi, Raneem Alnajim, Layan Alabdullatef, Rawan Alyahya, Jun Chen, Deyao Zhu, Ahmed Alsinan, and Mohamed Elhoseiny. Minigt-med: Large language model as a general interface for radiology diagnosis. *arXiv preprint arXiv:2407.04106*, 2024.
- [119] Juexiao Zhou, Xiaonan He, Liyuan Sun, Jiannan Xu, Xiuying Chen, Yuetan Chu, Longxi Zhou, Xingyu Liao, Bin Zhang, Shawn Afvari, et al. Pre-trained multimodal large language model enhances dermatological diagnosis using skingpt-4. *Nature Communications*, 15(1):5649, 2024.
- [120] Feizhong Zhou, Xingyue Liu, Qiao Zeng, Zhuhan Li, and Hanguang Xiao. Sigphi-med: A lightweight vision-language assistant for biomedicine. Available at SSRN 4988925, 2024.
- [121] Bihao Zhao, Sheng Wang, Jinchun Gu, Yitao Zhu, Lanzhuji Mei, Zixu Zhuang, Zhiming Cui, Qian Wang, and Dinggang Shen. Chatcad+: Toward a universal and reliable interactive cad using llms. *IEEE Transactions on Medical Imaging*, 43(11):3755–3766, 2024.
- [122] Weihao Gao, Zhuo Deng, Zhiyuan Niu, Fujun Rong, Chucheng Chen, Zheng Gong, Wenze Zhang, Daimin Xiao, Fang Li, Zhenjie Cao, et al. Ophglm: Training an ophthalmology large language-and-vision assistant based on instructions and dialogue. *arXiv preprint arXiv:2306.12174*, 2023.
- [123] Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira. Perceiver: General perception with iterative attention. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pages 4651–4664. PMLR, 18–24 Jul 2021.
- [124] Shezheng Song, Xiaopeng Li, and Shasha Li. How to bridge the gap between modalities: A comprehensive survey on multimodal large language model. *arXiv preprint arXiv:2311.07594*, 2023.
- [125] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 19730–19742. PMLR, 23–29 Jul 2023.
- [126] Linli Yao, Lei Li, Shuhuai Ren, Lean Wang, Yuanxin Liu, Xu Sun, and Lu Hou. Deco: Decoupling token compression from semantic abstraction in multimodal large language models. *arXiv preprint arXiv:2405.20985*, 2024.
- [127] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 26286–26296, 2024.
- [128] Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific Data*, 3(1):1–9, 2016.
- [129] Alistair EW Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, et al. Mimic-iv, a freely accessible electronic health record dataset. *Scientific Data*, 10(1):1, 2023.

- [130] Emily Herrett, Arlene M Gallagher, Krishnan Bhaskaran, Harriet Forbes, Rohini Mathur, Tjeerd Van Staa, and Liam Smeeth. Data resource profile: clinical practice research datalink (cprd). *International Journal of Epidemiology*, 44(3):827–836, 2015.
- [131] Lucy Lu Wang, Kyle Lo, Yoganand Chandrasekhar, Russell Reas, Jiangjiang Yang, Douglas Burdick, Darrin Eide, Kathryn Funk, Yannik Katsis, Rodney Kinney, et al. *CORD-19: The COVID-19 open research dataset*. In *Proceedings of the 1st Workshop on NLP for COVID-19 at ACL 2020*. Association for Computational Linguistics, July 2020.
- [132] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua 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)*, pages 2567–2577. Association for Computational Linguistics, November 2019.
- [133] Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In Gerardo Flores, George H Chen, Tom Pollard, Joyce C Ho, and Tristan Naumann, editors, *Proceedings of the Conference on Health, Inference, and Learning*, volume 174 of *Proceedings of Machine Learning Research*, pages 248–260. PMLR, 07–08 Apr 2022.
