

A Survey of Large Language Models for Healthcare: from Data, Technology, and Applications to Accountability and Ethics

Kai He, *Member, IEEE*, Rui Mao *Member, IEEE*, Qika Lin, Yucheng Ruan, Xiang Lan, Mengling Feng*, *Senior Member, IEEE* Erik Cambria, *Fellow, IEEE*

Abstract—The utilization of large language models (LLMs) in the Healthcare domain has generated both excitement and concern due to their ability to effectively respond to free-text queries with certain professional knowledge. This survey outlines the capabilities of the currently developed LLMs for Healthcare and explicates their development process, with the aim of providing an overview of the development roadmap from traditional Pretrained Language Models (PLMs) to LLMs. Specifically, we first explore the potential of LLMs to enhance the efficiency and effectiveness of various Healthcare applications highlighting both the strengths and limitations. Secondly, we conduct a comparison between the previous PLMs and the latest LLMs, as well as comparing various LLMs with each other. Then we summarize related Healthcare training data, training methods, optimization strategies, and usage. Finally, the unique concerns associated with deploying LLMs in Healthcare settings are investigated, particularly regarding fairness, accountability, transparency and ethics. Our survey provide a comprehensive investigation from perspectives of both computer science and Healthcare specialty. Besides the discussion about Healthcare concerns, we supports the computer science community by compiling a collection of open source resources, such as accessible datasets, the latest methodologies, code implementations, and evaluation benchmarks in the Github¹. Summarily, we contend that a significant paradigm shift is underway, transitioning from PLMs to LLMs. This shift encompasses a move from discriminative AI approaches to generative AI approaches, as well as a shift from model-centered methodologies to data-centered methodologies.

Index Terms—Large Language Model, Medicine, Healthcare Application

I. INTRODUCTION

Pretrained Language Models (PLMs) were primarily employed as a constituent part of Natural language Processing (NLP) systems [1]–[4], such as those used in Speech Recognition [5], [6], Metaphor Processing [7], Sentiment Analysis [8], [9], Information Extraction [10]–[12], and Machine Translation [13], [14]. However, with recent advancements, these PLMs are demonstrating an increasing capacity to function

* Corresponding author: Mengling Feng.

Kai He, Qika Lin, Yucheng Ruan, Xiang Lan and Mengling Feng are with the School of Saw Swee Hock School of Public Health and Institute of Data Science, National University of Singapore. E-mail: {kai_he, linqika, yuchengruan, ephlanx, ephfm}@nus.edu.sg.

Rui Mao and Erik Cambria are with the School of Computer Science and Engineering, Nanyang Technological University, Singapore, 639798. E-mail: {rui.mao, cambria}@ntu.edu.sg.

Manuscript received on September 30, 2023.

¹<https://github.com/KaiHe-better/LLM-for-Healthcare>

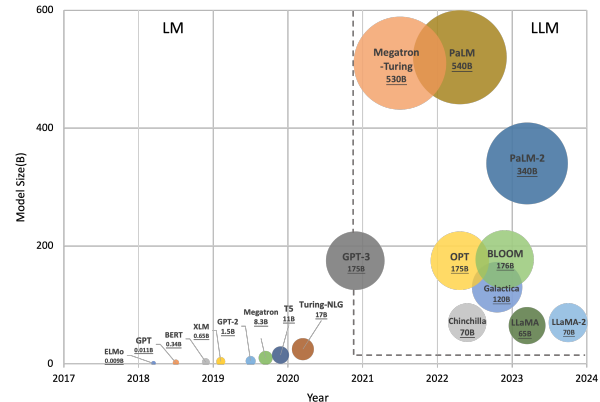


Fig. 1. The development from PLMs to LLMs. GPT-3 [17] marks a significant milestone in the transition from PLMs to LLMs, signaling the beginning of a new era.

as independent systems in their own right. Recently, OpenAI launched their Large Language Models (LLMs) ChatGPT and GPT-4, which shows superior performance in various NLP-related tasks, as well as scientific knowledge, such as Biology, Chemistry, and Medical exams [15]. Med-PaLM 2 [16] is Google’s LLMs, which are tailored to the medical domain. It is the first LLM that can achieve an “expert” level of performance on the MedQA dataset of US Medical Licensing Examination (USMLE²)-style questions, with an accuracy of over 85%.

A notable symbol of these advancements is the exponential growth in the sizes of PLMs. In the past five years, model sizes have increased by an astonishing 5,000 times, as depicted in Figure 1. Despite sharing many technical components, it is remarkable that simply scaling up these models leads to the emergence of novel behaviors, enabling qualitatively distinct capabilities [18]. In this context, the study [19] is relevant as it proposes a power-law positive relationship between model performance and three crucial factors: model size, dataset size, and amount of compute. We are currently witnessing a transition period where all three of these factors are skyrocketing, marking the evolution from PLMs to LLMs and opening up

²The United States Medical Licensing Examination (USMLE), a three-step examination program used to assess clinical competency and grant licensure in the United States. This is a main dataset for the evaluation of Healthcare LLMs.

previously inconceivable possibilities³.

Before LLMs, PLMs such as BERT [21] and RoBERTa [22] have gained a lot of attention. While early efforts have been made to develop PLMs-based neural networks for Healthcare [11], [23]–[27], these models are predominantly single-task systems that lack expressivity and interactive capabilities. This limits their usefulness for tasks like classification, regression, or segmentation in AI for Healthcare technologies [28], [29]. Additionally, PLMs face obstacles such as being difficult to explain, lacking adequate robustness, and requiring excessive amounts of data [30]–[32]. As a result, there is a disparity between what current models can accomplish and what is expected of them in real-world clinical workflows. However, recent advancements in LLMs have greatly improved these areas, facilitating deeper integration between LLMs and Healthcare. For instance, the emergent Chain-of-Thought (CoT) ability [33] provides the solution to the explainability challenge, while impressive few-shot even zero-shot ability [34] also alleviates the requirements of expensive medical annotations.

This paper considers GPT-3 [17] as a crucial milestone that signifies the start of the transition from PLMs to LLMs. GPT-3 is the first renowned LLM that has over 100 billion parameters, displays exceptional few-shot learning ability, and introduces in-context learning. Later, many other LLMs are proposed, including Megatron-LM [35], OPT [36], Chinchilla [37], Galactica [38], LLaMA [39], PaLM [40], and PaLM-2 [16]. These LLMs show distinguishing language understanding, generating ability, instruction following, reasoning ability and also common sense of the world [15], establishing them as fundamental models across diverse domains, including Finance [41], Education [42], and Healthcare [43].

These advanced improvements present a remarkable opportunity for LLMs to contribute significantly to Healthcare, such as the released LLMs of HuatuoGPT [44], Med-PaLM 2 [16], and Visual Med-Alpaca [45]. These studies have improved LLMs by tailoring them to the unique characteristics of the Healthcare field. For example, HuatuoGPT argues that, as an intelligent medical advice provider, LLMs should have the ability to actively ask questions for the patients rather than respond passively. Visual Med-Alpaca integrates with medical “visual experts” for multimodal biomedical tasks, enabling a wide range of tasks, from interpreting radiological images to addressing complex clinical inquiries.

Considering the immense potential of LLMs for Healthcare, we firmly believe that dedicating efforts to develop effective, ethical, and tailored LLMs for Healthcare is not only necessary but also imperative. Thus, this paper summarizes related studies in areas including algorithm development, potential healthcare applications, performance evaluation, as well as fairness, accountability, transparency, and ethics of LLMs. Limitations and future works are also discussed. Our goal is to update

readers on the latest developments in this field. Specifically, we recognize that different Healthcare scenarios require different capabilities from LLMs. For example, emotional comfort with patients needs more fluent conversations and empathy; hospital guide needs specific knowledge about the related buildings; and medical consultation needs more professional medical specialization. For computer science researchers, they require knowledge in selecting the appropriate LLM base model, suitable training data, and effective training strategies to find optimal solutions for diverse application scenarios. Our survey provide them with a comprehensive guidance for pursuit of achieving the best outcomes in their research endeavors. For medical researchers, we aspire for this survey to serve as a valuable resource aiding in the precise selection of LLMs aligning with their specific clinical requirements.

Comparing with existing studies about LLMs for Healthcare, they primarily concentrate on healthcare applications and often discuss the impacts without delving into the technical aspects of development and usage methods. In contrast, our survey represents the a comprehensive examination of LLMs specifically within the Healthcare domain, including detailed technology summarization, various Healthcare applications, and discussion about fairness, accountability, transparency, and ethics. For example, the surveys [46] only focus on medical or Healthcare applications of LLMs. They discussed the strengths and limitations of LLMs to improve the efficiency and effectiveness of clinical, educational and research work in medicine. The study [29] explores general applications of LLMs and emphasizes the potential future impacts they may have. The study [46] aims to support Healthcare practitioners in comprehending the rapidly evolving landscape of LLMs in the field of medicine, with a particular focus on highlighting both the potentials and pitfalls. However, they are not provided any detailed technological insights. Some former studies [47], [48] involved part of technological content, but they focus on general LLM developments and assessments [15] without specific adaptations and discussions for Healthcare. The studies of [30], [49] have focused on Healthcare PLMs rather than LLMs.

Besides our comprehensive investigation, the survey further analyze and summarize some development trends, including the current transition from PLMs to LLMs in the Healthcare domain. We provide a brief introduction to Healthcare PLMs as background information and then delve into the details of Healthcare LLMs, including technology details about how to develop and evaluate a private Healthcare LLM from scratch. Additionally, we analyze ethical concerns towards Healthcare LLMs, such as fairness, accountability, transparency, and ethics. Finally, we outline the distinct challenges that emerge when employing LLMs within the Healthcare domain. These challenges encompass augmenting medical knowledge, seamless integration of LLMs within healthcare procedures, interactions between patients and medical practitioners, and inherent issues associated with LLMs. Our contributions can be summarized as:

- We propose a comprehensive survey about LLMs for Healthcare. Our paper provides an overview of the development roadmap from PLMs to LLMs, which updates

³This paper defines LLMs as large models that after GPT-3, which include the ability to follow instructions and with at least 1 billion parameters. It is important to note that this definition is not rigid and serves as a way to distinguish recent LLMs from traditional PLMs studies, such as ELMo [20] and BERT [21]. The concept of PLMs is not a hypernym of LLM in our study. PLMs in our paper generally refers to the language model studies before LLMs appeared.

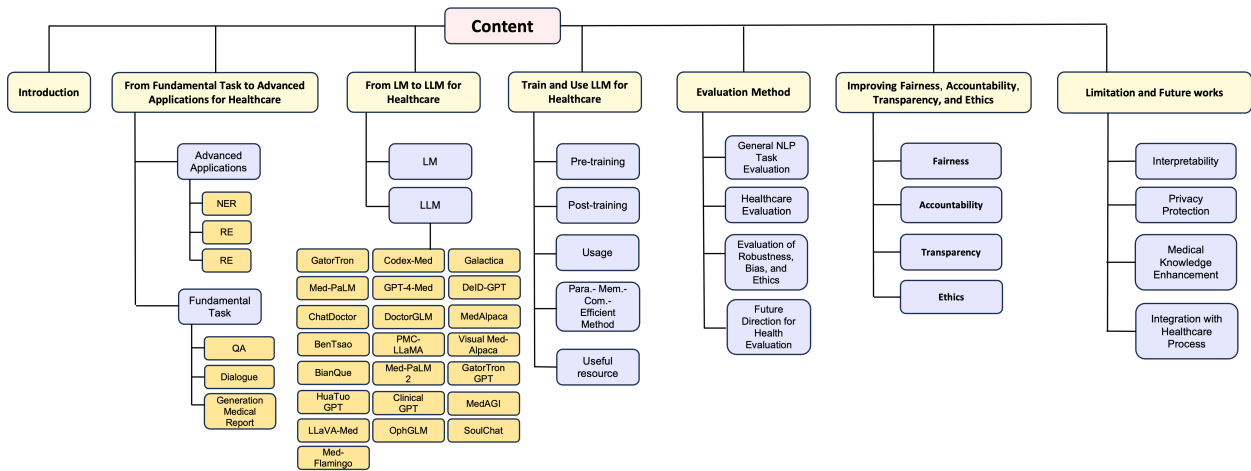


Fig. 2. The organizational framework for the content. Section III, Section IV, Section V are technology details, while Section II, Section VI and Section VII are more valued for Healthcare professionals.

readers on the latest advancements in this field.

- We have compiled a extensive list of publicly available data, training techniques, and evaluation systems for LLMs in Healthcare, which can be useful for those who plan to create their private Healthcare LLMs.
- We analyze numerous ethical considerations pertaining to the utilization of LLMs in the healthcare domain. These considerations encompass aspects such as robustness, toxicity, bias, fairness, accountability, transparency, ethics, as well as other constraints and prospective research areas. Our comprehensive analysis is anticipated to guide medical researchers in making informed choices when selecting LLMs suitable for their specific needs.

The overall structure of this paper is illustrated in Figure 2. Besides this Introduction section, Section II presents the applications of PLMs and LLMs in the Healthcare domain. Section III introduces and discusses the existing studies on PLMs and LLMs, highlighting their differences. The training and utilization of LLMs are described in Section IV. Evaluation methods for LLMs are discussed in Section V. Section VI focuses on the topics of fairness, accountability, transparency, and ethics specifically related to Healthcare LLMs. Lastly, Section VII provides the conclusion of the paper.

II. WHAT LLMs CAN DO FOR HEALTHCARE? FROM FUNDAMENTAL TASKS TO ADVANCED APPLICATIONS

Numerous endeavors have been made to apply PLMs or LLMs to Healthcare. In the early stages, the studies primarily focused on fundamental tasks, including medical Named Entity Recognition (NER), Relation Extraction (RE), Text Classification (TC), and Semantic Textual Similarity (STS), due to the challenges of accessing diverse medical datasets, the complexity of the medical domain, and limitations of the models’ capabilities [30]. Recently, the concept of Artificial General Intelligence (AGI) with Healthcare adaptation has been proposed [31], [50], which has led to more practical applications in various aspects of the Healthcare field. For instance, some online medical consultation systems [51], [52]

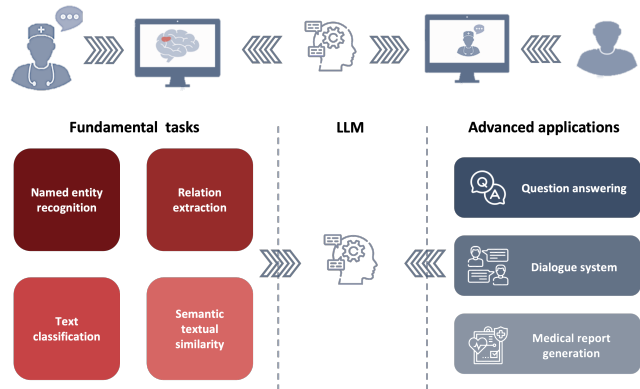


Fig. 3. LLMs for Healthcare: from fundamental task to advanced applications.

have been deployed, which can answer professional medical questions for patients and serve as guides in hospitals. Furthermore, some researchers explore the automatic generation of multimodal medical reports [53], [54]. The overall application framework of LLMs for Healthcare is shown in Figure 3. In the following sections, we analyze what LLMs can do for Healthcare in detail.

A. NER and RE for Healthcare

The initial step towards unlocking valuable information in unstructured Healthcare text data involves performing NER and RE. By extracting medical entities such as drugs, adverse drug reactions, proteins, chemicals, as well as predicting the relations between them, a multitude of useful functions can be achieved, including but not limited to Adverse Drug Event [55], Drug Drug Interaction [56], [57], and Chemistry Protein Reaction [58]. These two tasks also provide fundamental information for a range of other Healthcare applications, such as medical entity normalization and coreference [59], [60], medical knowledge base and knowledge graph construction [25], [61], and entity-enhanced dialogue [62], [63]. For example, by employing NER and RE tasks, the Healthcare

knowledge databases Drugbank⁴ [64] and UMLS [65] are constructed, which facilitate various applications in Intellectual Healthcare⁵ [66].

In the early stages of research on NER with PLMs, a significant portion of studies focused on sequence labeling tasks, as highlighted in previous research [32]. To accomplish this, PLMs-based approaches were employed to generate contextualized representations for individual tokens, coupled with a classification header such as a linear layer, BiLSTM, or CRF [67]–[69]. In the case of RE tasks, the extracted entity pairs' representations were typically fed into a classification header to determine the existence of relations between the given entities [10], [11], [70].

In the era of LLMs, NER and RE have been improved to work under more complex conditions and more convenient usages. One example is LLM-NERRE [71], which combines NER and RE to handle hierarchical information in scientific text. This approach has demonstrated the ability to effectively extract intricate scientific knowledge for tasks that require the use of LLMs. These tasks often involve complexities that cannot be effectively handled by typical PLMs such as BERT. Meanwhile, LLMs can finish medical NER and RE well even without further training. The study [72] employed InstructGPT [73] to perform zero- and few-shot information extraction from clinical text, despite not being trained specifically for the clinical domain. The results illustrated that InstructGPT can perform very well on biomedical evidence extraction [74], medication status extraction [75], and medication attribute extraction [75]. This observation supports the notion that LLMs can be applied with flexibility and efficiency, highlighting the adaptability, and showcasing their potential to contribute to advancements in Healthcare research and applications.

B. Text Classification for Healthcare

The aim of TC is to assign labels to text of different lengths, such as phrases, sentences, paragraphs, or documents. In Healthcare research, a large amount of patient data is collected in the electronic format, including disease status, medication history, lab tests, and treatment outcomes, which is a valuable source of information for analysis. However, these data can only be used with appropriate labels, while TC is one of the most commonly used technology. A research study [76] proposed several methods, based on hybrid Long Short-Term Memory (LSTM) and bidirectional gated recurrent units (Bi-GRU) to achieve medical TC. These methods were demonstrated effective in the Hallmarks dataset and AIM dataset [77] (Both these two datasets were sourced from biomedical publication abstracts). Another research study [78] used text classification to identify prescription medication mentioned in tweets and achieved good results using models

⁴Drugbank is a free and comprehensive online database that provides information on drugs and drug targets. The most recent version (5.0) includes 9591 drug entries, such as 2037 FDA-approved small molecule drugs, 241 FDA-approved biotech drugs, 96 nutraceuticals, and over 6000 experimental drugs.

⁵UMLS is a collection of controlled vocabularies used in biomedical sciences and Healthcare. It features a mapping structure that enables easy translation among different terminology systems, and serves as an extensive thesaurus and ontology of biomedical concepts.

like BERT, RoBERTa, XLNet, ALBERT, and DistillBERT with four proposed information fusion methods.

However, PLMs-based TC usually cannot satisfy explainable and reliable requirements in the Healthcare field, while LLMs-based TC mitigates these issues to some extent. For example, CARP [79] takes advantage of LLMs by introducing Clue And Reasoning Prompting to achieve better TC tasks. This study adopts a progressive reasoning strategy tailored to address the complex linguistic phenomena involved in TC. First, LLMs were prompted to find superficial clues like keywords, tones, and references. Then, a diagnostic reasoning process was induced for final decision-making. AMuLaP [80] is another example, which proposed Automatic Multi-Label Prompting for few-shot TC. By exploring automatic label selection, their method surpasses the GPT-3-style in-context learning method, showing significant improvements compared with previous PLMs-based results [81].

C. Semantic Textual Similarity for Healthcare

STS is a way to measure how much two phrases or sentences mean the same thing. In Healthcare, STS is often used to combine information from different sources, especially used for Electronic Health Records (EHR). The 2018 BioCreative/Open Health NLP (OHNLP) challenge [82] and the National NLP Clinical Challenges (n2c2) 2019 Track 1 show that STS can help reduce mistakes and disorganization in EHRs caused by copying and pasting or using templates. This means that STS can be used to check the quality of medical notes and make them more efficient for other NLP tasks [83]. The study [84] proposed a new method using ClinicalBERT, which was a fine-tuned BERT-based method. The proposed iterative multitask learning technique helps the model learn from related datasets and select the best ones for fine-tuning.

The study [85] applied pre-trained language models to the STS task and explored different fine-tuning and pooling strategies. They found that domain-specific fine-tuning has less impact on clinical STS than it does on general STS. The study [86] achieved the third-best performance on the STS task of 2019 N2C2, which demonstrated the efficiency of utilizing transformer-based models to measure semantic similarity for clinical text. GatorTron [87] is a clinical LLM, which formulated STS as a regression task without any fine-tuning. This LLM learned the sentence-level representations of the two pieces of text and adopted a linear regression layer to calculate the similarity score.

D. Question Answering for Healthcare

Traditionally, QA is a separate task that involves generating or retrieving answers for given questions. In Healthcare, QA can be very beneficial for medical professionals to find necessary information in clinical notes or literature, as well as providing basic Healthcare knowledge for patients. According to a report by the Pew Research Center [88], over one-third of American adults have searched online for medical conditions they may have. A strong QA system for Healthcare can significantly fulfill the consultation needs of patients.

Many studies [30], [49], [89] explored how to adapt general PLMs to answer Healthcare questions, including designing special pertaining task [90], fine-tuning on Healthcare data [91], and introducing external Healthcare knowledge base [92]. However, due to their limited language understanding and generation abilities [93], PLMs-based QA systems struggle to play a significant role in real-world Healthcare scenarios.

With the advent of powerful LLMs, prompt-based methods have been introduced to solve various tasks by formulating them as QA tasks, including NER [94], RE [10], and Sentiment Analysis [95]–[98]. In addition to these tasks, LLMs have significantly improved typical QA tasks in professional fields, such as Healthcare. For instance, Med-PaLM 2 [16], a medical domain LLM, achieved a score of up to 86.5% on the USMLE dataset, outperforming Med-PaLM [99] by over 19% and setting a new state-of-the-art. This LLM also approached or exceeded state-of-the-art performance across MedMCQA [100], PubMedQA [101], and MMLU clinical topics datasets [102]. In the study [103], the use of ChatGPT, Google Bard, and Claude for patient-specific QA from clinical notes was investigated. The accuracy, relevance, comprehensiveness, and coherence of the answers generated by each model were evaluated using a 5-point Likert scale on a set of patient-specific questions. Another study [104] proposed a retrieval-based medical QA system that leverages LLMs in combination with knowledge graphs to address the challenge.

E. Dialogue System for Healthcare

Chatbots have demonstrated promising potential to assist both patients and health professionals [105]–[107]. The implementation of Healthcare Dialogue Systems can decrease the administrative workload of medical personnel and mitigate the negative consequences resulting from a shortage of physicians [108]. Apart from the QA component, dialogue systems are generally classified into two categories: task-oriented and open-domain dialogue systems [109]. Task-oriented dialogue systems are designed to address specific issues for Healthcare, such as hospital guides or medication consultations. In contrast, open-domain dialogue systems prioritize conversing with patients without any specific tasks. These systems are usually used as chatbots to provide emotional support, or mental health-related applications [110], [111]. For example, the study of [112] shows that patients who participated in a telehealth project had lower scores for depression, anxiety, and stress, and experienced 38% fewer hospital admissions. However, this project adds to the workload of physicians who are already occupied with face-to-face medical practice. In addition to their existing responsibilities, they are required to provide remote telemedicine consultations, further increasing their workload. To maintain good results without overburdening physicians, automated dialogue systems are a promising technology for Healthcare.

In the early stages, the study of [113] proposed an ontology-based dialogue system that supports electronic referrals for breast cancer. This system can handle the informative responses of users based on the medical domain ontology.

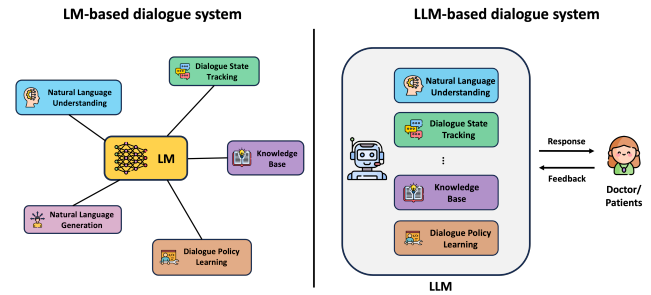


Fig. 4. The comparison between PLMs-based with LLMs-based dialogue system.

Another study KR-DS [114] is an end-to-end knowledge-routed relational dialogue system that seamlessly incorporates a rich medical knowledge graph into topic transitions in dialogue management. KR-DS includes a novel Knowledge-routed Deep Q-network (KR-DQN) to manage topic transitions, which integrates a relational refinement branch for encoding relations among different symptoms and symptom-disease pairs and a knowledge-routed graph branch for topic decision-making. In general, PLMs-based dialogue systems often comprise multiple sub-modules, like Nature Language Understanding, Dialogue Management, Nature Language Understanding, or Knowledge Introduction modules [109]. Each individual sub-module within the overall system has the potential to become a bottleneck, thereby restricting the system’s practical applications.

In the case of LLM-based dialogue systems, the original pipeline system can be transformed into an end-to-end system leveraging the capabilities of a powerful LLM [47], as shown in Figure 4. By utilizing an LLM, the remaining task involves aligning the system with human preferences and fine-tuning it for specific fields, without the need of many extra sub-modules, and achieving some advanced abilities that PLMs can hardly do. For example, a new approach [115] was proposed to detect depression, which involves an interpretable and interactive system based on LLMs. The proposed system not only provides a diagnosis, but also offers diagnostic evidence that is grounded in established diagnostic criteria. Additionally, users can engage in natural language dialogue with the system, which allows for a more personalized understanding of their mental state based on their social media content. Chatdoctor [116] is a specialized language model designed to overcome the limitations observed in the medical knowledge of prevalent LLMs like ChatGPT, by providing enhanced accuracy in medical advice. Chatdoctor adapted and refined LLaMA [39] using a large Healthcare dialogues dataset, and incorporating a self-directed information retrieval mechanism. This allows Chatdoctor to utilize real-time information from online sources to engage in conversations with patients. More LLMs for Healthcare can be seen in Section III-B.

F. Generation of Medical Reports from Images

Medical reports are of significant clinical value to radiologists and specialists, but the process of writing them can be tedious and time-consuming for experienced radiologists, and

error-prone for inexperienced ones. Therefore, the automatic generation of medical reports has emerged as a promising research direction in the field of Healthcare combined with AI. This capability can assist radiologists in clinical decision-making and reduce the burden of report writing by automatically drafting reports that describe both abnormalities and relevant normal findings, while also taking into account the patient’s history. Additionally, related models are expected to provide assistance to clinicians by pairing text reports with interactive visualizations, such as highlighting the region described by each phrase.

In an early stage, the study [117] proposed a data-driven neural network that combines a convolutional neural network with an LSTM to predict medical tags and generate a single sentence report, by employing a co-attention mechanism over visual and textual features. However, a single-sentence report is limited to real medical scenes. To generate multi-sentence reports, the study [118] proposed a multi-level recurrent generation model consisting of a topic-level LSTM and a word-level LSTM, and they also fused multiple image modalities by focusing on the front and later views.

Most recently proposed models for automated report generation rely on multimodal technology implemented by LLMs, which can support more advanced applications. For example, VisualGPT [119] utilizes linguistic knowledge from large language models and adapts it to new domains of image captioning in an efficient manner, even with small amounts of multimodal data. To balance the visual input and prior linguistic knowledge, VisualGPT employs a novel self-resurrecting encoder-decoder attention mechanism that enables the pre-trained language model to quickly adapt to a small amount of in-domain image-text data. ChatCAD [120] introduced LLMs into medical-image Computer Aided Diagnosis (CAD) networks. Their proposed framework leverages the capabilities of LLMs to enhance the output of multiple CAD networks, including diagnosis networks, lesion segmentation networks, and report generation networks, by summarizing and reorganizing information presented in natural language text format. Their results show that ChatCAD achieved significant improvements under various measures compared with the other two report-generation methods (R2GenCMN [121] and CvT2DistilGPT2 [122]). ChatCAD+ [123] is a multimodal system that addresses the writing style mismatch between radiologists and LLMs. The system is designed to be universal and reliable, capable of handling medical images from diverse domains and providing trustworthy medical advice by leveraging up-to-date information from reputable medical websites. ChatCAD+ also incorporates a template retrieval system that enhances report generation performance by utilizing exemplar reports, resulting in greater consistency with the expertise of human professionals. It should be noted that ChatCAD and ChatCAD+ are both integrated systems that utilize existing LLMs, rather than being LLMs themselves.

G. Summary

In addition to the conventional NLP tasks, LLMs play an integral role in specific sub-fields of Healthcare. One

notable example is the application of LLMs in advancing oncology research, where they contribute to scientific advancements and improve research efficiency. The studies [124]–[126] have emerged as the predominant learning paradigm in histopathology image analysis, offering valuable support for various tumor diagnosis tasks, including tumor detection, subtyping, staging, and grading. It is worth mentioning that these applications place significant emphasis on the multimodal capability of LLMs, as Healthcare data inherently consists of text, images, and time series data. By leveraging the strengths of LLMs, researchers and Healthcare professionals can harness the power of multiple modalities to improve diagnostic accuracy and patient care.

Apart from the aforementioned achievements, both general and Healthcare LLMs face several challenges that need to be addressed. These challenges encompass the effective structuring of high-quality data, the development of robust evaluation methods for assessing model output, and the seamless integration of LLMs into medical processes. For more detailed information on these challenges, please refer to Section VI and Section VII.

III. FROM PLMS TO LLMs FOR HEALTHCARE

Apart from the increasing model sizes, two significant developments from PLMs to LLMs are the transition from Discriminative AI to Generative AI and from model-centered to data-centered approaches.

During the PLMs period, published PLMs were primarily evaluated on Natural Language Understanding (NLU) tasks, such as mentioned NER, RE, and TC. These studies are grouped as discriminative AI, which concentrates on classification or regression tasks instead of generation tasks. In contrast, generative AI generates new content, often requiring the model to understand existing data (e.g., textual instructions) before generating new content. The evaluation tasks of generative AI are usually QA and conversation tasks.

The second perspective is the change from model-centered to data-centered. Before the rise of LLMs, previous research focused on improving neural architecture to enhance the encoding abilities of proposed models. As neural models became increasingly larger, the over-parameterization strategy [127] demonstrated promising abilities in learning potential patterns reserved in annotated datasets. Under such conditions, high-quality data played a more significant role in further enhancing various Healthcare applications [128], [129], namely, the transition from model-centered to data-centered direction. On the other hand, recent related developments present a multimodal trend, providing significant support to the data of EHRs, medical images, and medical sequence signals. Based on powerful LLMs, more existing and promising research and applications for Healthcare can be explored. Addressing the challenge of systematically collecting matched multimodal data holds significant importance. For such reason, we list detailed data usages and access links of each LLM in section III-B.

In the following sections, we first briefly introduce the focus of previous PLMs studies, and then more details about existing LLMs in the Healthcare field are provided. All content-related PLMs and LLMs are organized in chronological order.

TABLE I
BRIEF SUMMARIZATION OF EXISTING PLMS FOR HEALTHCARE.

Model Name	Base	Para. (B)	Features	Date	Link
BioBERT [91]	BERT	0.34	Biomedical Adaption	05/2019	GitHub
BlueBERT [130]	BERT	0.34	Biomedical Benchmark	06/2019	GitHub
MIMIC-BERT [131]	BERT	0.34	Clinical Concept Extraction	08/2019	-
BioFLAIR [132]	BERT	0.34	Less Computationally Intensive	08/2019	GitHub
Bio-ELECTRA-small [133]	ELECTRA	0.03	Training From Scratch	03/2020	-
AlphaBERT [134]	BERT	0.11	Character-level	04/2020	GitHub
Spanish-bert [135]	BERT	0.11	Spanish	04/2020	-
GreenCovidSQuADBERT [136]	BERT	0.34	CPU-only, CORP-19	04/2020	GitHub
BEHRT [137]	Transformer	-	Training From Scratch	04/2020	GitHub
BioMed-RoBERTa [138]	RoBERTa	0.11	Biomedical Adaption	05/2020	GitHub
RadBERT [139]	BERT	-	RadCore Radiology Reports	05/2020	-
CT-BERT [140]	BERT	0.34	COVID-19	05/2020	GitHub
French-BERT [141]	BERT	0.11	French Language Models	06/2020	-
FS-/RAD-/GER-BERT [142]	BERT	0.11	Chest Radiograph Reports	07/2020	GitHub
Japanese-BERT [143]	BERT	10.11	Japanese Clinical Narrative	07/2020	GitHub
MC-BERT [144]	BERT	0.11	Chinese Biomedical Benchmark	08/2020	GitHub
BioALBERT-ner [145]	ALBERT	0.18	Biomedical NER	09/2020	GitHub
BioMegatron [146]	Megatron	1.2	Training From Scratch	10/2020	GitHub
CharacterBERT [131]	BERT	0.11	Character-CNN module	10/2020	GitHub
ClinicalBERT [147]	BERT	0.11	For Predicting Hospital Readmission	11/2020	GitHub
Clinical XLNet [148]	XLNet	0.11	Temporal Information	11/2020	GitHub
Bio-LM [149]	RoBERTa	0.34	Biomedical Adaption	11/2020	GitHub
BioBERTpt [150]	BERT	0.11	Portuguese Clinical	11/2020	GitHub
RoBERTa-MIMIC [151]	RoBERTa	0.11	Clinical Concept Extraction	12/2020	GitHub
Clinical KB-ALBERT [152]	ALBERT	0.03	Introducing Medical KB	12/2020	GitHub
CHMBERT [153]	BERT	0.11	Chinese Medical, Cloud Computing	01/2021	-
PubMedBERT [154]	BERT	0.11	Training From Scratch	01/2021	Huggingface
ouBioBERT [155]	BERT	0.11	Up-sampling, Amplified Vocabulary	02/2021	GitHub
BERT-EHR [156]	BERT	0.11	Depression, Chronic Disease Prediction	03/2021	GitHub
ArabBERT [157]	BERT	0.11	Arabic Language	03/2021	GitHub
ABioNER [158]	BERT	0.11	Arabic NER	03/2021	-
ELECTRAMed [159]	ELECTRA	0.11	Biomedical Adaption	04/2021	GitHub
KcBioLM [160]	PubMedBERT	0.11	Introducing Medical KB	04/2021	GitHub
SINA-BERT [161]	BERT	0.11	Persian Language	04/2021	-
Med-BERT [162]	BERT	0.11	Slay Length Prediction	05/2021	GitHub
Galen [163]	RoBERTa	0.11	Spanish Language	05/2021	GitHub
SCIFIVE [164]	T5	0.77	Biomedical Text Generation	05/2021	GitHub
BioELECTRA [165]	ELECTRA	0.34	Training From Scratch	06/2021	GitHub
UnisBERT [152]	BERT	0.11	Introducing Medical KB	06/2021	GitHub
MedGPT [131]	GPT-2	1.5	Temporal Modelling	07/2021	-
MentalBERT [111]	BERT	0.11	Mental Healthcare	10/2021	huggingface
CODER [166]	mbERT	0.34	Cross-lingual, Introducing Medical KB	02/2022	GitHub
BioLinkBERT [167]	BERT	0.34	PubMed with Citation Links	03/2022	GitHub
BioALBERT [168]	ALBERT	0.03	Biomedical Adaption	04/2022	GitHub
BioBART [169]	BART	0.4	Biomedical NLP	04/2022	GitHub
SAPBERT [170]	BERT	0.11	Self-Alignment Pretraining	10/2022	GitHub
VPP [10]	BART	0.14	Soft prompt, Biomedical NER	03/2023	GitHub
KAD [171]	BERT	-	Multimodal, Chest Radiology Images	03/2023	GitHub

* For the column of Para. (B), only the largest size is listed.

TABLE II
SUMMARIZATION OF TRAINING DATA AND EVALUATION TASKS FOR EXISTING PLMS FOR HEALTHCARE.

Model Name	Method	Training Data	Eval task
BioBERT [91]	FT	PubMed, PMC	Biomedical NER, RE, QA
BlueBERT [130]	FT	PubMed, MIMIC-III	BLUE
MIMIC-BERT [131]	FT	MIMIC-III	Biomedical NER
BioFLAIR [132]	FT	PubMed	Bio NER
Bio-ELECTRA-small [133]	PT	PubMed	Biomedical NER
AlphaBERT [134]	FT	Discharge diagnoses	Extractive Summarization Task
Spanish-bert [135]	FT	Spanish	Spanish Clinical Case Corpus
GreenCovidSQuADBERT [136]	FT	CORP19, PubMed, PMC	NER, QA
BEHRT [137]	PT	CPRD, HES	Disease Prediction
BioMed-RoBERTa [138]	FT	BIOMED	CHEMPROT, RCT
RadBERT [139]	FT	Radiology Report Corpus	Report Coding, Summarization
CT-BERT [140]	FT	Tweet	COVID-19 Text Classification
French-BERT [141]	FT	French clinical documents	DEFT challenge
FS-/RAD-/GER-BERT [142]	FT,PT	Unstructured radiology reports	Chest Radiograph Reports Classification
Japanese-BERT [143]	FT	Japanese EHR	Symptoms Classification
MC-BERT [144]	FT	Chinese EHR	Chinese Biomedical Evaluation benchmark
BioALBERT-ner [145]	FT	PubMed, PMC	Biomedical NER
BioMegatron [146]	PT	PubMed	biomedical NER, RE, QA
CharacterBERT [131]	Bert	OpenWebText, MIMIC-III, PMC	Medical NER, NLI, RE, SS
ClinicalBERT [147]	FT	MIMIC-III	Hospital Readmission Prediction
Clinical XLNet [148]	FT	MIMIC-III	PMV, Mortality
Bio-LM [149]	FT	PubMed, PMC, MIMIC-III	18 Biomedical NLP Tasks
BioBERTpt [150]	FT	Private clinical notes, WMT16	SemClinBr
RoBERTa-MIMIC [151]	FT	i2b2 2010, 2012, n2c2 2018	i2b2 2010, 2012, N2C2 2018
Clinical KB-ALBERT [152]	FT	MIMIC-III, UMLS	MedNLI, i2b2 2010, 2012
CHMBERT [153]	FT	Medical text data	Disease Prediction
PubMedBERT [154]	PT	PubMed	BLURB
ouBioBERT [155]	FT	PubMed, Wikipedia	BLUE
BERT-EHR [156]	FT	General EHR	Myocardial Infarction, Breast Cancer, Liver Cirrhosis
ArabBERT [157]	PT	Arabic Wikipedia, OSIAN	Arabic SA, NER, QA
ABioNER [158]	FT	Arabic scientific literature	Arabic NER
ELECTRAMed [159]	FT	PubMed	Biomedical NER, RE, and QA
KcBioLM [160]	FT	PubMed	BLURB
SINA-BERT [161]	FT	Online Persian source	Persian QA, SA
Med-BERT [162]	FT	General EHR	Disease prediction
Galen [163]	FT	Private clinical cases	CodiEsp-D, CodiEsp-P, Catemist-Coding tasks
SCIFIVE [164]	T5	PubMed, PMC	Biomedical NER, RE, NLI, QA
BioELECTRA [165]	PT	PubMed, PMC	BLURB, BLUE
UnisBERT [152]	FT	MIMIC-III	6 BioNLP Tasks
MedGPT [131]	FT	MIMIC-III, private EHRs	Disorder Prediction
MentalBERT [111]	FT	Reddit	Depression Stress, Suicide Detection,
CODER [166]	FT	UMLS	MCSM, Medical RE
BioLinkBERT [167]	FT	PubMed	BLURB, USMLE
BioALBERT [168]	FT	PubMed, PMC, MIMIC-III	6 BioNLP Tasks
BioBART [169]	FT	PubMed	Biomedical EL, NER, QA, Dialogue, Summarization
SAPBERT [170]	FT	UMLS	MEL
VPP [10]	FT	PubMed	Biomedical NER
KAD [171]	FT	MIMIC-CXR	PadChest, ChestXray14, CheXpert and ChestX-Det10

☆ PMV means prolonged mechanical ventilation prediction. NER means Named Entity Recognition, NLI means Natural Language Inference, RE means Relation Extraction, SS means Sentence Similarity. MCSM means medical conceptual similarity measure [172]. MEL means medical entity linking. EL means Entity Linking. For clarity, we only list parts of representative evaluation tasks.

A. PLMs for Healthcare

While our survey primarily concentrates on LLMs for Healthcare, it is important to acknowledge that previous studies on PLMs have played a foundational role in the development of LLMs. In this section, we sum up the key research focus at a high level for PLMs, namely 1) enhancing neural architectures, and 2) utilizing more efficient pre-training tasks. These two points will be compared with the distinct study focus of LLMs in section III-B, to further support the transition from discriminative AI to generative AI and from model-centered to data-centered.