- [134] Sheng Zhang, Xin Zhang, Hui Wang, Lixiang Guo, and Shanshan Liu. Multi-scale attentive interaction networks for chinese medical question answer selection. *IEEE Access*, 6:74061–74071, 2018.
- [135] Asma Ben Abacha and Dina Demner-Fushman. A question-entailment approach to question answering. *BMC Bioinformatics*, 20:1–23, 2019.
- [136] Jianquan Li, Xidong Wang, Xiangbo Wu, Zhiyi Zhang, Xiaolong Xu, Jie Fu, Prayag Tiwari, Xiang Wan, and Benyou Wang. Huatuo-26m, a large-scale chinese medical qa dataset. *arXiv preprint arXiv:2305.01526*, 2023.
- [137] Guangtao Zeng, Wenmian Yang, Zeqian Ju, Yue Yang, Sicheng Wang, Ruisi Zhang, Meng Zhou, Jiaqi Zeng, Xiangyu Dong, Ruoyu Zhang, et al. Meddialog: Large-scale medical dialogue datasets. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, pages 9241–9250, 2020.
- [138] Olivier Bodenreider. The unified medical language system (umls): integrating biomedical terminology. *Nucleic Acids Research*, 32(suppl_1):D267–D270, 2004.
- [139] Odma Byambasuren, Yunfei Yang, Zhifang Sui, Damai Dai, Baobao Chang, Sujian Li, and Hongying Zan. Preliminary study on the construction of chinese medical knowledge graph. *Journal of Chinese Information Processing*, 33(10):1–9, 2019.
- [140] Marco Basaldella, Fanguy Liu, Ehsan Shareghi, and Nigel Collier. COMETA: A corpus for medical entity linking in the social media. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3122–3137. Association for Computational Linguistics, November 2020.
- [141] Jason J Lau, Soumya Gayen, Asma Ben Abacha, and Dina Demner-Fushman. A dataset of clinically generated visual questions and answers about radiology images. *Scientific Data*, 5(1):1–10, 2018.
- [142] Bo Liu, Li-Ming Zhan, Li Xu, Lin Ma, Yan Yang, and Xiao-Ming Wu. Slake: A semantically-labeled knowledge-enhanced dataset for medical visual question answering. In *2021 IEEE 18th International Symposium on Biomedical Imaging*, pages 1650–1654. IEEE, 2021.
- [143] Xuehai He, Yichen Zhang, Luntian Mou, Eric Xing, and Pengtao Xie. Pathvqa: 30000+ questions for medical visual question answering. *arXiv preprint arXiv:2003.10286*, 2020.
- [144] Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Roger G Mark, and Steven Horng. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific Data*, 6(1):317, 2019.
- [145] Dina Demner-Fushman, Marc D Kohli, Marc B Rosenman, Sonya E Shooshan, Laritza Rodriguez, Sameer Antani, George R Thoma, and Clement J McDonald. Preparing a collection of radiology examinations for distribution and retrieval. *Journal of the American Medical Informatics Association*, 23(2):304–310, 2016.
- [146] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 590–597, 2019.
- [147] Obioma Pelka, Sven Koitka, Johannes Rückert, Felix Nensa, and Christoph 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*, pages 180–189. Springer International Publishing, 2018.
- [148] Zhi Huang, Federico Bianchi, Mert Yuksekogonul, Thomas J Montine, and James Zou. A visual–language foundation model for pathology image analysis using medical twitter. *Nature Medicine*, 29(9):2307–2316, 2023.
- [149] Sanjay Subramanian, Lucy Lu Wang, Ben Bogin, Sachin Mehta, Madeleine van Zuylen, Sravanthi Parasa, Sameer Singh, Matt Gardner, and Hannaneh Hajishirzi. MedICaT: A dataset of medical images, captions, and textual references. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2112–2120. Association for Computational Linguistics, November 2020.