1. Improving Neural Architectures: in the early days of language modeling, task-specific models were the primary focus of neural architecture designs. PLMs, predominantly Word2Vec or GloVe, only a small part of parameters are used to generate static word embeddings in overall neural architecture. Following, the advent of ELMo introduced contextual embeddings and signaled a shift in neural architectures. Distinct from earlier models, ELMo’s word representations were contingent on entire sentence contexts, thereby allowing a word’s representation to dynamically change based on its context. At this stage, the language model generated dynamic word representations and these word representations hold nearly equal importance as the task-specific parameters in the design of neural architectures. With the Transformer architecture, BERT stood out as a game-changer in the neural architecture design. Almost all parameters are responsible for generating robust word representations. Subsequently, the PLMs training paradigm shifted from creating task-specific models to adopting a pre-train/fine-tune paradigm. Researchers started to place greater emphasis on improving PLMs that generate better word embeddings rather than crafting task-tailored models. This shift trend continues when we come to the LLM era, with more and more parameters for general and meaningful word representations.

2. More Efficient Pre-training Tasks: several previous studies [91], [131], [134] have demonstrated that pre-training can markedly bolster performance. This approach offers substantial benefits due to its capacity to enrich language understanding and improve model performance across various pretraining tasks [173]–[175]. Besides, there are specific tasks tailored for the Healthcare domain. The study [90] focused on adapting PLMs to the Healthcare field. The researchers introduced Biomedical Entity Masking as a technique to incorporate more medical knowledge into the model, thereby enhancing its performance in Healthcare-related applications.

Among the above pre-training tasks, masked language modeling and next-word prediction are the two most representative tasks, which also correspond to autoencoding and autoregressive PLMs. Autoencoding PLMs mask portions of the input and task the model with reconstructing the original sequence, thereby compelling the model to harness both left and right contextual information. Autoregressive PLMs are widely adopted in the GPT family, with a pre-training objective of predicting subsequent tokens in a sequence using preceding ones. In the PLMs era, autoencoding PLMs generally outperformed autoregressive PLMs. On the contrary, the majority

TABLE III
BRIEF SUMMARIZATION OF EXISTING LLMs FOR HEALTHCARE.

Model Name	Base	Para. (B)	Features	Date	Link
GatorTron [181]	Transformer	0.345, 3.9, 8.9	Training from scratch	06/2022	Github
Codex-Med [182]	GPT-3.5	175	CoT, Zero-shot	07/2022	Github
Galactica [38]	Transformer	1.3, 6.4, 30, 120	Reasoning, Multidisciplinary	11/2022	Org
Med-PaLM [99]	Flan-PaLM/PaLM	540	CoT, Self-consistency	12/2022	-
GPT-4-Med [183]	GPT-4	-	no specialized prompt crafting	03/2023	-
DeID-GPT [184]	GPT-4	-	De-identifying	03/2023	Github
ChatDoctor [116]	LLaMA	7	Retrieve online, external knowledge	03/2023	Github
DoctorGLM [185]	ChatGLM	6	Extra prompt designer	04/2023	Github
MedAlpaca [186]	LLaMA	7, 13	Adapt to Medicine	04/2023	Github
BenTsao [187]	LLaMA	7	Knowledge graph	04/2023	Github
PMC-LLaMA [188]	LLaMA	7	Adapt to Medicine	04/2023	Github
Visual Med-Alpaca [45]	LLaMA	7	multimodal generative model, Self-Instruct	04/2023	Github
BianQue [189]	ChatGLM	6	Chain of Questioning	04/2023	Github
Med-PaLM 2 [16]	PaLM 2	340	Ensemble refinement, CoT, Self-consistency	05/2023	-
GatorTronGPT [190]	GPT-3	5, 20	Training from scratch for medicine	05/2023	Github
HuatuogPT [44]	Bloomz	7	Reinforced learning from AI feedback	05/2023	Github
ClinicalGPT [191]	BLOOM	7	multi-round dialogue consultations	06/2023	-
MedAGI [192]	MiniGPT-4	-	multimodal, AGI	06/2023	Github
LLaVA-Med [193]	LLaVA	13	multimodal, self-instruct, curriculum learning	06/2023	Github
OphGLM [194]	ChatGLM	6	multimodal, Ophthalmology LLM	06/2023	Github
SoulChat [195]	ChatGLM	6	Mental Healthcare	06/2023	Github
Med-Flamingo [196]	Flamingo	80B	multimodal, Few-Shot generative medical VQA	07/2023	Github

of LLMs predominantly utilize the autoregressive approach, which has proven to be more effective.

For Healthcare PLMs, as observed in Tables I and II, a majority of the models utilize the discriminative approach, predominantly built upon the BERT architecture. The rationale behind this architectural choice is evident: many typical Healthcare applications are classification tasks. These tasks range from NER in the biomedical domain to more specific challenges such as disease prediction and relation extraction. In addition, the methodology of fine-tuning (FT) stands out as the prevalent training methodology. This trend suggests a broader implication: while general pre-trained models offer a foundational grasp of language, they require refinement through domain-specific data to excel in the applications of Healthcare. The choice of training datasets provides further support to the models' intent of achieving a holistic understanding of the medical domain.

Unlike recent LLMs, LLMs have the advantage of eliminating the need for FT and can directly infer at various downstream tasks. Moreover, the core research focus does not primarily revolve around improving neural architectures and developing more efficient pre-training tasks.

B. LLMs for Healthcare

With the surge in general LLM research [47], [48], there has also been a notable development of LLMs specifically tailored for the Healthcare field. In contrast to the emphasis on neural architecture designs [176], [177], pretraining tasks [178], and training strategies [179], [180] in previous PLMs research, the studies on LLMs for Healthcare greater emphasis on the collection of diverse, precise, and professional Healthcare data, and also data security and privacy protection.

In the following sections, we present an overview and analysis of the published LLMs designed for Healthcare. For the sake of convenience, we have compiled the pertinent information in Table III and Table IV, facilitating easy comparisons.

1) *GatorTron*: GatorTron [181], an early LLM developed for the Healthcare domain, aims to investigate how systems utilizing unstructured EHRs can benefit from clinical LLMs with billions of parameters. This LLM is trained from scratch, utilizing over 90 billion tokens, including more than 82 billion words of de-identified clinical text. The GatorTron-base model consists of 24 transformer blocks, similar to the architecture of the BERT large model. The GatorTron-medium model has been scaled up to 3.9 billion parameters (10 times the base setting), and the GatorTron-large model has been scaled up to 8.9 billion parameters, similar to BioMegatron [146] (which has 8.3 billion parameters). After training, GatorTron was systematically evaluated on five clinical NLP tasks, including clinical concept extraction, medical RE, Semantic Textual Similarity (STS), medical Natural Language Inference (NLI), and medical QA.

GatorTron's performance on various clinical NLP tasks has been evaluated. For clinical concept extraction, GatorTron was tested on i2b2 2010 [197], i2b2 2012 [198], and n2c2 2018 [199], achieving F1 measures of 89.96%, 80.91%, and 90.00%, respectively. Regarding medical Relation Extraction (RE), GatorTron-large was tested on n2c2 2018 [199] and achieved an F1 measure of 96.27%. For Semantic Textual Similarity (STS) and medical Natural Language Inference (NLI), GatorTron-large achieved 88.96% and 90.20% Pearson correlation and Accuracy, respectively, in the n2c2 2019 dataset [200] and MedNLI [201]. Regarding medical QA [202], GatorTron-large attained 74.08% and 97.19% on the emrQA Medication and emrQA Relation tasks.

In summary, GatorTron was an early attempt to investigate the impact of increasing LLM size on Healthcare tasks, which follows that the Megatron-Turing NLG model [203] was scaled up to 530 billion parameters, and the GPT-3 [17] model was developed with 175 billion parameters for general domain tasks. The results obtained from GatorTron demonstrated significant improvements for sentence-level and document-level NLP tasks, such as Medical NLI and Medical QA, but

TABLE IV
SUMMARIZATION OF TRAINING DATA AND EVALUATION TASKS FOR EXISTING LLMs FOR HEALTHCARE.

Model Name	Method	Training Data	Eval datasets
GatorTron [181]	PT	Clinical notes	CNER, MRE, MQA
Codex-Med [182]*	ICL	-	USMLE, MedMCQA, PubMedQA
Galactica [38]	PT, IFT	DNA sequence	MedMCQA, PubMedQA, Medical Genetics
Med-PaLM [99]	IPT	Medical data	MultiMedQA, HealthSearchQA
GPT-4-Med [183]*	ICL	-	USMLE, MultiMedQA
DeID-GPT [184]*	ICL	-	i2b2/UTHealth de-identification task
ChatDoctor [116]	IFT	Patient-doctor dialogues	iCliniq
DoctorGLM [185]	IFT	Chinese medical dialogues	-
MedAlpaca [186]	IFT	Medical dialogues and QA	USMLE, Medical Meadow
BenTsao [187]	IFT	Medical knowledge graph, Medical QA	Customed medical QA
PMC-LLaMA [188]	IFT	Biomedical academic papers	PubMedQA, MedMCQA, USMLE
Visual Med-Alpaca [45]	PT, IFT	medical QA	-
BianQue [189]	IFT	medical QA	-
Med-PaLM 2 [16]	IFT	-	MultiMedQA, Long-form QA
GatorTronGPT [190]	PT	Clinical and general text	PubMedQA, USMLE, MedMCQA, DDI, BC5CDR, KD-DTI
HuatuogPT [44]	IFT	Instruction and Conversation Data	CmedQA, webmedQA, and Huatuo26M
ClinicalGPT [191]	IFT+RLHF	Medical dialogues and QA, EHR	MedDialog, MEDQA-MCMLE, MD-EHR, cMedQA2
MedAGI [192]	IFT	Public medical datasets and images	SkinGPT-4, XrayChat, PathologyChat
LLaVA-Med [193]	IFT	multimodal biomedical instruction	VQA-RAD, SLAKE, PathVQA
OphGLM [194]	IFT	Knowledge graphs, medical dialogues	Fundus diagnosis pipeline tasks [194]
SoulChat [195]	IFT	Long text, empathetic dialogue	-
Med-Flamingo [196]	IFT	Image-caption/tokens pairs	VQA-RAD, Path-VQA, Visual USMLE

✧ * means the study focuses on evaluating the Healthcare LLM, rather than proposing a new LLM. PT means pre-training, ICL means In-context-learning (no parameters updated), IFT means instruction fine-tuning, and IPT means instruction prompt tuning. IPT comes from [99]. It should be noted that this concept is slightly different from Instruction fine-tuning or supervised fine-tuning.

TABLE V
DESIGNED CoT PROMPTS FOR HEALTHCARE QA.

#1 – Let’s think step by step
#2 – Let’s think step by step like a medical expert
#3 – Let’s use step-by-step inductive reasoning, given the medical nature of the question
#4 – Let’s differentiate using step-by-step reasoning like a medical expert
#5 – Let’s derive the differential diagnosis

only moderate improvements for phrase-level tasks, such as Clinical NER and Medical RE. According to this, their results show that larger transformer models are more beneficial for sentence-level and document-level NLP tasks.

2) *Codex-Med*: Codex-Med [182] aimed to investigate the effectiveness of GPT-3.5 models. Specifically, the performance of Codex [204] and InstructGPT [73] was investigated in their ability to answer and reason about real-world medical questions. To evaluate their effectiveness, two multiple-choice medical exam question datasets, namely USMLE [205] and MedMCQA [100], as well as a medical reading comprehension dataset called PubMedQA [101] were utilized. These datasets served as benchmarks to assess the language models’ comprehension and accuracy in addressing medical-related queries.

The study also explored three different prompting scenarios, namely CoT, in-context-learning (ICL, adding question-answer exemplars), and retrieval augmentation (injecting Wikipedia passages into the prompt). Additionally, the study investigated how scaling inference-time computing enabled Codex 5-shot CoT to be calibrated and achieve human-level performance on the three medical datasets.

According to previous studies [102], human experts achieved performance of 87.0%, 90.0%, and 78.0% accuracy on the USMLE, MedMCQA, and PubMedQA datasets, respectively. The Codex-Med study found that Codex (code-davinci-002) 5-shot with CoT achieved 60.2%, 62.7%, and 78.2% accuracy on the same datasets, while SOTA results (at that

time) after fine-tuning were 50.32%, 52.93%, and 78.20%. The study also designed five CoT prompts, as shown in Table V, which improved InstructGPT’s accuracy by 0.7%, 2.2%, and 3.5% on the USMLE, MedMCQA, and PubMedQA datasets, respectively. Furthermore, the study conducted an error analysis of CoT results and found that most of the incorrectly answered questions were due to CoTs containing reasoning errors (86%) or a lack of knowledge (74%). Misunderstanding the questions or context was less frequently observed (50%).

Generally, the main focus of Codex-Med was to investigate the efficacy of GPT 3.5 for Healthcare QA tasks using zero/few-shot learning and CoT prompting, without new LLMs or related technology proposed. The study revealed that the general LLM can significantly outperform fine-tuned BERT baselines for Healthcare QA tasks. In addition to the robust results, the study also identified a form of bias where the ordering of answer options affects predictions. However, the study acknowledged that many other biases, such as those related to gender or race, may also impact predictions, including those hidden in the training data.

3) *Galactica*: Aiming to solve the problem of information overload in the scientific field, Galactica was proposed to store, combine, and reason about scientific knowledge, including Healthcare. Galactica [38] was trained on a large corpus of papers, reference material, and knowledge bases to potentially discover hidden connections between different research and bring insights to the surface. Unlike other PLMs and LLMs, which rely on an un-curated crawl-based paradigm, Galactica’s training data consists of 106 billion tokens from high-quality sources, such as papers, reference material, and encyclopedias. This allows for the exploration of purposefully designed LLMs with a clear understanding of what enters the corpus, similar to expert systems that have normative standards. Galactica is built on a Transformer architecture in a decoder-only setup, utilizing the GeLU Activation [206], a 2048-length context

window, no biases in any of the dense kernels or layer norms, and Learned Positional Embeddings [207]. The study proposed five versions of Galactica, namely Galactica 125M, 1.3B, 6.7B, 30B, and 120B, which were tested on various scientific tasks, with a particular emphasis on evaluating Healthcare-related tasks.

Galactica demonstrated impressive results on various Healthcare-related tasks. Specifically, on PubMedQA [101], Galactica achieved a score of 77.6%, surpassing the state-of-the-art result of 72.2% [167]. On MedMCQA dev [100], Galactica achieved a score of 52.9% compared to the state-of-the-art result of 41.0% [154]. Furthermore, on BioASQ [208] and USMLE [205], Galactica's performance was close to the state-of-the-art results achieved by fine-tuned models (94.8% and 44.6%) [167].

Galactica emphasizes the importance of dataset design for LLMs. In response to this, the study curated a high-quality dataset and engineered an interface to interact with the body of knowledge. As a result, Galactica performs exceptionally well in knowledge-intensive scientific tasks, achieving promising results on PubMedQA and MedMCQA.

4) *Med-PaLM*: Med-PaLM [99] is a variant of PaLM [40] by employing instruction prompt tuning. Instruction prompt tuning is a parameter-efficient approach for aligning LLMs to new domains using a few exemplars proposed in the study [99]. Instead of using a hard prompt that is specific to each medical dataset, instruction prompt tuning used in this study employs a soft prompt as an initial prefix that is shared across multiple datasets. The soft prompt is then followed by a task-specific human-engineered prompt that includes instructions and/or few-shot exemplars, which may include CoT examples, along with the actual question and/or context. The authors of the study contend that current medical question answering benchmarks [205] are restricted to evaluating classification accuracy or automated natural language generation metrics (e.g., BLUE [209]), and do not allow for the thorough analysis necessary for real-world clinical applications. For such reason, they proposed MultiMedQA benchmark, consisting of LiveQA TREC 2017 [210], MedicationQA [211], PubMedQA [101], MMLU [102], MedMCQA [100], USMLE [205] and HealthSearchQA [99].

The original study [40] introduced PaLM, a densely-activated decoder-only transformer language model, which was trained using Pathways [212], a large-scale ML accelerator orchestration system that enables efficient training across TPU pods. The PaLM training corpus comprises 780 billion tokens, including a mix of web pages, Wikipedia articles, source code, social media conversations, news articles, and books. Further, this study [99] utilized instruction-tuned [213] to create Flan-PaLM, which was then fine-tuned using instruction prompt tuning to align it more closely with the medical domain, resulting in Med-PaLM. In the study [99], Flan-PaLM was evaluated on MedMCQA, USMLE, and PubMedQA, resulting in scores of 57.6%, 67.6%, and 79.0%, respectively.

It should be noted that this study also proposed a framework for human evaluation, which consists of 12 aspects, including Scientific consensus, Extent of possible harm, Likelihood of possible harm, Evidence of correct comprehension, Evidence

of correct retrieval, Evidence of correct reasoning, Evidence of incorrect comprehension, Evidence of incorrect retrieval, Evidence of incorrect reasoning, Inappropriate/incorrect content, Missing content, and Possibility of bias. For evaluating Med-PaLM [99], clinicians were asked to rate answers provided to questions in the HealthSearchQA, Live QA and Medication question answering datasets. Following the proposed human evaluation framework, Flan-PaLM, Med-PaLM, and clinicians achieved 61.9%, 92.6%, and 92.9% consensus, respectively. They argued that human evaluation reveals important limitations of today's models, reinforcing the importance of both evaluation frameworks and method development in creating safe, helpful LLM models for clinical applications.

5) *GPT-4-Med*: In the study [183], the authors provide a thorough evaluation of GPT-4, a state-of-the-art LLM, on medical competency examinations and benchmark datasets. Despite not being specifically trained or engineered for clinical tasks, GPT-4 is a general-purpose model that was analyzed on two sets of official practice materials for USMLE. The evaluation also included the MultiMedQA suite of benchmark datasets to test performance on various aspects of medical knowledge and reasoning. The results of the evaluation demonstrated that GPT-4, even in a zero-shot setting, significantly outperformed earlier models, achieving an average score of 86.65% and 86.7% on the Self-Assessment and Sample Exam of the USMLE tests, respectively. This is compared to the scores of 53.61% and 58.78% obtained by GPT-3.5. As the details of GPT-4 are not publicly available, we will not discuss them in detail.

6) *DeID-GPT*: The digitization of Healthcare has allowed for the sharing and reuse of medical data, but it has also raised concerns regarding confidentiality and privacy. As a result, there is a pressing need for effective and efficient solutions for de-identifying medical data, particularly in free-text formats. In the study [184], the authors developed a novel de-identification framework called DeID-GPT, which utilizes GPT-4 to automatically identify and remove identifying information. Compared to existing medical text data de-identification methods, DeID-GPT demonstrated the highest accuracy and remarkable reliability in masking private information from unstructured medical text while preserving the original structure and meaning of the text. This study is among the first to utilize ChatGPT and GPT-4 for medical text data processing and de-identification, providing insights for further research and solution development on the use of LLMs such as ChatGPT/GPT-4 in Healthcare. However, as with GPT-4-Med, we cannot discuss the details for this study.

7) *ChatDoctor*: The primary goal of the ChatDoctor [116] was to address the limitations of existing LLMs, including ChatGPT, in terms of their medical knowledge accuracy. To achieve this, the project first utilized a generic conversation model LLaMA and trained it using 52,000 instruction-following data from Stanford University's Alpaca [214]. Subsequently, the project collected a dataset of 100,000 patient-physician conversations (HealthcareMagic-100k) from an online medical consultation website (www.Healthcaremagic.com). The LLaMA model was initially fine-tuned with Alpaca's data to acquire basic con-

versation skills. Then, the model was further refined using the HealthcareMagic-100k dataset to improve its medical knowledge accuracy.

However, ChatDoctor didn't provide enough evaluations, except for some QA examples. For such reason, we cannot make further discussion and analysis for ChatDoctor.

8) *DoctorGLM*: DoctorGLM [185] is a Chinese LLM for Healthcare, which represents an effort to expand the use of LLMs beyond the English language and to explore a viable and affordable pipeline for creating customized medical LLMs. To achieve this, DoctorGLM was fine-tuned using ChatGLM-6B [215], a bilingual model capable of proficiently processing both English and Chinese. The most significant reason why choosing the GLM model as the base model is due to the unique scaling property, which enables INT4 quantization and effective inference on a single RTX 3060 (12G), making it more efficient and cost-effective for hospitals to deploy their medical dialogue models based on their in-house data. This breakthrough in Healthcare language modeling has significant implications for improving the efficiency and affordability of medical dialogue models.