- [150] Weixiong Lin, Ziheng Zhao, Xiaoman Zhang, Chaoyi Wu, Ya Zhang, Yanfeng Wang, and Weidi Xie. Pmc-clip: Contrastive language-image pre-training using biomedical documents. In *Medical Image Computing and Computer Assisted Intervention*, pages 525–536. Springer, 2023.
- [151] Sheng Zhang, Yanbo Xu, Naoto Usuyama, Jaspreet Bagga, Robert Tinn, Sam Preston, Rajesh Rao, Mu Wei, Naveen Valluri, Cliff Wong, et al. Large-scale domain-specific pretraining for biomedical vision-language processing. *arXiv preprint arXiv:2303.00915*, 2(3):6, 2023.
- [152] Tao Tu, Anil Palepu, Mike Schaeckermann, Khaled Saab, Jan Freyberg, Ryutaro Tanno, Amy Wang, Brenna Li, Mohamed Amin, Nenad Tomasev, et al. Towards conversational diagnostic ai. *arXiv preprint arXiv:2401.05654*, 2024.
- [153] Xi Yang, Aokun Chen, Nima PourNejatian, Hoo Chang Shin, Kaleb E Smith, Christopher Parisien, Colin Compas, Cheryl Martin, Anthony B Costa, Mona G Flores, et al. A large language model for electronic health records. *NPJ Digital Medicine*, 5(1):194, 2022.
- [154] Ruixiang Tang, Xiaotian Han, Xiaoqian Jiang, and Xia Hu. Does synthetic data generation of llms help clinical text mining? *arXiv preprint arXiv:2303.04360*, 2023.

- [155] Tongtong Wu, Linhao Luo, Yuan-Fang Li, Shirui Pan, Thuy-Trang Vu, and Gholamreza Haffari. Continual learning for large language models: A survey. *arXiv preprint arXiv:2402.01364*, 2024.
- [156] Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*, 2022.
- [157] Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open problems and fundamental limitations of reinforcement learning from human feedback. *Transactions on Machine Learning Research*, 2023. Survey Certification.
- [158] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020.
- [159] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- [160] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- [161] Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. Toxicity in chatgpt: Analyzing persona-assigned language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1236–1270. Association for Computational Linguistics, December 2023.
- [162] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, page 311–318. Association for Computational Linguistics, 2002.
- [163] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81. Association for Computational Linguistics, 2004.
- [164] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016.
- [165] Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. A diversity-promoting objective function for neural conversation models. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119. Association for Computational Linguistics, 2016.
- [166] Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *2015 IEEE Conference on Computer Vision and Pattern Recognition*, pages 4566–4575, 2015.
- [167] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*, 2020.
- [168] Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. Is ChatGPT a good NLG evaluator? a preliminary study. In *Proceedings of the 4th New Frontiers in Summarization Workshop*, pages 1–11. Association for Computational Linguistics, December 2023.
- [169] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Advances in Neural Information Processing Systems*, 2024.
- [170] Peter Szolovits, Ramesh S Patil, and William B Schwartz. Artificial intelligence in medical diagnosis. *Annals of Internal Medicine*, 108(1):80–87, 1988.
- [171] J Yuan, P Bao, Z Chen, M Yuan, J Zhao, J Pan, Y Xie, Y Cao, Y Wang, Z Wang, et al. Advanced prompting as a catalyst: Empowering large language models in the management of gastrointestinal cancers. *The Innovation*, 521, 2023.
- [172] Lingxuan Zhu, Weiming Mou, and Rui Chen. Can the chatgpt and other large language models with internet-connected database solve the questions and concerns of patient with prostate cancer and help democratize medical knowledge? *Journal of Translational Medicine*, 21(1):269, 2023.
- [173] Fenglin Liu, Tingting Zhu, Xian Wu, Bang Yang, Chenyu You, Chenyang Wang, Lei Lu, Zhangdaihong Liu, Yefeng Zheng, Xu Sun, et al. A medical multimodal large language model for future pandemics. *NPJ Digital Medicine*, 6(1):226, 2023.
- [174] Peter Lee, Sebastien Bubeck, and Joseph Petro. Benefits, limits, and risks of gpt-4 as an ai chatbot for medicine. *New England Journal of Medicine*, 388(13):1233–1239, 2023.
- [175] Stephen R Ali, Thomas D Dobbs, Hayley A Hutchings, and Iain S Whitaker. Using chatgpt to write patient clinic letters. *The Lancet Digital Health*, 5(4):e179–e181, 2023.
- [176] Sajjan B Patel and Kyle Lam. Chatgpt: the future of discharge summaries? *The Lancet Digital Health*, 5(3):e107–e108, 2023.
- [177] Reece Alexander James Clough, William Anthony Sparkes, Oliver Thomas Clough, Joshua Thomas Sykes, Alexander Thomas Steventon, and Kate King. Transforming healthcare documentation: harnessing the potential of ai to generate discharge summaries. *BJGP Open*, 2024.
- [178] Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. Capabilities of gpt-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*, 2023.
- [179] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. The dawn of llms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 9(1):1, 2023.
- [180] Zhichao Yang, Zonghai Yao, Mahbuba Tasmin, Parth Vashisht, Won Seok Jang, Feiyun Ouyang, Beining Wang, Dan Berlowitz, and Hong Yu. Performance of multimodal gpt-4v on usmle with image: Potential for imaging diagnostic support with explanations. *medRxiv*, pages 2023–10, 2023.
- [181] Sal Khan. Harnessing gpt-4 so that all students benefit. a nonprofit approach for equal access. *Khan Academy Blog*, 2023.
- [182] Duolingo Team. Introducing duolingo max, a learning experience powered by gpt-4. *Retrieved March*, 15:2023, 2023.

- [183] Mert Karabacak, Burak Berksu Ozkara, Konstantinos Margetis, Max Wintermark, and Sotirios Bisdas. The advent of generative language models in medical education. *JMIR Medical Education*, 9:e48163, 2023.
- [184] Hyunsu Lee. The rise of chatgpt: Exploring its potential in medical education. *Anatomical sciences education*, 17(5):926–931, 2024.
- [185] Qing Lyu, Josh Tan, Michael E Zapadka, Janardhana Ponnatapura, Chuang Niu, Kyle J Myers, Ge Wang, and Christopher T Whitlow. Translating radiology reports into plain language using chatgpt and gpt-4 with prompt learning: results, limitations, and potential. *Visual Computing for Industry, Biomedicine, and Art*, 6(1):9, 2023.
- [186] Zhiyong Han, Fortunato Battaglia, Abinav Udaiyar, Allen Fooks, and Stanley R Terlecky. An explorative assessment of chatgpt as an aid in medical education: Use it with caution. *Medical Teacher*, pages 1–8, 2023.
- [187] Sangzin Ahn. The impending impacts of large language models on medical education. *Korean Journal of Medical Education*, 35(1):103, 2023.
- [188] Alastair C van Heerden, Julia R Pozuelo, and Brandon A Kohrt. Global mental health services and the impact of artificial intelligence–powered large language models. *JAMA Psychiatry*, 80(7):662–664, 2023.
- [189] Munmun De Choudhury, Sachin R Pendse, and Neha Kumar. Benefits and harms of large language models in digital mental health. *arXiv preprint arXiv:2311.14693*, 2023.
- [190] Anna Stock, Stephan Schlögl, and Aleksander Groth. Tell me, what are you most afraid of? exploring the effects of agent representation on information disclosure in human-chatbot interaction. In *International Conference on Human-Computer Interaction*, pages 179–191. Springer, 2023.
- [191] Ana Paula Chaves and Marco Aurelio Gerosa. How should my chatbot interact? a survey on social characteristics in human–chatbot interaction design. *International Journal of Human–Computer Interaction*, 37(8):729–758, 2021.