DoctorGLM was trained using a database of medical dialogues in Chinese, which was derived from the ChatDoctor [116] dataset by utilizing the ChatGPT API for translation. To facilitate the fine-tuning process on an A100 80G GPU, they employed the LoRA technique [216] that resulted in faster inference times, making it easier for researchers and developers to utilize LLMs. During the inference stage, the model uses a prompt designer module to pre-process the user's input. This module extracts relevant keywords, such as the name of the disease or symptoms, from the user's input and generates a brief description based on a professional disease knowledge library containing 3,231 detailed disease documents. However, it should be noted that this extra information may also mislead the LLMs, because the descriptions from patients usually are non-professional and can be imprecise. If inputs contain too much context relying on patient statements, the model may overlook other potential factors. While they did not provide detailed information on their evaluation, except an accuracy of 67.6% on the USMLE (without any specific citation provided in their original paper).

Furthermore, DoctorGLM disclosed its computation cost. Their training process can handle roughly 80,000 single question-answer pairs per hour per GPU. Assuming that three epochs are required, and the cloud computing server of an A100 GPU costs about 5 USD per hour, the total training time needed is 3.75 hours, which amounts to a cost of approximately 18.75 USD to fine-tune DoctorGLM on 100,000 QA pairs. In terms of inference, DoctorGLM only needs around 13 GB of GPU memory, and it can be performed on a consumer-level GPU like an RTX 3090. This implies a total cost (inference PC) of about 1500 USD. This information will be very helpful for the people who plan to estimate their training costs.

9) *MedAlpaca*: Different from the general domain, Healthcare data inherently has sensitive and imperative needs for privacy safeguards. For such reason, non-transparent models with unclear data management practices are ill-suited for

medical applications. To tackle this challenge and avert unauthorized data transfers, MedAlpaca [186] employed an open-source policy that enables on-site implementation, aiming at mitigating privacy concerns.

MedAlpaca is built upon the LLaMA [39] with 7 and 13 billion parameters. They present Medical Meadow, a collection of medical tasks that are compiled for fine-tuning and evaluating the performance of LLMs in the context of medicine. Generally, Medical Meadow consists of two main categories, including Instruction Fine-Tuning formats and generally crawled Internet text. More details about training data can be seen in Section IV-E3. As for training strategy, MedAlpaca also implemented LoRA [216] for weight updates to adapt the LLM to specific tasks. Besides, they employed 8-bit matrix multiplication for the feed-forward and attention projection layers [217], along with an 8-bit optimizer [218] to further reduce the memory requirements.

MedAlpaca's performance was evaluated in a zero-shot setting across the USMLE Step 1, Step 2, and Step 3 self-assessment datasets, achieving accuracies of 47.3%, 47.7%, and 60.2%, respectively. However, with the application of LoRA and model quantization, the impact on MedAlpaca's performance was evident. The accuracy of Step 1, Step 2, and Step 3 decreased to 25.0%, 25.5%, and 25.5% for MedAlpaca-13b-LoRA, respectively. Additionally, for MedAlpaca-13b-LoRA-8bit, the accuracy further declined to 18.9%, 30.3%, and 28.9%.

10) *BenTsao*: BenTsao [187] (formerly known as HuaTuo, with the name change on May 12, 2023) is a LLaMA-based LLM that has been supervised-fine-tuned using generated QA instances. The model places a strong emphasis on ensuring the accuracy of facts in its responses, which is crucial in the biomedical domain. To accomplish this objective, two types of medical knowledge were utilized in constructing BenTsao: (1) structured medical knowledge such as medical knowledge graphs, and (2) unstructured medical knowledge such as medical guidelines. For medical knowledge graphs, they gathered diverse instructional data from CMeKG, a Chinese medical knowledge graph.

In terms of evaluation, they introduced a novel metric called SUS, which considers Safety, Usability, and Smoothness in evaluating PLMs in the biomedical domain. The SUS scale ranges from 1 (not acceptable) to 3 (good), with a score of 2 indicating an acceptable response. Five annotators with medical backgrounds evaluated the randomly mixed responses of the models using SUS. For Safety, Usability, and Smoothness, LLaMA received scores of 2.93, 1.21, and 1.58; Alpaca received scores of 2.64, 2.05, and 2.30; ChatGLM received scores of 2.59, 1.93, and 2.41, while BenTsao received scores of 2.88, 2.12, and 2.47, respectively.

11) *PMC-LLaMA*: PMC-LLaMA [188] is an open-source language model that by tuning LLaMA-7B on a total of 4.8 million biomedical academic papers for further injecting medical knowledge, enhancing its capability in the medical domain. PMC-LLaMA starts with the S2ORC [219] Datasets with 81.1M English-language academic papers and filters them with PubMed Central (PMC)-id. As a result, there are around 4.9M papers left, that are highly related to medical knowledge

totaling over 75B tokens.

Preliminary evaluations of PMC-LLaMA were conducted on three Healthcare QA datasets, namely PubMedQA, MedMCQA, and USMLE. According to their reports, PMC-LLaMA-7B achieved accuracies of 44.70%, 50.54%, and 69.5% on the USMLE test set, MedMCQA, and PubMedQA, respectively. However, when LoRA was applied, the accuracies decreased to 30.64%, 34.33%, and 68.20% on USMLE, MedMCQA, and PubMedQA, respectively.

12) *Visual Med-Alpaca*: Visual Med-Alpaca [45] is an open-source biomedical foundation model that originates from the University of Cambridge. It is designed to efficiently handle multimodal biomedical tasks by integrating with medical “visual experts”. The model is built upon the LLaMa-7B architecture [39] and trained using a collaboratively curated instruction set comprising contributions from both the GPT-3.5-Turbo language model and human experts. By incorporating plug-and-play visual modules and undergoing a few hours of instruction-tuning, Visual Med-Alpaca demonstrates versatility in performing various tasks, including the interpretation of radiological images and addressing complex clinical inquiries. Moreover, the model can be easily replicated as it requires only a single consumer GPU.

The biomedical instruction set for Visual Med-Alpaca was created through a multi-step process. Initially, medical questions were extracted from diverse medical datasets sourced from the BigBIO repository [220]. In order to enhance the dataset’s diversity and comprehensiveness, a self-instruct approach was adopted within the biomedical domain. This involved collecting inquiries from various medical question-and-answer datasets, namely MEDIQA RQE, MedQA, MedDialog, MEDIQA QA, and PubMedQA. These inquiries were used to prompt GPT-3.5-Turbo to generate corresponding answers. To ensure the quality of the instruction set, multiple rounds of human filtering and editing were conducted, resulting in a final dataset comprising 54,000 high-quality question-answer pairs.

13) *BianQue*: BianQue 1.0 [189] is a Chinese LLM for Healthcare, fine-tuned by a combination of instructions and multiple rounds of questioning dialog. In the medical field, it has been observed that doctors often require multiple rounds of questioning to make informed decisions, as opposed to a simple “command-and-response” model. Patients may not initially provide complete information during consultations, necessitating doctors to ask further questions before reaching a diagnosis and providing appropriate recommendations. In light of this, BianQue-1.0 was proposed, aiming to enhance the questioning capability of AI systems to simulate the consultation process followed by doctors. We define this capability as the “questioning” aspect of the “observation, sniffing, questioning, and cutting” process.

Considering the existing Chinese language model architecture, parameter count, and computational requirements, BianQue-1.0 based ClueAI/ChatYuan-large-v2 [221] as the baseline model. They fine-tuned the model for 1 epoch using eight NVIDIA RTX 4090 graphics cards, resulting in BianQue-1.0. For training, They created a hybrid dataset comprising Chinese medical QA commands and multi-turn

dialogues, named BianQueCorpus. BianQueCorpus merged various existing open-source Chinese medical QA datasets, including MedDialog-CN [222], IMCS-V2 [223], CHIP-MDCFNPC [224], MedDG [225], cMedQA2 [226], and Chinese-medical-dialogue-data⁶. This mixed dataset consisted of over 9 million samples and required approximately 16 days to complete one epoch of training. By combining these datasets, BianQue was able to examine the characteristics of both single-round and multiple-round interactions, as well as the questioning patterns employed by doctors.

Based on BianQueCorpus, BianQue 1.0 was updated to BianQue 2.0 on 6 June 2023. BianQue 2.0 chose ChatGLM-6B as the initialization model, and employed Instruction Fine-Tuning training of the full amount of parameters. Different from the BianQue-1.0 model, BianQue-2.0 expands the data such as drug instruction instruction, medical encyclopedic knowledge instruction, and ChatGPT distillation instruction, which strengthens the model’s suggestion and knowledge query capability

14) *Med-PaLM 2*: Med-PaLM 2 [16] is an updated version of LLMs for Healthcare, building upon Google’s Med-PaLM. It incorporates domain-specific medical Instruction Fine-Tuning, similar to how Med-PaLM is built upon PaLM using medical-based instructions [40]. It is worth noting that Med-PaLM 2 is based on PaLM 2 [227], which is a “smaller” PLM with 340B parameters. In contrast, PaLM, despite being part of the same series, is a “larger” model with 540B parameters. Interestingly, this represents a rare exception where a “smaller” LLM outperforms a “larger” LLM within the same series [16].

Med-PaLM 2 was evaluated on multiple-choice QA, including MedQA [205], MedMCQA [100], PubMedQA [101] and MMLU clinical topics [102] datasets, and long-form questions sampled from MultiMedQA [99]. For multiple-choice QA datasets, Med-PaLM 2 scored up to 86.5% on the USMLE dataset, compared with 67.2% from the Med-PaLM model. For MedMCQA and PubMedQA, Med-PaLM 2 achieved 72.3% and 75.0% accuracy, compared with 73.7% and 80.4% from GPT-4-base 5-shot. Med-PaLM 2’s long-form answers are evaluated by physicians and laypeople based on criteria including alignment with medical consensus, reading comprehension, knowledge recall, reasoning, inclusion of irrelevant content, omission of important information, potential for demographic bias, possible harm extent, and possible harm likelihood. Med-PaLM 2’s answers are often preferred over answers from physicians and the original Med-PaLM model. Besides, Med-PaLM 2 uses ensemble refinement as a new prompting strategy. This involves generating multiple reasoning paths and conditioning them to refine the final answer.

15) *GatorTronGPT*: GatorTronGPT [190] is a clinical generative LLM designed with a GPT-3 architecture comprising 5 or 20 billion parameters. It utilizes a vast corpus of 277 billion words, consisting of a combination of clinical and English text. The training data used for GatorTronGPT comprises de-identified clinical text sourced from the University of

⁶<https://github.com/Toyhom/Chinese-medical-dialogue-data>

Florida (UF) Health, along with 195 billion diverse English words obtained from the Pile dataset [87], [228]. Notably, GatorTronGPT was trained from scratch using the GPT-3 architecture. The study aimed to explore how the text generation capabilities of GatorTronGPT can contribute to medical research and Healthcare advancement.

GatorTronGPT underwent evaluation on two Healthcare-related tasks: biomedical RE and QA. In the biomedical RE task, the datasets DDI [229], BC5CDR [230], and KD-DTI [231] were utilized. GatorTronGPT achieved F1-measure scores of 50%, 49.4%, and 41.9% on these datasets. Regarding QA, GatorTronGPT (20B) attained accuracy scores of 77.6%, 45.1%, and 42.9% on the PubMedQA [101], MedMCQA [100], and USMLE [205] datasets. These evaluations demonstrate the performance of GatorTronGPT in these specific Healthcare tasks.

To examine the utility of text generation in the clinical domain, the study [190] applied GatorTronGPT to generate 20 billion words of synthetic clinical text, which were used to train synthetic NLP models, denoted as GatorTronS ('S' stands for synthetic). The GatorTronS was demonstrated that it has significant ability in clinical concept extraction and medical RE, which outperforms GatorTron [181].

16) *HuatuoGPT*: HuatuoGPT [44] is a Chinese LLM designed specifically for medical consultation purposes. Its training approach incorporates a combination of distilled data from ChatGPT and real-world data obtained from doctors during the supervised fine-tuning stage. The study highlights that ChatGPT responses are often detailed, well-presented, and informative, while they lack the ability to perform like a doctor, particularly in areas such as integrative diagnosis. To address this limitation, real-world data from doctors were introduced as supplementary training data. By incorporating real-world medical expertise into the training process, HuatuoGPT aims to enhance its performance and ensure its responses align more closely with the expectations and requirements of medical professionals in a consultation setting.

A reward model was subsequently trained to align the language model with both the distilled data and the real-world data, following a Reinforcement Learning from AI Feedback (RLAIF) approach. RLAIF is employed to reward the generation of responses that possess two important qualities: being patient-friendly (learned from ChatGPT, characterized by improved presentation quality, informative content, the ability to follow instructions, and fluent conversation) and doctor-like (learned from doctors, exhibiting professional and interactive diagnostic capabilities). Technically, HuatuoGPT utilizes LLMs to score the generated responses. These scores are based on criteria such as correctness, richness of information, logical consistency, and diagnostic ability. By incorporating these evaluation metrics, the model aims to align itself with the strengths of both ChatGPT and doctors, creating responses that are not only patient-friendly but also exhibit the expertise and interactive qualities expected from medical professionals.

HuatuoGPT underwent evaluation using three Chinese QA datasets: cMedQA2 [232], webMedQA [233], and Huatuo-26M [234]. The evaluation metrics employed included BLEU, GLEU, ROUGE, and DISTINCT, which were used to assess

the quality and distinctiveness of the generated responses. Additionally, GPT-4 was utilized to review the quality of the model outputs. In the experiment results, HuatuoGPT outperformed BenTsao [187] in a set of 100 multi-turn dialogues, as determined by GPT-4. This indicates that HuatuoGPT demonstrated superior performance and generated higher-quality responses compared to the BenTsao model in the evaluated dialogues.

17) *ClinicalGPT*: ClinicalGPT [191] is a Chinese LLM explicitly designed and optimized for clinical scenarios. By incorporating extensive and diverse real-world data, such as medical records, domain-specific knowledge, and multi-round dialogue consultations in the training process, ClinicalGPT is better prepared to handle multiple clinical tasks. Furthermore, a comprehensive evaluation framework was introduced that includes medical knowledge question-answering, medical exams, patient consultations, and diagnostic analysis of medical records.

ClinicalGPT employs a training strategy inspired by the T5 model [235] to leverage the text generation capabilities of language models for various tasks. Reinforcement learning techniques are employed to enhance the fine-tuned models, aiming to generate high-quality and helpful outputs while improving the generation of medical texts. This aids in accurately describing and treating patient conditions. The training and evaluation data utilized by ClinicalGPT consist of Chinese medical question-and-answer datasets, including cMedQA2 [232], cMedQA-KG [191], and MEDQA-MCMLE [191]. Additionally, multi-turn medical conversation datasets such as MedDialog [236] and electronic health record (EHR) datasets like MD-EHR [191] are incorporated. These diversified datasets contribute to the training and evaluation of ClinicalGPT, enabling it to generate accurate and valuable medical information for a range of Healthcare-related tasks.

18) *MedAGI*: MedAGI [192] has been developed in response to the increasing number of domain-specific professional multimodal LLMs being created in the medical field. It can be regarded as a paradigm to unify domain-specific medical LLMs with the lowest cost and a possible path to achieving medical AGI, rather than it is a LLM. Its primary objective is to automatically select appropriate medical models by analyzing user queries through our innovative adaptive expert selection algorithm. This eliminates the need for retraining, regardless of the introduction of new models. Consequently, MedAGI presents itself as a future-proof solution in the constantly evolving medical domain. To evaluate the performance of MedAGI, a comprehensive study was conducted across three distinct medical domains: dermatology diagnosis, X-ray diagnosis, and analysis of pathology pictures. The results unequivocally showcased MedAGI's remarkable versatility and scalability, consistently delivering exceptional performance across diverse domains.

19) *LLaVA-Med*: LLaVA-Med [193] is a cost-efficient approach for training a vision-language conversational assistant based LLaVA [237] that can answer open-ended research questions of biomedical images. The key idea is to leverage a large-scale, broad-coverage biomedical figure-caption dataset extracted from PubMed Central, use GPT-4 to self-instruct

open-ended instruction-following data from the captions, and then fine-tune a large general-domain vision-language model using a novel curriculum learning method. The architecture of LLaVA-Med is the same as LLaVA, which consists of a LLaMA as an encoder and a CLIP [238] as vision encoder. The one impressive point is that training LLaVA-Med only need eight A100 GPUs in less than 15 hours, and the model still exhibits excellent multimodal conversational capability and can follow open-ended instruction to assist with inquiries about a biomedical image.

As for training data, LLaVA-Med proposed a novel data generation pipeline to create biomedical multimodal instruction-following data (image, instruction, output) with PMC-15M [239], the largest biomedical image-text datasets. The whole training stage consisted of three stages, including biomedical concept feature alignment, end-to-end instruction-tuning, and fine-tuning for downstream tasks. Following, LLaVA-Med was evaluated on three biomedical medical visual question answering (VQA) datasets, VQA-RAD [240], SLAKE [241], and PathVQA [242]. On these datasets, LLaVA-Med achieved 84.19%, 85.34%, and 91.21% accuracy with the LLaVA vision encoder under the closed-set predictions setting.

20) *OphGLM*: OphGLM [194] was a large multimodal model designed specifically for ophthalmic applications. It introduced visual capabilities into LLMs, enabling it to serve as an ophthalmic language and vision assistant. The first major advancement of OphGLM involved utilizing fundus images as a starting point to develop a pipeline for disease assessment and diagnosis, as well as lesion segmentation, thus enabling the model to perform common ophthalmic disease diagnosis. Additionally, OphGLM constructed a novel dataset for ophthalmic multimodal instruction-following and dialogue fine-tuning. This dataset was created using disease-related knowledge data and publicly available real-world medical dialogues, enhancing the model's ability to understand and respond to ophthalmic-specific instructions.

To ensure a dataset that closely resembles real-world QA scenarios and enhances the interactive experience, OphGLM implemented a two-stage strategy for constructing medical conversations. In the first stage, OphGLM focused on constructing fine-tuning data based on instructions. For this purpose, genuine doctor-patient dialogues related to ophthalmic diseases were extracted from the MedDialog dataset [236]. To enable ChatGPT to simulate a medical expert, a set of prompts was designed to extract patients' intentions from publicly available doctor-patient dialogues. OphGLM aimed to provide professional and detailed medical explanations. In the second stage, OphGLM created the fine-tuned fundus dialog data, which involved five steps: (1) Generating prompts using real-world medical-patient conversations and knowledge graphs; (2) Developing medical knowledge-based instructions and conversations using the ChatGPT interface; (3) Conducting data cleaning to refine and prepare instances; (4) Eliminating duplicate data by validating against existing datasets; (5) Assessing instance quality through manual review and GPT4. Finally, new instructions and conversations were incorporated into the fundus dialog pool, enhancing the diversity and

richness of the dataset.

21) *SoulChat*: SoulChat [195] is an instruction-tuned LLM specifically designed for mental health applications in the Chinese language. Its primary focus is on fostering empathy and understanding. Through an investigation of existing AI-based counseling platforms, SoulChat identified a gap in the counseling process. It observed that users seeking online psychological help often provide a lengthy self-description, and in response, the AI counselor delivers a lengthy reply, missing the gradual process of confabulation that occurs in actual counseling sessions. In contrast, in the actual counseling process, there are multiple rounds of communication between the user and the counselor, in which the counselor guides the user through the process of confabulation and provides empathy.

To address this issue, SoulChat took a proactive approach by constructing a comprehensive dataset. They created over 150,000 instances of single-round long text counseling instructions and corresponding answers, totaling more than 500,000 responses. Additionally, they employed ChatGPT and GPT4 to generate approximately 1 million instances of multi-round response data. During pre-experiments, SoulChat discovered that counseling models solely driven by single-round long texts tended to produce lengthy responses that bored users and lacked the ability to guide them toward confiding. On the other hand, models solely driven by multiple rounds of counseling conversation data weakened their capability to provide effective suggestions. To strike a balance, SoulChat adopted a hybrid approach, combining single-round and multiple-round instances to construct the SoulChatCorpus. This corpus encompasses over 1.2 million samples, ensuring a rich and diverse dataset that captures the strengths of both single-round and multi-round counseling approaches.

22) *Med-Flamingo*: Med-Flamingo [196] is a vision-language model specifically designed to handle interleaved multimodal data comprising both images and text. Building on the achievements of Flamingo [243], one of the pioneering vision-language models known for its contextual learning and few-shot learning abilities, Med-Flamingo further enhances these capabilities for the medical domain. It achieves this by pre-training diverse multimodal knowledge sources across various medical disciplines, thereby unlocking few-shot generative medical VQA capabilities.

This study also proposed a novel dataset MTB that enables the pre-training of a multimodal few-shot learner for the general medical domain. Data of MTB consist of chopped cleaned text and images, collected from a set of 4,721 medical textbooks. Besides, PMC-OA dataset [244] was also employed, which consists of 1.6M image-caption pairs collected from PubMedCentral's OpenAccess subset.

In terms of evaluation, the study on Med-Flamingo introduced a unique evaluation dataset called Visual USMLE. This dataset combines medical VQA with complex, cross-specialty medical reasoning, resembling the format of the USMLE test. Visual USMLE consists of 618 USMLE-style questions that not only incorporate images but also include a case vignette and potential tables of laboratory measurements. In addition, the original USMLE test format was modified from multiple-

TABLE VI
THE PERFORMANCE SUMMARIZATION FOR DIFFERENT HEALTHCARE
LLMS ON THREE POPULAR DATASETS.

	USMLE(%)	MedMCQA(%)	PubMedQA(%)
Finetuned BERT	44.62 [167]	43.03 [154]	72.20 [167]
Galactica	44.60	77.60	77.60
PMC-LLaMA	44.70	50.54	69.50
GatorTronGPT	42.90	45.10	77.60
DoctorGLM	67.60	-	-
MedAlpaca	60.20	-	-
Codex	60.20	62.70	78.20
Med-PaLM	67.60	57.60	79.00
Med-PaLM 2	86.50	72.30	81.80
GPT-4	86.70	73.66	80.40
Human	87.00	90.00	78.00

choice to open-ended. This adjustment increases the difficulty and realism of the benchmark, as the models are required to independently generate differential diagnoses and potential procedures, rather than selecting from a limited set of answer choices. Finally, Med-Flamingo was evaluated on three VQA datasets (VQA-RAD [240], Path-VQA [245], Visual USMLE) with designed clinical evaluation score, BERT similarity score, and Exact-match score [244].

C. Summary

In this section, we present an overview of existing PLMs and LLMs in the Healthcare domain, highlighting their respective research focuses. Furthermore, we provide a comprehensive analysis of the performance of these LLMs on benchmark datasets such as USMLE, MedMCQA, and PubMedQA. The summarized results of these evaluations can be found in Table VI. The intention behind this analysis is to showcase the progress in Healthcare QA development and offer a clear comparison between different Healthcare-focused LLMs. In conclusion, two of the most robust LLMs identified in this analysis are Med-PaLM 2 and GPT-4. It is important to note that while GPT-4 is a general-purpose LLM, Med-PaLM 2 is specifically designed for Healthcare applications. Additionally, it is worth highlighting that the gap between LLM performance and human performance has significantly narrowed, indicating remarkable progress in the development of LLMs for Healthcare-related tasks.

As mentioned earlier, one notable difference between PLMs and LLMs is that PLMs are typically discriminative AI models, while LLMs are generative AI models. Although there are some auto-regressive PLMs like GPT-1 and GPT-2 also evaluated with classification tasks, and auto-encoder PLMs have been more prominent during the PLMs period. As for LLMs, with their powerful capabilities, they have successfully unified various Healthcare tasks as QA tasks or dialogue tasks with the generative way.

From a technological perspective, most PLMs studies focus on improving neural architectures and designing more efficient pre-training tasks. On the other hand, LLM studies primarily emphasize data collection, recognizing the importance of data quality and diversity due to the over-parameterization strategy employed in LLM development. This aspect becomes even more crucial when LLMs undergo Instruction Fine-Tuning to

align with human desires. A study [16] reveals that the selection of mixed ratios of different training data significantly impacts the performance of LLMs. However, these mixed ratios of pretraining and Instruction Fine-Tuning, often referred to as a “special recipe” from different strong LLM developers, are rarely publicized. Therefore, apart from instruction fine-tuning, we anticipate the emergence of more exciting and innovative methods for training LLMs, particularly those designed to handle unique features of Healthcare data.

In terms of the investigated Healthcare LLMs mentioned above, most of them are derived from general LLMs, except for GatorTron, Galactica, and GatorTronGPT. For these LLMs, IFT approach is the most commonly utilized training technique. Many LLMs make use of instruction data to fine-tune their models, including Galactica, MedAGI, OphGLM, MedAlpaca, BenTsao, PMC-LLaMA, BianQue, Med-PaLM 2, GatorTronGPT, and ClinicalGPT. However, compared to IFT, RLHF/RLAIF is less commonly employed, with only MedAlpaca and HuatuoGPT utilizing this technology. The main reason for this limited application of RLHF/RLAIF is believed to be the lack of sufficient stability, as mentioned in the study [246]. From this part of the survey content, we have identified two emerging trends. Firstly, there is a growing exploration of multi-model approaches, including LLaVA-Med, MedAGI, OphGLM, Visual Med-Alpaca, and Med-Flamingo. Secondly, Chinese Healthcare LLMs are rapidly developing, with examples such as DoctorGLM, ClinicalGPT, SoulChat, BenTsao, BianQue, and HuatuoGPT. In addition to the development of Healthcare LLMs, there are also studies investigating the use of general LLMs for health-related tasks, such as Codex-Med, GPT-4-Med, and DeID-GPT. Regarding techniques about LLM optimization, LoRA [216], ZeRO [247], and model quantization [217] are the three most commonly employed methods. These optimization technologies are discussed in detail in Section IV-D.

Finally, it is worth noting that many Healthcare LLM papers provide details about the prompts they used. This observation demonstrates the prompt brittleness, as different prompts can have a significant impact on the model’s performance. Modifications in the prompt syntax, sometimes in ways that are not intuitive to humans, can lead to significant changes in the model’s output [248]. This instability is more matters for Healthcare than other general applications.

IV. TRAIN AND USE LLM FOR HEALTHCARE

In this section, we review the training and usage of LLM for Healthcare. First, we introduce the pre-training methods from PLMs and post-training methods from LLMs. Then, the usage of LLMs, including fine-tuning, in-context learning, CoT, and AI-agent. To achieve the promising usage of LLMs, an efficient training frame and data are necessary. Thus we also summarize the commonly used training data for Healthcare LLM and efficient training framework. The whole content structural arrangement is shown in Figure 5.

A. Pre-training Methods

1) *Masked Language Modeling*: the concept of masked language modeling (MLM) is first introduced with the release

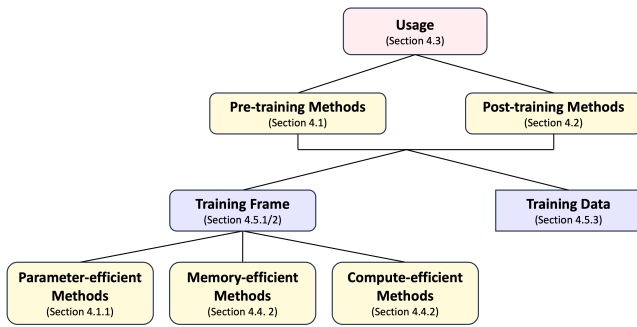


Fig. 5. The structural content arrangement for the section Train and Use LLM for Healthcare.

of the BERT [249] as shown in Eq. 1.

$$loss_{MLM} = \sum_i \log P(x_i | x_{i-k}, x_{i-k+1}, \dots, x_{i+k-1}, x_{i+k}; \theta), \quad (1)$$

where k is the window size of context, and the conditional probability P is modeled by a neural network with parameters θ .

Following the success of BERT, MLM has emerged as a widely adopted approach in NLP research. It has been extended and improved upon in subsequent pretrained language models, including RoBERTa [22], ALBERT [250], ERNIE [251] and DeBERTa [252]. These models build upon the foundations laid by BERT, refining and expanding MLM techniques for further advancements in NLP tasks.

2) *Next Word Prediction*: next word prediction is a language modeling task to predict the next word or sequence of words given the input context. This task is a core component of GPT series models, utilizing statistical patterns and linguistic structures to generate accurate predictions based on the context provided. Specifically, in next-word prediction, the objective is to assign probabilities to all possible words in the vocabulary and select the word with the highest probability as the prediction for the next word as Eq. 2.

$$loss_{PLM} = \sum_i \log P(x_i | x_{i-k}, x_{i-k+1}, \dots, x_{i-1}; \theta). \quad (2)$$

The difference between MLM and next-word prediction is that the former can see bi-direction inputs while the latter can only use information before the current input token i . It also has proven to be a fundamental pertaining task for both PLMs and LLMs.

3) *Sequence-to-sequence MLM*: sequence-to-sequence MLM, which is an extension of traditional MLM, has been adapted for text generation tasks like machine translation, text summarization, and question answering. It is first introduced in T5 (Text-to-Text Transformer) by the study [235], which present a unified framework for transforming various text-based language problems into a text-to-text format.

In traditional MLM, the training goal is to predict masked words in a single sequence. However, sequence-to-sequence MLM extends this objective to predict masked tokens in both the input and output sequences, simultaneously. This approach enables the model to learn the relationships and

dependencies between the input and output sequences, which is particularly advantageous for text-to-text tasks. By jointly modeling dependencies in both input and output, the model can better understand the contexts and generate more accurate translations or summaries.

The application of sequence-to-sequence MLM has also been extended to the biomedical domain, as demonstrated in SCIFIVE [164], a domain-specific T5 model that has been pretrained on biomedical corpora. SCIFIVE outperforms many compared baselines, highlighting the potential of sequence-to-sequence MLM in biomedical text generation tasks.

4) *Replaced Token Detection*: replaced token detection was implemented upon the launch of the ELECTRA model [253]. In contrast to the conventional approach such as BERT, where input corruption involves replacing certain tokens with [MASK] and subsequently training a model to reconstitute the original tokens, ELECTRA incorporates the concept of generative adversarial networks (GANs). This entails substituting selected tokens with plausible alternatives drawn from another generator network. This method demonstrates superior efficiency in sample utilization for training, as compared to the conventional masked language modeling technique.

5) *Sentence Boundary Detection*: the pre-training objective of sentence boundary detection constitutes initially introduced in SpanBERT [254]. This method takes into consideration the original boundary of text spans, introducing enhancements aimed at more effectively addressing tasks necessitating the modeling of inter-span text relationships. Unlike the conventional BERT, where random tokens are masked and the pre-training objective focuses on predicting these tokens for contextual understanding, certain tasks such as NER or RE require a more holistic consideration of relationships spanning entire textual segments, rather than focusing on individual tokens.

SpanBERT tackles this challenge by extending BERT’s training objective via two key modifications: firstly, contiguous random spans are masked, as opposed to individual random tokens; secondly, a novel training objective is introduced, referred to as sentence boundary detection. This objective is designed to predict the complete masked span given the observed tokens within its boundaries. Consequently, the model becomes proficient in capturing contextual information among words within a given span. This approach notably has benefits for tasks involving the identification of entities, relations, or other structured information embedded within texts.

6) *Next Sentence Prediction*: next sentence prediction (NSP) is a pre-training objective originating from the BERT, together with MLM. With NSP in pre-training, models tend to recognize semantic and syntactic correlations among sentences, enhancing their ability to generate more contextual responses.

The NSP task is conventionally formulated as a binary classification problem, wherein a pair of sentences - a context sentence and a next sentence - is presented to the model. The model’s objective is to predict whether the next sentence logically follows the context sentence in the original text.

However, empirical study in RoBERTa [22] has demonstrated that the model performance could be significantly

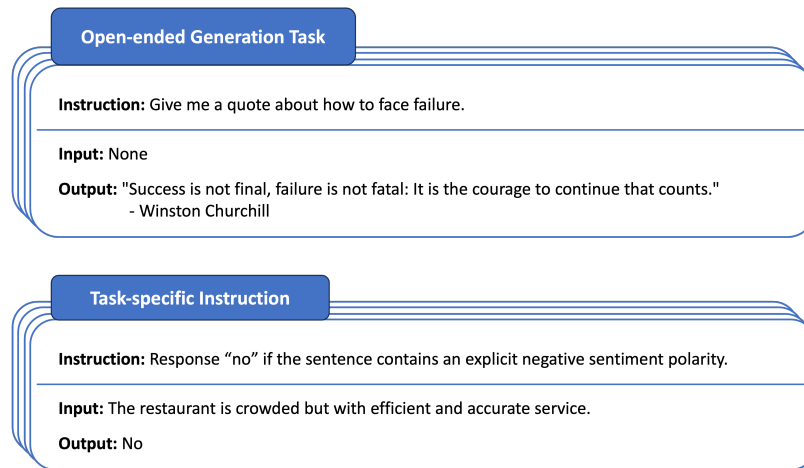


Fig. 6. The examples of instructions demonstrations. For open-ended generation task, there can just instructions without inputs. For task-specific instruction, a LLM needs respond to specific inputs.

improved without NSP pre-training. Recent models such as GPT-2, GPT-3, XLNet, and T5 have also explicitly removed the NSP objective in their pretraining step.

7) *Sentence Order Prediction*: sentence order prediction (SOP) is an effective pre-training objective introduced in ALBERT [250], which involves training the model to predict the correct order of sentences within a pair of sentences. More specifically, during the pretraining phase, a model is presented with pairs of sentences, and the objective is to determine whether the sentences are in the correct order or if they should be swapped. This objective encourages the model to learn contextual relationships on the sentence level.

SOP is an attempt to address the ineffectiveness of NSP, which focuses solely on binary sentence ordering without considering the nuances of document structure.

B. Post-training Methods

1) *From predicting tokens to follow instructions – Instruction Fine-Tuning and Supervised Fine-tuning*: through the pretraining process, we can obtain a strong but uncontrolled model, which can perform precise token predictions but is insufficient to follow the user’s instructions in a useful way. For such reason, the study [255] proposed Instruction Fine-Tuning (IFT), which involves fine-tuning the base model on demonstrations of written directions using diverse sets of tasks, along with traditional NLP tasks such as sentiment analysis, text classification, and summarization.

The used instructions demonstrations consist of three key components: the instruction itself, the inputs, and the outputs. The inputs are optional, like open-ended generation with ChatGPT, and solely rely on the instructions. When both inputs and outputs are included, they form an instance, and there can be multiple instances of inputs and outputs for a given instruction. The aim is to enhance the ability of instruction following. Figure 6 shows examples of instructions demonstrations.

After IFT, PLMs fine-tuned by instructions may not always generate useful and safe responses. These behaviors include being evasive by consistently providing unhelpful responses such as “I’m sorry, I don’t understand”, or generating unsafe

responses to user inputs on sensitive topics. To address and mitigate such behavior, the process of Supervised Fine-tuning (SFT) is employed, which involves fine-tuning the base language model using high-quality human annotated data with a focus on ensuring helpfulness and harmlessness⁷.

2) *Reinforced Learning from Human Feedback (RLHF)*: RLHF is employed in recent LLM studies, including general LLMs [214], [256] and medical LLMs [44], [186]. The goal of RLHF is to train AI systems to align with human goals, which remains the same as SFT. Actually, RLHF can be regarded as a cost-effective alternative to the SFT method with two differences: (1) SFT utilizes data from human responses for training, aiming to bring the model closer to human-like behavior without involving a direct comparison process. On the other hand, the RLHF process begins with training a reward model to rank, where different rewards (high or low) are assigned during the reinforcement learning stage (the rewards are scaled to have positive and negative values rather than both being positive). The introduction of a comparison process in RLHF helps guide the output of the model to align more closely with human behavior. (2) When considering the same amount of data, collecting data for SFT is generally more challenging compared to RLHF. Moreover, each piece of SFT data contains more information or training value than a piece of RLHF data in terms of ranking.

According to the definition from the study [257], RLHF refers to methods that combine three interconnected steps: feedback collection, reward modeling, and policy optimization (the pre-training process and SFT are regarded as an optional step 0, which initiate LLMs to perform RLHF).

First, RLHF employed an initiated LLM π_θ generated some prompts x_i (questions or instructions), and then used human \mathcal{H} manually respond to these prompts with the hypothesis that \mathcal{H} follow the map function f which has consistent with the

⁷It should notice that the concepts of SFT and IFT are closely related without a strict boundary, and they are not strictly distinguished yet. A slight difference exists in recent literature, where SFT is frequently applied to safety topics rather than training the ability of instruction following. With such conditions, SFT is typically performed after the IFT stage.

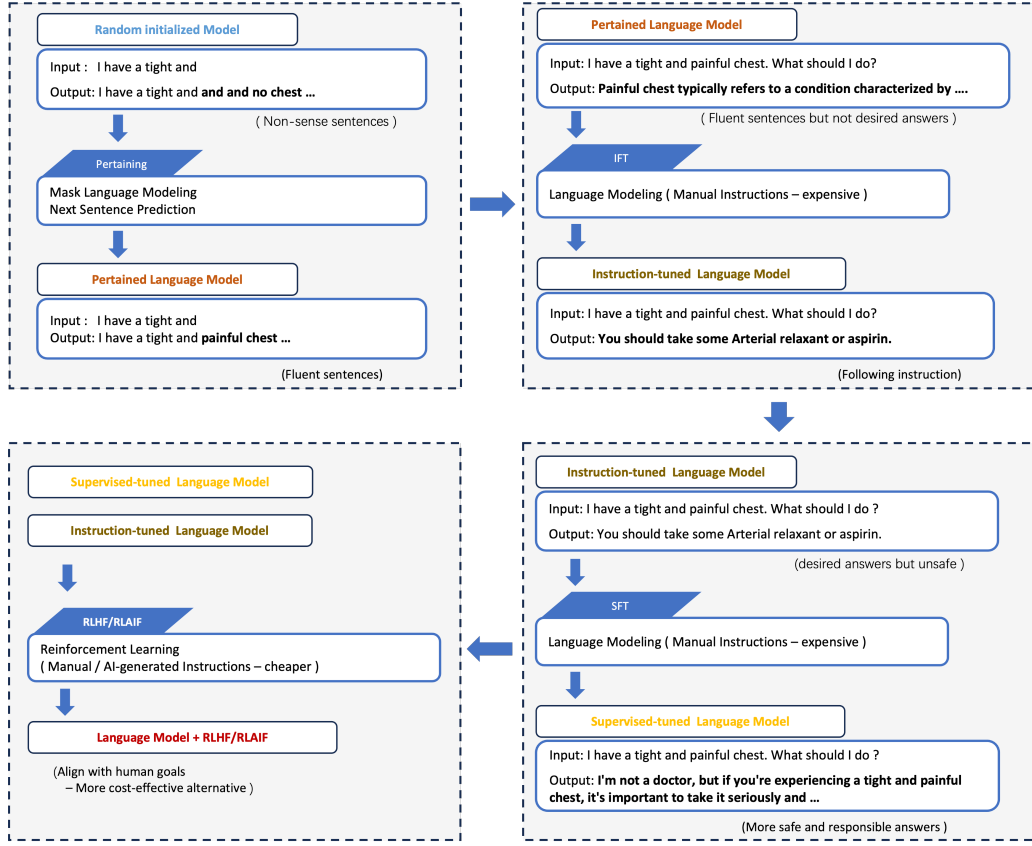


Fig. 7. The relations among Pretraining, IFT, SFT, RLHF, and RLAIF. The integration of these components can form a pipeline that facilitates the training of a LLM, enhancing its ability to effectively follow instructions.

required reward function $r_{\mathcal{H}}$. The human response is regarded as human feedback y_i . Additionally, different humans have different feedback, which can be integrated into random noise ϵ_i . The whole step of feedback collection to get a y_i can be donated as:

$$y_i = f(\mathcal{H}, x_i, \epsilon_i). \quad (3)$$

Second, RLHF needs to train a reward model \hat{r}_{\emptyset} to fit the required reward function $r_{\mathcal{H}}$. Given the collected pairs of prompts and related human feedback \mathcal{D} , the \hat{r}_{\emptyset} is trained by minimizing the loss following

$$\mathcal{L}(\emptyset) = \sum_{i=1}^n l(\hat{r}_{\emptyset}(x_i), y_i) + \lambda_r(\emptyset), \quad (4)$$

where $\lambda_r(\emptyset)$ is a regularizer, l is a chosen loss function and a cross-entropy loss is the most common choice.

Third, policy optimization means using the fitted reward model \hat{r}_{\emptyset} to fine-tune the base LLM π_{θ} with reinforcement learning. \hat{r}_{\emptyset} is further trained by maximizing the reward following

$$\mathcal{R}(\theta') = \mathbb{E}_{x \sim \pi_{\theta'}} [\hat{r}_{\emptyset}(x) + \lambda_{\beta}(\theta, \theta', x)], \quad (5)$$

where λ_{β} is a regularizer.

3) *From Human Feedback to AI Feedback*: for IFT, data usually came from human-written instructions [255]. To further expand the instruction dataset, a large community effort of hand-crafted instructions is conducted [258]. However, manual

labor is expensive and cannot support LLMs with continuous improvements. To deal with such problems, some studies explored self-instruct, which aims to instruct an LLM by itself or by other LLMs [129], [259]–[261].