- [192] Long Bai, Guankun Wang, Mobarakol Islam, Lalithkumar Seenivasan, An Wang, and Hongliang Ren. Surgical-vqla++: Adversarial contrastive learning for calibrated robust visual question-localized answering in robotic surgery. *Information Fusion*, 113:102602, 2025.
- [193] Ranjit Barua. Innovations in minimally invasive surgery: The rise of smart flexible surgical robots. In *Emerging Technologies for Health Literacy and Medical Practice*, pages 110–131. IGI Global, 2024.
- [194] Lalithkumar Seenivasan, Mobarakol Islam, Gokul Kannan, and Hongliang Ren. Surgicalgpt: End-to-end language-vision gpt for visual question answering in surgery. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 281–290. Springer, 2023.
- [195] Yihan Cao, Yanbin Kang, Chi Wang, and Lichao Sun. Instruction mining: Instruction data selection for tuning large language models. In *First Conference on Language Modeling*, 2024.
- [196] Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srinivasan, Tianyi Zhou, Heng Huang, and Hongxia Jin. Alpapasus: Training a better alpaca with fewer data. In *The Twelfth International Conference on Learning Representations*, 2024.
- [197] Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. Chain-of-verification reduces hallucination in large language models. In *Findings of the Association for Computational Linguistics*, pages 3563–3578. Association for Computational Linguistics, August 2024.
- [198] Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. Retrieval augmentation reduces hallucination in conversation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3784–3803. Association for Computational Linguistics, November 2021.
- [199] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. DINOv2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2024.
- [200] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023.
- [201] Zhuofan Zong, Bingqi Ma, Dazhong Shen, Guanglu Song, Hao Shao, Dongzhi Jiang, Hongsheng Li, and Yu Liu. MoVA: Adapting mixture of vision experts to multimodal context. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [202] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, pages 4582–4597. Association for Computational Linguistics, August 2021.
- [203] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- [204] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR, 2019.
- [205] Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu Zhang, Bo Zhang, Xiaolin Wei, et al. Mobilevlm: A fast, reproducible and strong vision language assistant for mobile devices. *arXiv preprint arXiv:2312.16886*, 2023.
- [206] Zhengqing Yuan, Zhaoxu Li, Weiran Huang, Yanfang Ye, and Lichao Sun. TinyGPT-v: Efficient multimodal large language model via small backbones. In *2nd Workshop on Advancing Neural Network Training: Computational Efficiency, Scalability, and Resource Optimization*, 2024.
- [207] Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, Kranthi Kiran GV, et al. RWKV: Reinventing RNNs for the transformer era. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14048–14077. Association for Computational Linguistics, December 2023.
- [208] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. In *First Conference on Language Modeling*, 2024.

- [209] Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma. Investigating the catastrophic forgetting in multimodal large language model fine-tuning. In *Conference on Parsimony and Learning*, 2023.
- [210] Junhao Zheng, Qianli Ma, Zhen Liu, Binquan Wu, and Huawen Feng. Beyond anti-forgetting: Multimodal continual instruction tuning with positive forward transfer. *arXiv preprint arXiv:2401.09181*, 2024.
- [211] Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. Editing large language models: Problems, methods, and opportunities. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10222–10240. Association for Computational Linguistics, December 2023.
- [212] Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong. Transformer-patcher: One mistake worth one neuron. *arXiv preprint arXiv:2301.09785*, 2023.
- [213] Tom Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. Aging with grace: Lifelong model editing with discrete key-value adaptors. In *Advances in Neural Information Processing Systems*, volume 36, pages 47934–47959. Curran Associates, Inc., 2023.
- [214] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.
- [215] Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. In *The Eleventh International Conference on Learning Representations*, 2023.
- [216] Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun Ma, and Jie Yu. Pmet: Precise model editing in a transformer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18564–18572, 2024.
- [217] Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650, 2021.