Among these studies, BenTsao generated instruction data based on a medical knowledge graph CMeKG [262]. They first sampled knowledge instances from the knowledge graph and then generated 8,000 instructions based on the specific knowledge with the OpenAI API. These automatically generated instructions were employed to fine-tune their base model.

HuatuoGPT [44] incorporated a combination of distilled and real-world data to enhance the instruction following ability, encompassing medical instruction data as well as medical conversation data. Subsequently, the researchers employed Reinforcement Learning with AI Feedback (RLAIF) to effectively leverage the strengths of both data types. Specifically, HuatuoGPT employed real instructions and conversations as training data, extracting multiple responses from their fine-tuned model. These responses were subsequently evaluated by an LLM, such as ChatGPT, based on criteria such as informativeness, coherence, adherence to human preferences, and factual accuracy, taking into account actual diagnoses provided by real doctors. Then, these paired response data are used to train the reward model, using the fine-tuned model as its backbone for better generalization.

Baize [261] proposed a method called self-distillation with feedback (SDF) as an alternative to RLHF. They utilized

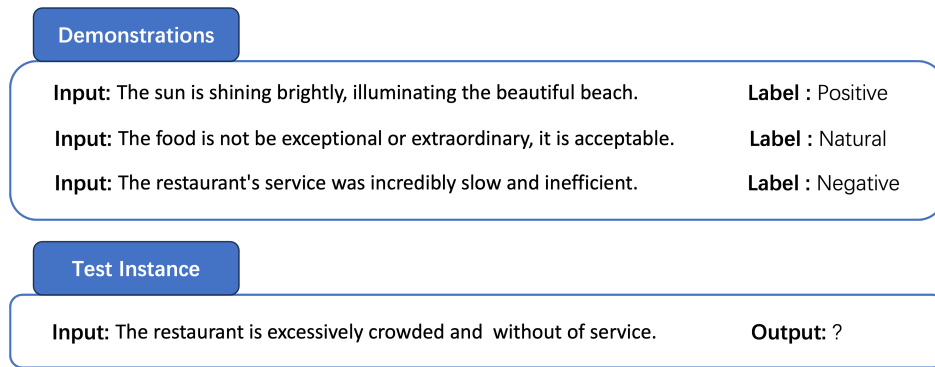


Fig. 8. An In-context Learning example for sentiment analysis task.

ChatGPT to automatically generate a high-quality multi-turn chat corpus. Initially, they collected a total of 111.5k dialogues through self-chat, and then an additional 47k dialogues in the medical domain were used to fine-tune a specialized Baize model for Healthcare. Following the SFT stage, the resulting Baize model was used to generate four responses for the collected instructions, which were then utilized in the following SDF process. Notably, Baize's data collection pipeline differed from Alpaca's single-turn self-instruct, as Baize focused on enhancing the model's multi-turn chat capability by leveraging high-quality chat transcripts obtained from ChatGPT.

Humpback [129] utilizes a substantial amount of unlabeled data to construct a superior instruction-tuning dataset through the implementation of an iterative self-training algorithm. The proposed approach consists of two key components: self-augment and self-curate. These components are designed to automatically generate high-quality training examples, thereby enhancing the performance of the model. In the self-augment phase, Humpback needs a seed instruction set and a web corpus, and they train a backward model to generate instructions for unlabelled data to create candidate training data. Acknowledging the potential lack of quality in these generated data, Humpback leverages self-curate to assess their quality. Consequently, the model can focus on self-training using only the most reliable (instruction, output) pairs.

Besides the above, the study [259] is a purely model-generated large IFT dataset. Vicuna [263] tried to learn the ChatGPT responses when interactively chatting with humans.

4) *Summary:* in Section IV-A and Section IV-B, we talk about how to train PLMs and LLMs. We first provide a comprehensive overview of diverse pre-training techniques, elucidating their operational mechanisms. Different pre-training methodologies contribute substantially to endowing these models with a foundational grasp of linguistic knowledge including grammar, syntax, and semantics. These methodologies facilitate to build the representations for words, phrases, and sentences and capture their contextual interrelationships.

Furthermore, different pre-training methods enable transfer learning, which is one of the most significant advantages of PLMs. The knowledge acquired during pretraining can be subsequently fine-tuned for specific tasks, thereby reducing the amount of data and training time when applied to various downstream tasks. This transfer of knowledge greatly im-

proves the efficiency and effectiveness of models in Healthcare applications. By taking into account the efficacy of these pre-training tasks, it becomes evident that they contribute to the development of robust models capable of comprehending context at various levels, encompassing individual words to entire paragraphs. For example, NSP and SOP prove instrumental in capturing sentence-level relationships while MLM and NWP facilitate the acquisition of contextual information at the word level. This deep contextual understanding is essential for tasks that involve complex language comprehension, such as question answering, sentiment analysis, and text summarization. It is noteworthy that these pre-training methodologies are universally applicable and readily adaptable to the domain of Healthcare. However, it is imperative to acknowledge that the quality and effectiveness of these pre-training techniques wield substantial influence over the overall capabilities and usefulness of language models in Healthcare-related applications.

When it comes to LLMs, the associated training methods primarily focus on IFT, SFT, RLHF, and RLAIF (with AI Feedback). The relations among them are summarized in Figure 7. Different from pretraining described in Section IV-A, whose goal is minimizing the distance between training data with generated data, IFT aims to change the model from autoregressive prediction to own the ability of instruction following. Namely, the objective of IFT is to adjust the model's output to closely align with the response to a given instruction, rather than precisely predicting the next token. Further, SFT provided LLM with alignment with human goals, not only with precise instruction responses but also with controllability for those responses. RLHF offers greater flexibility when compared to SFT. By employing RLHF, we have the capability to quantify and incorporate diverse properties into the desired model output. This is achieved through the learning of a reward function, allowing us to incentivize qualities such as truthfulness, non-toxicity, and helpfulness to humans.

Typically, the IFT/SFT phase involves the utilization of high-quality and costly training data. On the other hand, RLHF relies on relatively lower-quality and less expensive data. If the IFT/SFT phase is skipped, the transition from pre-training directly to RLHF may result in a substantial gap, as the training data for RLHF may not be sufficient to achieve desirable fine-tuning outcomes. Furthermore, the

training process of RLHF can be sensitive to parameter settings at times. However, incorporating IFT/SFT can be seen as an extension of pre-training in the RLHF stage, potentially alleviating the issue of unstable training and leading to more stable outcomes.

As for RLAIIF, this technology can be considered as a more cost-effective alternative to RLHF. The study [264] is a more detailed survey about how to learn from AI feedback. They examine and categorize a diverse range of recent studies that employ these strategies and encompass various stages, such as training, generation, and post-hoc correction. Additionally, they provide a comprehensive overview of the significant applications of these approaches and conclude by delving into potential future directions and the challenges they may entail.

C. Usage

1) *From Fine-tuning to In-context Learning*: in the PLMs era, the most common scenarios of applying PLMs to various down-stream tasks is fine-tuning a general PLM, such as BERT [8] and BART [266] with domain-specific or task-specific data. After parameters are updated by fine-tuning, these PLMs can achieve various goals.

On the contrary, when comes to LLMs, it is hard to tune their parameters due to required GPU memories and training time cost. In-context learning (ICL) is a promising technology to deal with such problems, which just concatenates some demonstration examples with inputs and feeds into LLMs, without any parameter updates. Figure 8 shows an example of ICL.

As shown in Figure 8, the whole input consists of demonstrations and a test input. Each of these demonstrations is an input-label pair. For example, “Input: The sun is shining brightly, illuminating the beautiful beach. Label: Positive” is a demonstration. There can be any number of demonstrations, as long as not exceed the total input length required by the model. Output in Figure 8 is the final prediction we want to obtain from LLMs.

Additionally, there are several implicit concepts that require attention when utilizing ICL. These include the input distribution, label space, demonstration format, and input-label mapping. Regarding input distribution, it pertains to whether the input sentences in the demonstrations and the test instances originate from the same domain, such as news or medical corpora. Label space refers to whether the labels assigned to the demonstrations share the same semantic space as the labels assigned to the test instances. The demonstration format encompasses the manner in which demonstrations are structured. The most common format involves an input sentence accompanied by a related label, although it is also can utilize solely an input sentence or a label as a demonstration. Input-label mapping pertains to the appropriateness of assigning a suitable label to the input sentences in the demonstrations.

The study [265] conducted some interesting experiments about the above concepts and discussed what really works in ICL as shown in Figure 9. The findings indicate that the final results are significantly influenced by input distribution, label space, and demonstration format. However, concerning

input-label mapping, the current results suggest that as long as there is the appropriate label space, the accuracy of the labels themselves does not have a substantial impact on the outcomes.

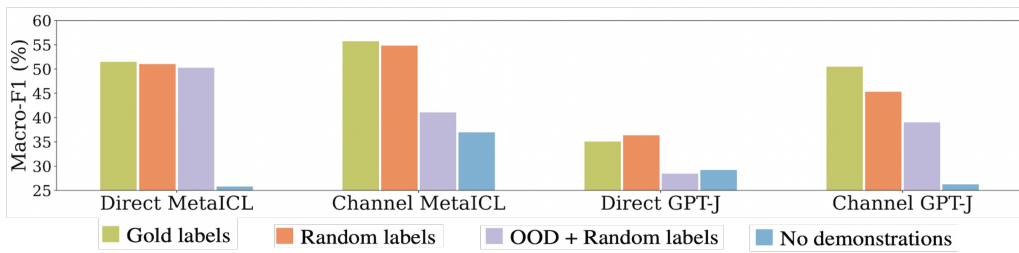
2) *From System 1 Deep Learning To System 2 Deep Learning – Chain-of-Thought*: according to the report by Bengio et al. [267] presented at the 2019 Conference on Neural Information Processing Systems (NeurIPS), two distinct categories of Deep Learning systems exist: System 1 and System 2. System 1 encompasses the current applications of deep learning, including image recognition, face recognition, machine translation, sentiment classification, speech recognition, and autonomous driving. On the other hand, System 2 represents the future potential of deep learning, involving tasks such as reasoning, planning, and other logic-based and reasoning-oriented activities.

System-1 tasks in the field of NLP have been largely resolved, demonstrating significant progress. However, progress in System-2 tasks has been limited until recently when the emergence of advanced LLMs triggered a significant shift. The study [33] proposed the CoT prompting, which found it can significantly improve the reasoning performance of LLM by adding a series of intermediate steps as shown in Figure 10.

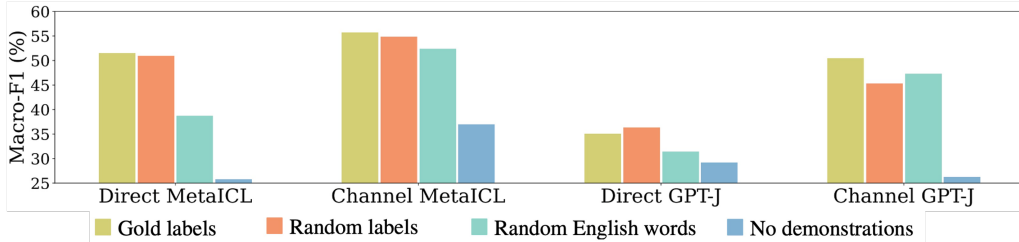
Furthermore, the study [268] found that by just adding a sentence “Let’s think step by step”, the reasoning ability of LLMs can be significantly boosted. For example, adding this simple sentence can raise accuracy from 17.7% to 78.7% on MultiArith [269] dataset, and from 10.4% to 40.7% on GSM8K [270] dataset. Later, there are many CoT studies [44], [189], [195] aiming to enhance the logical reasoning ability of LLM in various Healthcare applications by exploring different prompting.

3) *AI Agents*: The core idea behind recent AI agents is to build autonomous agent systems that utilize LLMs as their central controllers. These systems consist of several components, including Planning, Memory, Tool Use, and Action, as described in the study [271]. The Planning component plays a crucial role in breaking down complex tasks into smaller and manageable sub-goals. This enables the agent to handle large tasks more efficiently by tackling them step by step. The Memory component provides the agent with the ability to store and retrieve information over extended periods. It typically utilizes an external vector store and fast retrieval mechanisms, allowing the agent to retain relevant knowledge and recall it as needed. With the Planning and Memory components in place, AI agents can take actions and interact with external tools. AutoGPT⁸ is an example of such an autonomous agent system. It leverages GPT-4 to autonomously develop and manage operations. When provided with a topic, AutoGPT can think independently and generate steps to implement the given topic, along with implementation details. This shows the agent’s ability to plan, utilize its memory, and take appropriate actions to accomplish tasks in an autonomous manner. Relevantly, AgentBench [272] proposed a benchmark to comprehensively evaluate the abilities of LLMs as agents, from as operating systems, web browsing, web shopping, house-holding, lateral

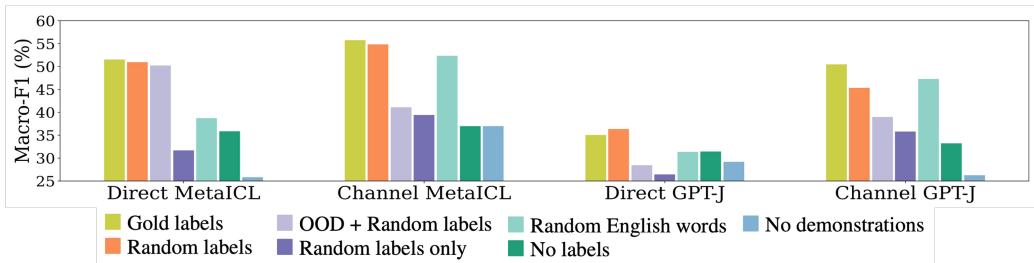
⁸<https://github.com/Significant-Gravitas/Auto-GPT>



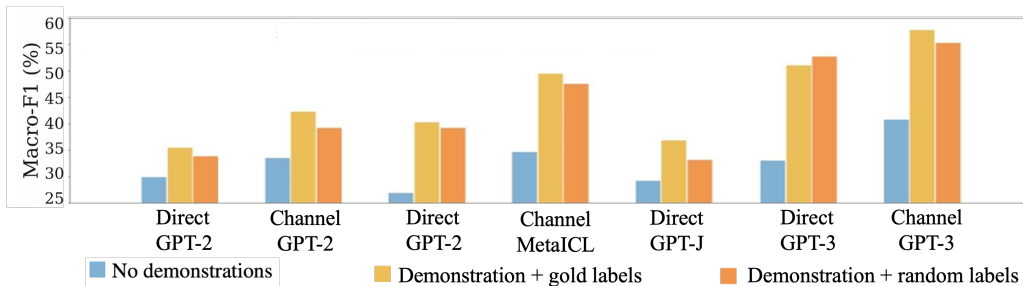
(a) Impact of the inputs distribution. OOD means out-of-distribution. The results show that the input distribution has significant effects on the final performance.



(b) Impact of the label space. The results show that the label space have significant effects on the final performance.



(c) Impact of the demonstration format. The results show that the demonstration format has significant effects on the final performance.



(d) Impact of the input-label mapping. The results show that input-label mapping only has slight effects on the final performance.

Fig. 9. What Makes In-Context Learning Work? ★ The figures all comes from the study [265]. We perform the proper layout and arrangement for discussions ★. We only list the classification task (x-axis) here and sub-figure (d) shows parts of the original results for clarity.

thinking puzzles, digital card games, knowledge graphs, and databases. LangChain⁹ is one of the most popular libraries to build AI agents systems, which can help to combine LLMs with other sources of computation or knowledge.

As far as we know, AI agents have not been widely adopted in the Healthcare field. However, we anticipate the development of more capable AI agent systems in this domain. For instance, it is possible to train specialized models for different medical processes, such as hospital guidance, auxiliary diagnosis, drug recommendation, and prognostic follow-

up. These relatively small models can be integrated into a comprehensive AI medical system, where an LLM serves as the central controller. Additionally, specialized disease systems can be established for each department within the Healthcare system. The LLM can play a crucial role in determining which specialized disease systems should be involved in a particular case. This helps in effectively allocating resources and providing specialized care.

Overall, the vision is to leverage AI agents and LLMs to create comprehensive and specialized AI systems in Healthcare, covering various medical processes and enabling efficient decision-making and patient care.

⁹<https://github.com/langchain-ai/langchain/blob/master>

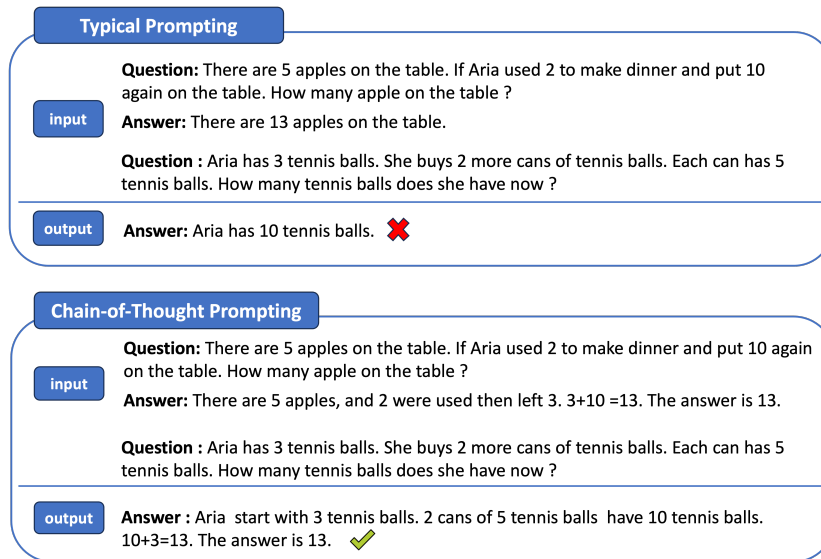


Fig. 10. An example of Chain-of-Thought (CoT). CoT is the sequential and logical prompts, which can help LLM split a complex problem into multiple simple steps.

4) *Summary:* In the era of PLMs, the most common practice is to fine-tune general PLMs for specific purposes. However, this approach requires additional computational resources for deployment and is often limited in its scope of usage. Conversely, in the era of LLMs, the focus shifts towards utilizing a powerful LLM without the need for parameter updates. Instead of fine-tuning, LLMs are typically used in conjunction with ICL. This means that by providing demonstration examples, LLMs can effectively perform various functions in Healthcare applications. Moreover, by providing step-by-step demonstration examples, LLMs can significantly enhance their logical reasoning abilities. This allows them to make more informed decisions and provide intelligent responses in Healthcare scenarios.

In summary, in the LLM era, the emphasis is on using powerful LLMs without fine-tuning, leveraging ICL, and providing step-by-step demonstrations to enhance logical reasoning capabilities. This approach offers promising possibilities for applying LLMs in Healthcare applications and offers more interpretability.

D. Parameters-, Memory-, and Compute-efficient Methods

1) *Parameters-efficient Methods:* As the model parameter size gets bigger and bigger, the cost of doing full fine-tuning on the downstream task dataset is getting higher and higher. To alleviate this problem, a series of parameters-efficient tuning methods are proposed to help pre-trained LLMs efficiently adapt to a variety of downstream tasks. These are very practical methods when adopting general LLMs to the Healthcare field.

In general, there are three main typical methods used in parameters-efficient optimizations: Adapters, Prefix Tuning, and LoRA. Adapter methods [273]–[275] involve inserting smaller neural network modules into the intermediate layers of PLMs or LLMs. During fine-tuning, only the parameters of the adapter modules are trained while keeping the rest of the model parameters fixed.

Prefix Tuning [10], [276] is another approach where a trainable prefix is added to the input sequence or hidden layers. These added prefixes do not correspond to real tokens and are free parameters that can be learned. Prefix Tuning fixes the pre-training parameters of PLMs or LLMs, optimizes only the task-specific prefixes, and requires only one copy of a small number of prefixes for each task to be stored during deployment.

The aforementioned approaches exhibit limitations. The Adapter method introduces additional inference latency when incorporating an adapter module, while Prefix Tuning reduces the sequence length available for downstream task processing due to the allocation of a portion of the sequence length for prefixes. To overcome these issues, LoRA [216] presents a superior approach for achieving parameter-efficient fine-tuning while avoiding the aforementioned problems. LoRA's core concept involves approximating the parameter update of a full-rank weight matrix with a low-rank matrix, thereby necessitating training only a small ascending-dimensions matrix and a small descending-dimensions matrix. Notably, LoRA offers several advantages, including the absence of introduced inference latency, a significant reduction in video memory consumption, and the ability to customize for diverse tasks.

2) *Compute-efficient and Memory-efficient Methods: Parallelism.* Generally, when we train LLMs, the parameters of models, gradients, and optimized states take up the Video Random Access Memory (VRAM) of GPUs. When one single GPU cannot satisfy training requirements, Data Parallelism (DP), Model Parallelism (MP), and Pipeline Parallelism (PP) are three compute-efficient and memory-efficient strategies.

DP involves replicating model parameters on each device to achieve the compute-efficient goal. During each step of the training process, a mini-batch of data is evenly divided across all the data parallel processes. This means that each process performs forward and backward propagation on a distinct subset of data samples. Afterwards, the gradients are

averaged across all the processes and used to locally update the model parameters. This approach is specifically designed for scenarios where there is a large amount of data and relatively small PLMs. The backward of this approach is introducing redundancy in terms of memory and computational resources.

When an LLM can not fit in the VRAM, MP [277] allows to put different layers of an LLM into different devices. MP is operator-level parallelism, which utilizes the properties of certain operators to split the operator across multiple devices for computation. However, it is worth noting that not all operators are splittable. Splittable operators need to satisfy that (1) they are parallelizable, (2) and that one input is the model parameter itself. If we consider MP as splitting the LLM vertically, PP takes a different approach by horizontally partitioning the model. Each partition is then executed on a separate device, and micro-batching is employed to conceal any pipeline bubble [278]. The major drawback of MP and PP is the significant amount of communication time required between different devices, which can be regarded as memory-efficient but compute-inefficient methods.

ZeRO. Based on the parallelism scenarios mentioned above, a series of ZeRO-related studies are introduced [247], [279], [280], presents a set of memory optimization techniques. This series includes ZeRO, ZeRO-Offload, and ZeRO-Infinity, which aim to eliminate redundant parameters, utilize CPU and Random Access Memory (RAM), and introduce NVMe for improved performance.

The ZeRO [247] comprised three stages: ZeRO-1, ZeRO-2, and ZeRO-3, each corresponding to the partitioning of different components of the model. Specifically, ZeRO-1 partitions the optimizer states, ZeRO-2 partitions both the gradients and optimizer states, and ZeRO-3 partitions all model states.

However, the ZeRO [247] is essentially a parallelism solution. Different from this, the core idea of ZeRO-Offload [279] is using RAM as the supplement of VRAM. Originally, with a single V100, a 1.4B model can be trained with PyTorch and the throughput is 30 TFLOPS, and with ZeRO-Offload augmentation, a 10B model can be trained and the throughput is 40 TFLOPS.

ZeRO-Offload is more focused on single card scenarios, while ZeRO-Infinity [280] is typical of the industrial field and goes for very large-scale training. Both designed for Offload, ZeRO-Infinity pays more attention to communications between multiple GPUs.

All the above functions are integrated into the library DeepSpeed in Huggingface¹⁰. More details and related training tools can be seen in Section IV-E.

Quantization. The definition of quantization is approximating the weights or activation values represented by the high bit widths (Float32) with the low bit widths (INT16, INT8, INT4), in which the form of the external display is the discretization of continuous values. It should be noted that quantization acceleration needs compatible hardware support.

Quantization can be employed to further downsize LLMs (If an un-quantified LLM is quantified to int4, its overall size was reduced to a quarter) and accelerate the computational

TABLE VII
HEALTHCARE DATA CAN BE USED TO TRAIN LLMs.

Data	Type	size	Link
MMIC-III	EHR	58,976 hospital admissions for 38,597 patients	Homepage
MMIC-IV	EHR	covering a decade of admissions between 2008 and 2019	Homepage
CPRD [281]	EHR	over 2,000 primary care practices and include 60 million patients	Homepage
PMC	Scientific Literature	PubMed citations and abstracts of biomedical literature	Data Link
RCT [282]	Scientific Literature	8 million full-text article records	Data Link
MS-2 [283]	Scientific Literature	470,402 abstract	Data Link
CDSR [284]	Scientific Literature	7,805 abstract	Data Link
SumPubMed [285]	Scientific Literature	33,772 abstract	Data Link
The Pile	Scientific Literature	825 GB English text	Data Link
SZORC [286]	Scientific Literature	63,709 abstract	Data Link
CORD-19 [287]	Scientific Literature	1M papers	Data Link
MeQSum [288]	Medical Question Summarization	1000 instances	Data Link
CFQ-Sum [289]	Medical Question Summarization	1507 instances	Data Link
UMLS	Knowledge Base	2M entities for 900K concepts	Homepage
COMETA [290]	Web Data (social media)	800K Reddit posts	Homepage
MedDialog [291]	Dialogue	3.66 million conversations	Homepage
CovidDialog [292]	Dialogue	603 consultations	Homepage
Medical Flashboards [186]	Dialogue	3395 instances	Data Link
Wikisde [186]	Dialogue	67704 instances	Data Link
Wikisde Patient Information [186]	Dialogue	5942 instances	Data Link
MEDQA [293]	Dialogue	2208 instances	Data Link
CORD-19 [287]	Dialogue	1056660 instances	Data Link
MMMLU [287]	Dialogue	3787 instances	Data Link
Pubmed Causal [294]	Dialogue	2446 instances	Data Link
ChatDoctor [295]	Dialogue	215000 instances	Data Link
Alpaca-EN-AN [214]	English Instructions	52K instructions	Data Link
Alpaca-CH-AN [214]	Chinese Instructions	52K instructions	Data Link
ShareGPT	Conversations	61653 long conversations	Data Link
WebText	Web Data	40 GB of text	Data Link
OpenWebText	Web Data	38 GB of text	Data Link
Colossal Clean Crawled Corpus	Web Data	806 GB of text	Data Link
OpenLI	EHR, Multimodel	3.7 million images from about 1.2 million papers	Homepage
OpenLI [296]	Multimodel	526 reports and 7,470 images	Homepage
ROCO [297]	Multimodel	81,000 radiology images and corresponding captions	Homepage
MedCAT [298]	Multimodel	17,000 images includes captions	Homepage
PMK-CA [244]	Multimodel	1.6M image-caption pairs	Homepage
CheXpert [299]	Multimodel	224,316 chest radiographs with associated reports	Homepage
PatChest [300]	Multimodel	160,000 images with related text	Homepage
MMIC-CXR	Multimodel	27,835 imaging studies for 64,588 patients	Homepage
PMC-15M [239]	Multimodel	15 million Figure-caption pairs	Homepage
OpenPath [301]	Multimodel	208,414 pathology images related descriptions	Homepage

✧ Although there are datasets available for Instruction Fine-Tuning, such as MultiMedQA and the USMLE test, we have opted not to include them in this list. These datasets are typically employed for evaluation purposes rather than serving as primary resources for Instruction Fine-Tuning.

efficiency. The inclusion of this feature holds significance not only in the deployment of Healthcare LLMs but also in providing substantial support to diverse Healthcare mobile devices equipped with AI cores. For example, the study [217] utilized 8-bit matrix multiplication combined with an 8-bit optimizer for the feed-forward and attention projection layers. This approach enabled the utilization of OPT-175B/BLOOM on a single server equipped with consumer GPUs.

E. Useful Resources

1) *OpenBMB*: OpenBMB (Open Lab for Big Model Base)¹¹ founded by TsinghuaNLP & ModelBest Inc, aiming to build foundation models and systems towards AGI. OpenBMB has published four main projects, namely CPM-Live, BMInf, BMTrain, and BMCook. CPM-Live¹² focuses on live training LLMs and includes three LLM training plans as milestones: CPM-Ant, CPM-Ant+, and CPM-Bee. The project provides real-time recording of training dynamics on the official website, which includes information such as loss function, learning rate, learned data, throughput, gradient size, cost curve, and mean and standard deviation of internal model parameters. This enables users to diagnose any issues during the training process more efficiently.

BMCook and BMInf toolkits allow users to utilize LLMs (specifically CPM-Ant) with limited computational resources. With BMInf, users can perform large model inference on a single GPU, including consumer graphics cards like the GTX 1060, replacing the need for a compute cluster. The compressed models (7B, 3B, 1B, 300M) provided by these toolkits can be adapted to various low-resource scenarios.

The BMTrain toolkit facilitates the efficient training of large models by leveraging distributed computing resources. The

¹¹<https://github.com/OpenBMB>

¹²<https://openi.pcl.ac.cn/OpenBMB/CPM-Live>

¹⁰https://huggingface.co/docs/transformers/main/main_classes/deepspeed

training of CPM-Ant took 68 days and cost 430,000 RMB, which is only 1/20th of Google’s cost for training the T5-11B model, estimated to be around \$130 million. Besides, OpenBMB’s solution contributes to approximately 1/10th of the carbon dioxide emissions compared to training the T5-11B model.

2) *DeepSpeed Chat*: DeepSpeed Chat [302] is a fast, affordable, and scalable open-source framework that enables end-to-end RLHF training to generate various chatGPT-like models. There are four core features of DeepSpeed Chat, including easy-breezy training, a high-performance system, accessible LLM support, and universal acceleration backed for RLHF. DeepSpeed Chat implements the pattern of Instruct-GPT [73], which includes SFT, training the reward model, and the final reinforcement learning-based tuning. By using DeepSpeed Chat, the above process can be easily achieved.

Furthermore, DeepSpeed Chat incorporates a DeepSpeed Hybrid Engine, facilitating a smooth transition of RLHF between inference and training phases. This functionality optimally harnesses a spectrum of optimizations tailored for either training or inference processes.

3) *Training Data*: as mentioned earlier, the transition from PLMs to LLMs brings a significant shift from a model-centered approach to a data-centered approach. Increasing the volume of pre-training data has become a key factor in enhancing the general capabilities of LLMs. In line with this, we have gathered and organized various datasets for training Healthcare LLMs, as presented in Table VII. Besides the medical training data, we also list three Github projects which integrate many general instruction and RLHF training data, including Awesome Instruction Datasets¹³, Awesome-text/visual-instruction-tuning-dataset¹⁴, and Awesome-instruction-tuning¹⁵. Our aim is to assist those interested in training or fine-tuning Healthcare LLMs in easily identifying the appropriate datasets.

In general, the most common sources of data for Healthcare LLMs include EHR, scientific literature, web data, and public knowledge bases. When considering the data structure, QA and dialogue data are the most frequently encountered. Additionally, apart from the conventional text data used in LLMs, it is crucial to acknowledge the significance of multimodal data. Given that the Healthcare domain inherently involves text, images, and time series data, multimodal LLMs offer a promising direction for further research. We anticipate that multimodal LLMs will receive expedited attention in future studies. Following, we briefly introduce some representative data set to provide a general view.

EHR. The Medical Information Mart for Intensive Care III dataset (MIMIC III) is widely recognized as one of the most widely used EHR datasets. It encompasses a comprehensive collection of data from 58,976 unique hospital admissions involving 38,597 patients who were treated in the intensive care unit at the Beth Israel Deaconess Medical Center between 2001 and 2012. Furthermore, the dataset includes 2,083,180 de-identified notes that are associated with these admissions.

MIMIC III provides valuable and extensive information for research and analysis in the field of Healthcare, which facilitates many PLMs and LLMs developments, such as MIMIC-BERT [131], GatorTron [181], and MedAGI [192].

Scientific Literature. PubMed is a freely accessible search engine that provides access to the MEDLINE database, which contains references and abstracts related to life sciences and biomedical topics. It serves as a comprehensive resource with over 32 million citations for biomedical literature, including content from MEDLINE, life science journals, and online books. These citations may also include links to full-text content available on PubMed Central and publisher websites. The PubMed abstracts alone contain approximately 4.5 billion words, while the full-text articles available on PubMed Central (PMC) contribute around 13.5 billion words. These datasets consist of high-quality academic and professional text, making them particularly suitable for training Healthcare LLMs. Various PLM and LLM models, such as BioBERT [91], BioELECTRA [303], GatorTron [181], and MedAlpaca [186], have been trained using PubMed data. PubMed’s vast collection of biomedical literature serves as a valuable foundation for advancing research and development in the Healthcare domain.

Web Data. Web data includes any text we can obtain from the Internet. Social media is one of the most commonly used data types. Reddit is a popular online platform that combines social news aggregation, content rating, and discussion features. Users can contribute various types of content, including links, text posts, images, and videos. The platform is organized into user-created boards called “communities” or “sub-reddits”, covering a broad range of topics. Popular posts with more up-votes rise to the top of their respective sub-reddits and can even make it to the site’s front page. Overall, Reddit offers a diverse and dynamic space for users to engage in discussions, share content, and explore a wide range of interests. The study [304] crawled health-themed forums on Reddit to form COMETA corpus as LLMs training data. Tweets are also usually employed to collect data, and COVID-twitter-BERT [140], Twitter BERT [305], and TwHIN-BERT [306] are trained with these data.

Public Knowledge Bases. There exist many Healthcare-related knowledge bases, such as UMLS [307], CMeKG [308], BioModels [309], and DrugBank [310]. Among them, UMLS is one of the most popular, which is a repository of biomedical vocabularies developed by the US National Library of Medicine. The UMLS has over 2 million names for 900,000 concepts from more than 60 families of biomedical vocabularies, as well as 12 million relations among these concepts. Based on this structured data, the US Medical Licensing Examination (USMLE) is organized and usually employed to test Healthcare LLMs.

CMeKG [308] is a Chinese medical knowledge graph that has been constructed by referring to authoritative international medical standards and a wide range of sources, including clinical guidelines, industry standards, medical textbooks, and other medical texts. This knowledge graph serves as a comprehensive resource for medical information. Building upon the CMeKG, HuaTuo [187] utilizes diverse instructional data for its instruction tuning process. HuaTuo employs Reinforcement

¹³<https://github.com/jianzhnie/awesome-instruction-datasets>

¹⁴<https://github.com/yaodongC/awesome-instruction-dataset>

¹⁵<https://github.com/zhilizju/Awesome-instruction-tuning>

Instruction: " In the following paragraph, several chemicals and diseases are mentioned. Your task is to identify the chemical-disease relations and present them in the format of (Subject, Object). "

Input: " To test the validity of hypothesis that hypomethylation of DNA plays an important role in the initiation of carcinogenic process, 5-azacytidine (5-AzC) (10 mg/kg), an inhibitor of DNA methylation, was given to rats during the phase of repair synthesis induced by the three carcinogens, benzo[a]-pyrene (200 mg/kg), N-methyl-N-nitrosourea (60 mg/kg) and 1,2-dimethylhydrazine (1,2-DMH) (100 mg/kg). "

Output: " (benzo[a]-pyrene, initiation of carcinogenic process),
 (5-AzC, initiation of carcinogenic process),
 (1,2-DMH, initiation of carcinogenic process),
 (N-methyl-N-nitrosourea, initiation of carcinogenic process) "

Fig. 11. An example of instruction instance. The instance comes from the study [311].

TABLE VIII
 THE STATISTICS OF COMPUTATION COST FOR EXISTING HEALTHCARE LLM.

Model Name	Total data size	epoch	Batch size	GPU type	GPU number	GPU time
Visual Med-Alpaca	54k data points	3	128	A100-80G	4	2.51 hours
GatorTron	>90 billion words	10	-	A100	992	6 days
Galactica	-	-	-	A100-80G	128	-
ChatDoctor	100k conversations	3	192	A100	6	3 hours
DoctorGLM	3.5G	1	4	A100-80G	1	8 hours
PMC-LLaMA	75B tokens	5	128	A100	8	7 days
Visual Med-Alpaca	44.8MB* (without images)	-	128	A100-80G	4	2.51 hours
BianQue 1.0	9 million samples	1	-	RTX 4090	8	16 days
GatorTronGPT	277B tokens	-	1,120/560	A100-80G	560	26 days
HuatuoGPT	226,042 instances	3	128	A100	8	-
LLaVA-Med	15 million figure-caption pairs	-	-	A100	8	15 hours
Med-Flamingo	1.3M image-caption pairs	-	400	A100-80G	8	6.75 days

TABLE IX
 ESTIMATED FLOPS AND TRAINING TOKENS FOR DIFFERENT MODEL SIZES.

Parameters	FLOPs	FLOPs (in Gopher unit)	Tokens
400 Million	1.92e+19	1/29, 968	8.0 Billion
1 Billion	1.21e+20	1/4, 761	20.2 Billion
10 Billion	1.23e+22	1/46	205.1 Billion
67 Billion	5.76e+23	1	1.5 Trillion
175 Billion	3.85e+24	6.7	3.7 Trillion
280 Billion	9.90e+24	17.2	5.9 Trillion
520 Billion	3.43e+25	59.5	11.0 Trillion
1 Trillion	1.27e+26	221.3	21.2 Trillion
10 Trillion	1.30e+28	22515.9	216.2 Trillion

★This estimation comes from the study [37]★. Gopher is another LLM study [312] used to compare.

Learning from AI Feedback (RLAIF) to refine its instructions and enhance its performance. The combination of CMeKG and HuaTuo demonstrates the application of public knowledge bases for developing Healthcare LLMs. More details can be seen in Section IV-B3.

Data for Instruction Fine-Tuning. The aforementioned data typically consists of general text that is commonly used for pretraining PLMs or LLMs. However, when transitioning from PLMs to LLMs, instruction data becomes crucial in order to equip LLMs with the capability of following instructions effectively. Unlike PLMs, which primarily focus on next-word prediction, LLMs place greater emphasis on responding to specific instructions.

To illustrate, an instruction instance is presented in Figure 11. In this example, the LLM is tasked with identifying

chemical-disease relations and understanding that its response should align with the given instruction, rather than predicting the next word. By leveraging a sufficient amount of instruction data for fine-tuning, an LLM can appropriately generate the desired output, as demonstrated in Figure 11. This emphasizes the importance of instruction-based training for LLMs to achieve accurate and contextually relevant responses.

4) *Summary:* in Section IV-D and Section IV-E, we present a comprehensive overview of two fundamental resources crucial for LLMs – the training framework and data. Specifically, Section IV-D2 highlights compute-efficient and memory-efficient methods, such as Parallelism, ZeRO, and Quantization, that have been proven to substantially reduce the overall cost associated with LLM training or fine-tuning. These cutting-edge technologies hold significant value as they effectively lower the entry barrier for researchers and practitioners interested in exploring the realm of LLMs. Building upon these advancements, numerous training frameworks have emerged, offering integrated solutions that encompass a wide range of acceleration techniques and enhanced support, thereby facilitating more convenient and streamlined LLM training processes.

The choice of a suitable training framework holds great significance in accelerating the development of LLMs. Among the available options, DeepSpeed-chat is a training framework proposed by Microsoft, which has gained recognition. Additionally, PyTorch’s official distributed training tool, Accelerate, offers stability and ease of use for small to medium-sized training tasks. Another noteworthy open-source LLM training framework is veGiantModel, developed by ByteDance, which provides valuable support in LLM development. You can find more information about veGiantModel at their GitHub repository¹⁶.

When it comes to the data used for training LLMs, the volume often surpasses the capacity of human teams to manually perform quality checks. Consequently, data collection processes heavily rely on heuristic rules for selecting data sources and applying filters. In the context of LLM training,

¹⁶<https://github.com/volcengine/veGiantModel>

TABLE X
THE GENERAL EVALUATION OF LLMs.

Categories	Studies	Evaluation Tasks and Conclusions
Content generation	[313]	Question answering. ChatGPT is a knowledgeable but inexperienced solver, and GPTs can achieve good accuracy in commonsense question answering.
	[314]	Affective computing. ChatGPT is a good generalist model without any specialized training.
	[315]	Text summarization. ChatGPT occasionally produces longer summaries than the input, and the fine-tuned BART performs better than the zero-shot ChatGPT.
	[316]	Dialogue. Claude and ChatGPT are optimized for chat applications.
	[317]	Translation. ChatGPT and GPT-4 achieve superior performance compared to commercial translation systems.
Logical reasoning	[318]	Code synthesis. ChatGPT and GPT-4 achieve superior performance.
	[319]	Mathematical reasoning. ChatGPT performed much worse than high-ability students.
	[320]	Deductive, inductive and abductive reasoning. LLMs perform worse than the fine-tuned state-of-the-art model.

☆ The general evaluation mainly contains two categories: content generation and logical reasoning.

TABLE XI
THE HEALTHCARE EVALUATION OF LLMs.