- [218] Haoran Li, Dadi Guo, Wei Fan, Mingshi Xu, Jie Huang, Fanpu Meng, and Yangqiu Song. Multi-step jailbreaking privacy attacks on ChatGPT. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4138–4153, December 2023.
- [219] Safiye Turgay, İlker İltter, et al. Perturbation methods for protecting data privacy: A review of techniques and applications. *Automation and Machine Learning*, 4(2):31–41, 2023.
- [220] Emilio Ferrara. Should chatgpt be biased? challenges and risks of bias in large language models. *arXiv preprint arXiv:2304.03738*, 2023.
- [221] Yifan Yang, Xiaoyu Liu, Qiao Jin, Furong Huang, and Zhiyong Lu. Unmasking and quantifying racial bias of large language models in medical report generation. *Nature Medicine*, 4:176, 2024.
- [222] Hadas Kotek, Rikker Dockum, and David Sun. Gender bias and stereotypes in large language models. In *Proceedings of The ACM Collective Intelligence Conference*, pages 12–24, 2023.
- [223] Ruibo Liu, Chenyan Jia, Jason Wei, Guangxuan Xu, and Soroush Vosoughi. Quantifying and alleviating political bias in language models. *Artificial Intelligence*, 304:103654, 2022.
- [224] Allison Lahnala, Charles Welch, Béla Neuendorf, and Lucie Flek. Mitigating toxic degeneration with empathetic data: Exploring the relationship between toxicity and empathy. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4926–4938. Association for Computational Linguistics, July 2022.
- [225] Minrui Xu, Hongyang Du, Dusit Niyato, Jiawen Kang, Zehui Xiong, Shiwen Mao, Zhu Han, Abbas Jamalipour, Dong In Kim, Xuemin Shen, Victor C. M. Leung, and H. Vincent Poor. Unleashing the power of edge-cloud generative ai in mobile networks: A survey of aigc services. *IEEE Communications Surveys & Tutorials*, 26(2):1127–1170, 2024.
- [226] Zheng Lin, Guanqiao Qu, Qiyuan Chen, Xianhao Chen, Zhe Chen, and Kaibin Huang. Pushing large language models to the 6g edge: Vision, challenges, and opportunities. *arXiv preprint arXiv:2309.16739*, 2023.
- [227] Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik Siu Chan, Xuhai Xu, Daniel McDuff, Cynthia Breazeal, and Hae Won Park. Adaptive collaboration strategy for llms in medical decision making. *arXiv preprint arXiv:2404.15155*, 2024.
- [228] Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming Li, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark Gerstein. MedAgents: Large language models as collaborators for zero-shot medical reasoning. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 599–621. Association for Computational Linguistics, August 2024.
- [229] Nimeesha Chan, Felix Parker, William Bennett, Tianyi Wu, Mung Yao Jia, James Fackler, and Kimia Ghobadi. MedtSllm: Leveraging llms for multimodal medical time series analysis. *arXiv preprint arXiv:2408.07773*, 2024.
- [230] Min Hu, Lei Liu, Xiaohua Wang, Yiming Tang, Jiaoyun Yang, and Ning An. Parallel multiscale bridge fusion network for audio–visual automatic depression assessment. *IEEE Transactions on Computational Social Systems*, 11(5):6830–6842, 2024.
- [231] Jian Chen, Yuzhu Hu, Qifeng Lai, Wei Wang, Junxin Chen, Han Liu, Gautam Srivastava, Ali Kashif Bashir, and Xiping Hu. Iifdd: Intra and inter-modal fusion for depression detection with multi-modal information from internet of medical things. *Information Fusion*, 102:102017, 2024.
- [232] Lang He, Mingyue Niu, Prayag Tiwari, Pekka Marttinen, Rui Su, Jiewei Jiang, Chenguang Guo, Hongyu Wang, Songtao Ding, Zhongmin Wang, Xiaoying Pan, and Wei Dang. Deep learning for depression recognition with audiovisual cues: A review. *Information Fusion*, 80:56–86, 2022.