Categories	Studies	Models	Scenarios	#Num	Conclusions
Medical Ex.	[321]	ChatGPT	Primary Care	674	Average performance of ChatGPT is below the mean passing mark in the last 2 years.
	[322]	ChatGPT	Medical licensure	220	ChatGPT performs at the level of a third-year medical student.
	[323]	ChatGPT	Medical licensure	376	ChatGPT performs at or near the passing threshold.
Medical Q&A.	[324]	ChatGPT	Physician queries	284	ChatGPT generates largely accurate information to diverse medical queries.
	[325]	ChatGPT, GPT-4, Bard, BLOOMZ	Radiation oncology	100	Each LLM generally outperforms the non-expert humans, while only GPT-4 outperforms the medical physicians.
	[103]	ChatGPT, Claude	Patient-specific EHR	–	Both models are able to provide accurate, relevant, and comprehensive answers.
	[326]	ChatGPT	Bariatric surgery	151	ChatGPT usually provides accurate and reproducible responses to common questions related to bariatric surgery.
	[327]	ChatGPT	Genetics questions	85	ChatGPT does not perform significantly differently than human respondents.
	[328]	ChatGPT	Fertility counseling	17	ChatGPT could produce relevant, meaningful responses to fertility-related clinical queries.
Medical Gen.	[329]	GPT-3.5, GPT-4	General surgery	280	GPT-3.5 and, in particular, GPT-4 exhibit a remarkable ability to understand complex surgical clinical information.
	[330]	GPT-3.5, GPT-4	Dementia diagnosis	981	GPT-3.5 and GPT-4 cannot outperform traditional AI tools in dementia diagnosis and prediction tasks.
	[331]	ChatGPT	Gastroenterology	20	ChatGPT would generate relevant and clear research questions, but not original.
Medical Ce.	[332]	ChatGPT, GPT-4	Radiology report	138	ChatGPT performs well and GPT-4 can significantly improve the quality.
	[333]	ChatGPT	Benchmark tasks	34.4K	Zero-shot ChatGPT outperforms the state-of-the-art fine-tuned models in datasets that have smaller training sets.
	[334]	ChatGPT	Clinical and research	–	ChatGPT could potentially exhibit biases or be susceptible to misuse.

☆ The Healthcare evaluation of LLMs includes Medical examination (Ex.), medical question answering (Q&A), medical generation (Gen.), and medical comprehensive evaluation (Ce.).

there are various data challenges to address, including the high cost of Healthcare data, near-duplicates, contamination in benchmark data, personally identifiable information, and the mixture of domains during pre-training and fine-tuning tasks.

Based on the above information, one of the primary concerns in developing an LLM – the computational cost, is involved. By considering the training framework, data requirements, and the size of the LLM itself, an estimation of the overall computational cost can be obtained. We have summarized the relevant computation costs from existing studies in Table VIII. Table IX comes from the study [37], which estimates the relation among the model size, the dataset size, and the training FLOPs when we need to train an LLM from scratch. These data can serve as a helpful reference for those seeking to estimate the expenses associated with LLM development.

V. EVALUATION METHOD

Presently, there is a wide range of LLMs available for general NLP tasks and Healthcare applications. Selecting the appropriate model as a benchmark for intelligent applications is of utmost importance. Consequently, evaluating the performance of LLMs holds significant value for both the NLP and

Healthcare communities. According to the survey from the study [15], most LLMs were evaluated based on downstream tasks, where the tasks can be categorized as testing language and reasoning ability, and scientific knowledge. They found that LLMs present notable proficiency in comprehending and generating human language, facilitating interactive exchanges with users through dialogues. This enables them to effectively address a wide array of Natural Language Processing (NLP) tasks and furnish elucidative responses. Nevertheless, it is important to note that their present capabilities do not categorize them as comprehensive AI systems. They still face performance limitations, particularly when compared to expert models, across multiple domains that necessitate domain-specific knowledge. On the other hand, the state-of-the-art LLMs demonstrate commendable performance in grasping general scientific knowledge and are capable of generating open-ended responses to science-related inquiries. Nevertheless, they are susceptible to errors, particularly when tackling questions that necessitate intricate multi-step reasoning. The exceptional proficiency in language presents a hurdle for users to accurately evaluate the factual correctness of information, thereby giving rise to a spectrum of ethical considerations.

In this section, we will begin by introducing studies on

the evaluation of general NLP tasks. Subsequently, we will review studies focusing on Healthcare evaluation, discussing aspects such as robustness, bias, and ethics. Finally, we will conclude by highlighting future directions for health evaluation and providing a summary.

A. General NLP tasks Evaluation

To provide a comprehensive exposition of LLM evaluation studies in NLP tasks, we propose two evaluation categories in response to the need for enhanced intelligence, namely content generation and logical reasoning. The typical evaluation studies and their main conclusions are summarized in Table X. The evaluation study for content generation is to generate answers for general NLP tasks, such as question answering, affective computing, text summarization, dialogue, translation, and code synthesis. The generation form typically manifests as either natural language or code. Overall, the results demonstrate that LLMs have made significant advancements in such general NLP tasks.

However, the question remains whether LLMs can still achieve superior performance, considering that complex logical reasoning tasks demand planning, abstraction, and inference abilities. Under this circumstance, several studies for evaluating logical reasoning abilities are proposed. From the reasoning form, they can be divided into deductive, inductive, abductive, and mathematical views. The relevant results all show that LLMs do not perform well in logical reasoning. For example, ChatGPT incorrectly answers almost all questions about probability and statistics, permutation and combination, and geometry [319]. The overall performance is much worse than high-ability students for mathematical reasoning. Meanwhile, the study [320] provides a comprehensive evaluation of LLMs for deductive, inductive and abductive reasoning, and demonstrates LLMs perform worse than the fine-tuned state-of-the-art model. Meanwhile, LLMs often exhibit noticeable logical flaws and hallucinations, which pose significant challenges for the practical application of LLMs in logical reasoning scenarios.

B. Healthcare Evaluation

Different from general NLP tasks, the field of Healthcare is characterized by its high level of specialization. Evaluating LLMs in this domain necessitates assessing their capacity to comprehend and utilize medical knowledge and terminology. The evaluation process may involve designing test cases tailored to specific tasks and challenges within the medical field. According to the different forms of evaluation, we categorize the current relevant work into four folds: medical examination, medical question answering, medical generation, and medical comprehensive evaluation, which are summarized in Table XI. The medical examination form involves verifying model performance through standard medical tests or examinations. Differently, medical question answering involves utilizing questions posed or collected by human experts to make assessments. Medical generation focuses on generating new medical descriptions or knowledge based on a given input. The studies on medical comprehensive evaluation aim

to provide assessments across various application scenarios rather than focusing on a single aspect.

In the form of medical examination, the study [321] evaluated the strengths and weaknesses of ChatGPT in primary care using the Membership of the Royal College of General Practitioners Applied Knowledge Test (AKT). It is observed that ChatGPT's average performance (60.17%) is below the mean passing mark in the last 2 years (70.42%), demonstrating further development is required to match the performance of qualified primary care physicians. The study [322] evaluated ChatGPT's performance on the medical licensing exams utilizing AMBOSS¹⁷ and the National Board of Medical Examiners (NBME), which shows that ChatGPT performs at the level of a third-year medical student on the question sets examined and its responses to questions provide interpretable context to justify models written response in most cases. Similarly, the study [323] tested the performance characteristics of ChatGPT on USMLE. They certified that ChatGPT is able to perform several intricate tasks relevant to handling complex medical and clinical information, as ChatGPT performed at or near the passing threshold of 60% accuracy.

The ability of QA plays a crucial role in the application of models, and as a result, there have been numerous studies focusing on the evaluation of medical question answering. To explore the accuracy and completeness of ChatGPT for medical queries, the study [324] collected 284 medical questions from 33 physicians across 17 specialties. After that, these physicians graded ChatGPT-generated answers to these questions for accuracy, showing that ChatGPT achieved relatively high accuracy and completeness scores. Utilizing 100 multiple-choice questions on radiation oncology physics created by an experienced medical physicist, the study [325] investigated LLMs' capacity in answering radiation oncology physics questions. Four LLMs (ChatGPT, GPT-4, Bard¹⁸, and BLOOMZ¹⁹) are utilized to compare with medical physicists and non-experts. The results demonstrate that all these LLMs generally outperform the non-expert humans and only GPT-4 outperforms the medical physicists. However, it is not allowed for GPT-4 to improve performance when scoring based on a majority vote across trials, while a team of medical physicists is able to greatly outperform GPT-4 using a majority vote. The study [325] investigated the use of LLMs (ChatGPT and Claude) for patient-specific question answering from EHRs. On both experiment settings of one question per session and one topic per session, ChatGPT and Claude are able to provide accurate, relevant, and comprehensive answers to general questions, specific questions, and nonanswerable questions. To examine the accuracy and reproducibility of LLMs in answering patient questions regarding bariatric surgery, the study [326] gathered 151 questions from nationally regarded professional societies and health institutions as well as Facebook support groups. Using ChatGPT, accurate and reproducible responses to common questions could be provided. The study [327] assessed ChatGPT in the field of genetics involving 13,636

¹⁷<https://www.amboss.com/>

¹⁸<https://bard.google.com/>

¹⁹<https://github.com/bigscience-workshop/xmft>

responses to 85 questions. Although ChatGPT is significantly better on memorization-type questions versus critical-thinking questions, it does perform significantly differently than human respondents. And it would generate plausible explanations for both correct and incorrect answers. The study [328] explored ChatGPT's ability for fertility counseling and only 6.12% ChatGPT factual statements were categorized as incorrect. The study showcases the capacity of LLMs to generate pertinent and meaningful responses to clinical queries related to fertility. Nevertheless, there are certain limitations to consider, including the challenges in providing reliable source citations and the unpredictable potential for generating fabricated information. The study [329] aimed to assess the performance of GPT-3.5 and GPT-4 in understanding complex surgical clinical information and its potential implications for surgical education and training. GPT-3.5 achieved an overall accuracy of 46.8%, while GPT-4 demonstrated a significant improvement with an overall accuracy of 76.4%. GPT-3.5 and GPT-4 specifically have a remarkable ability to understand complex surgical clinical information. The study [330] explored the potential of GPT-3.5 and GPT-4 to surpass traditional AI tools in dementia diagnosis. The experimental results, obtained from two real clinical datasets, indicate that while LLMs show promise for future advancements in dementia diagnosis. They currently do not outperform traditional AI tools in terms of performance.

The evaluation of medical generation can provide further insights into the level of control that LLMs have over medical knowledge. It is significant to pinpoint the most pressing and important research questions. To this end, the study [331] evaluated the potential of chatGPT for identifying research priorities in gastroenterology from four key topics. Several experienced experts reviewed and rated the generated research questions. It seems ChatGPT would generate relevant and clear research questions. However, the generated questions were not considered original. The study [332] investigated the feasibility of using ChatGPT and GPT-4 to translate radiology reports into plain language. According to the evaluation by radiologists, ChatGPT performs well and can successfully translate radiology reports into plain language with an average score of 4.27 in the five-point system. Further, GPT-4 can significantly improve the quality of translated reports.

Several studies evaluate the comprehensive capability of LLMs. For example, the study [333] provides a comprehensive evaluation of ChatGPT's zero-shot performance on various benchmark biomedical tasks, i.e., relation extraction, document classification, question answering, and summarization. Zero-shot ChatGPT achieves comparable performance to fine-tuned generative transformers such as BioGPT and BioBART. Additionally, when evaluated on datasets with limited training data, zero-shot ChatGPT outperforms these fine-tuned models. These results indicate that ChatGPT exhibits a high degree of specialization even within the biomedical domain. The study [334] conducted a concise investigation to assess the potential applications of ChatGPT in four clinical and research scenarios: support of clinical practice, scientific production, misuse in medicine and research, and reasoning about public health topics. The study draws the following conclusions: ChatGPT demonstrates the capability to offer valuable sug-

gestions; it accurately identifies the context and summarizes findings; potential misuse is identified; and ChatGPT exhibits significant potential for expediting scientific progress.

C. Evaluation of Robustness, Bias, and Ethics

To assess how well a model performs when faced with uncertainties, perturbations, or unexpected inputs, researchers have been studying robustness evaluation techniques. For instance, in the field of general NLP tasks, studies have explored the robustness of LLMs in areas such as semantic parsing [335] and vision-language tasks [336]. In the Healthcare domain, the evaluation of LLMs' robustness is relatively limited. One notable example is the evaluation of ChatGPT's robustness in translating radiology reports [332]. In this work, the original radiology reports were divided into 25 key information points, and the correctness and completeness of each point were evaluated in a point-by-point manner in the translated reports. The overall translation quality was found to be satisfactory for only 55.2% of the translated points, indicating ample room for improvement in the robustness of LLMs in Healthcare settings.

LLMs are generated through training on extensive text datasets, which can inherently contain various biases and imbalances. When the model is consistently exposed to specific biases or particular points of view during training, it tends to learn and reflect those biases, leading to biased outputs during text generation. In the manual evaluation process, the presence of biases can also arise due to the diverse academic backgrounds and perspectives of the experts involved. Each expert may have their own subjective interpretation or evaluation criteria, which can introduce deviations in the evaluation results [323].

Furthermore, during the evaluation process, LLMs may require the uploading of personal privacy data, such as patient-specific EHR [325]. This introduces a significant privacy risk that demands careful attention. Consequently, ethic issues related to data privacy and protection [324], [333] must be thoroughly considered in the evaluation

D. Future Directions for Health Evaluation

The study [15] found that present evaluation methodologies heavily rely on prompt engineering and established benchmark datasets. Different prompt formulations can lead to contrasting evaluation outcomes. Furthermore, the assessment of expert systems frequently hinges on utilizing (in-domain) datasets that were originally employed for training those systems. An ambiguity persists regarding potential inadvertent exposure of the scrutinized data, such as publicly available datasets and established scientific knowledge, during the training of Large Language Models (LLMs). These aspects could introduce bias into the comparison between LLMs and their corresponding baselines, impeding a fair assessment.

According to the current studies of Healthcare evaluation, we conclude the following four future directions.

Increase the evaluation of faithfulness. Healthcare professionals and patients place significant trust in the accuracy and reliability of information provided by LLMs. However, due to

the unique nature of the medical domain, there is a risk that LLMs may generate false knowledge or hallucinations, which could potentially lead to serious accidents or harm. Therefore, evaluating the faithfulness of LLMs becomes crucial in order to identify instances where these models may generate hallucinations and mitigate their impact.

Towards comprehensive and multitask evaluation. The current evaluation practices predominantly concentrate on assessing the performance of LLMs on one specific medical task, which might not provide a comprehensive understanding of their capabilities across the entire medical applications. Consequently, there is a clear need for a multitask evaluation system that can comprehensively evaluate the performance of LLMs across various medical tasks.

Towards multi-dimensional evaluation. While current evaluation efforts have primarily centered around accuracy, there is a growing recognition of the need for a multidimensional evaluation framework. It should consider various aspects beyond accuracy, such as the correctness of interpretation, robustness, hallucination ratio, content redundancy, biased description, and ICL capability.

Increase privacy protection in the evaluation process. Medical applications inherently involve sensitive data privacy concerns that surpass those of other NLP tasks. Consequently, safeguarding privacy during the evaluation process becomes of utmost importance. One potential solution to address this challenge is the adoption of federated learning approaches [337], which enable the implementation of large-scale evaluation systems while preserving privacy.

E. Summary

In conclusion, while LLMs demonstrate strong performance in general NLP tasks, they often fall short when it comes to tackling complex logical reasoning problems. For Healthcare evaluation, LLMs tend to perform below or just meet the threshold in medical examination scenarios. For the medical question answering part that has the most studies, LLMs exhibit underperformance in genetics questions and dementia diagnosis. For medical generation and comprehensive evaluation, LLMs usually perform well and have a positive impact despite the existence of non-original generations, biases, or instances of misuse. From these evaluation studies, we have discovered that LLMs hold significant potential for various applications in the health field. However, there are several pressing issues that need to be addressed to enhance their utilization in this domain.

VI. IMPROVING FAIRNESS, ACCOUNTABILITY, TRANSPARENCY, AND ETHICS

Fairness, accountability, transparency, and ethics are four important concerns in the AI domain. According to the study [338], fairness holds paramount significance in guaranteeing that AI does not perpetuate or exacerbate established societal disparities; Accountability plays an important role in ensuring that individuals responsible for the conception and execution of AI can be held answerable for their decisions; Transparency assumes a critical role in ensuring that AI

remains open to scrutiny and amenable to audits for possible biases or inaccuracies; Ethics, similarly, assumes a pivotal role in guaranteeing that AI is constructed and utilized in manners that align with prevailing social values and norms.

In the Healthcare domain, we believe that these four aspects are even more critical because the primary focus is on patient well-being and safety. In this context, the utmost importance lies in ensuring patients receive optimal Healthcare marked by equitable access to medical services. Additionally, the transparent and trustworthy nature of Healthcare decisions, the accountability in delivering accurate medical diagnoses and treatments, the safeguarding of patient confidentiality, and the adherence to elevated ethical standards emerge as distinct and noteworthy considerations, setting Healthcare apart from AI applications in other domains and more.

In the following subsections, we will survey the common fairness, accountability, transparency, and ethics issues related to using AI for Healthcare. Then, we will propose possible mitigation for these issues.

A. Fairness

Fairness within the context of LLMs and NLP refers to the principle of equitably treating all users and preventing any form of unjust discrimination. This essential concept revolves around the mitigation of biases, aiming to guarantee that the outcomes produced by an AI system do not provide undue advantages or disadvantages to specific individuals or groups. These determinations should not be influenced by factors such as race, gender, socioeconomic status [17], or any other related attributes, e.g., different input languages [339] and processing tasks [340], striving for an impartial and balanced treatment of all users. This fundamental tenet aligns with the broader objective of promoting equality and inclusivity within the applications of LLMs and NLP.

In an empirical study, the study [340] found that PLMs may generate biased outcomes given different tasks, prompts, and label word selection methods. They evaluated both small and large versions of four PLMs, showing that PLMs can yield huge accuracy gaps in sentiment analysis and emotion detection tasks, even though the prompts, label word selection, and input text have been well controlled. This finding goes against human intuition because cognitively, sentiment and emotions are divisions of subjective expressions into different granularities. Polarized sentiment can be thought of as a 2-dimensional summary of positive and negative emotions. The accuracy gap between sentiment and emotion classification tasks indicates that the performance of PLMs can be significantly impacted by how the label space is divided in a specific task. For sentiment analysis whose label space is evenly divided into two classes, e.g., positive and negative classes, PLMs tend to achieve better performance than the tasks whose label space is unevenly divided, e.g., angry, fear, sad, and joyful emotions.

The study [39] noticed that when comparing sentiment scores, the fine-tuned LLaMA 2-Chat exhibits a more positive sentiment compared to pretrained versions, whereas ChatGPT tends to generate responses with a more neutral sentiment.

In terms of gender, LLMs tend to express a more positive sentiment towards American female actresses than male actors. Regarding race, Asian Americans and Hispanic/Latino Americans tend to have relatively higher sentiment scores compared to other racial subgroups. In the religious ideology domain, Islam and Sikhism groups display the most significant increase in sentiment scores after fine-tuning. In the political ideology domain, both Liberalism and Conservatism groups tend to have the most positive sentiment scores, while Fascism group scores are predominantly negative. Lastly, in the profession domain, there is notably positive sentiment towards occupational categories like “Corporate titles” and “Computer”, while sentiment is most neutral towards “Professional driver types”.

The biases from LLMs can be attributed to the uneven distribution of demographic attributes in pre-training corpora [39]. Such an argument also holds for the Healthcare sector [341]. As an example, CNNs trained on publicly accessible chest X-ray datasets tend to exhibit underdiagnosis tendencies in marginalized communities, including female patients, Black patients, Hispanic patients, and those covered by Medicaid insurance [342]. These specific patient groups often experience systemic underrepresentation within the datasets, resulting in biased algorithms that may be susceptible to shifts in population demographics and disease prevalence. Furthermore, several global disease classification systems display limited intra-observer consensus, implying that an algorithm trained and assessed in one country may undergo evaluation under a dissimilar labeling framework in another country [343], [344].

Current common practices to improve AI fairness in the Healthcare domain focus on pre-processing, in-processing, and post-processing [341]. Importance weighting is a pre-processing technique, which involves adjusting the significance of less frequent samples from protected subgroups. Similarly, resampling endeavors to rectify sample-selection bias by acquiring more equitable subsets of the initial training dataset and can be naturally employed to address the underrepresentation of specific subgroups. In the case of tabular-structured data, methods like blinding, data transformation, and others can be utilized to directly remove proxy variables that encode protected attributes.

To alleviate the impact of confounding variables, an anti-discrimination component (an in-processing technique) can be integrated into the model to discourage the learning of discriminatory attributes related to a protected attribute. For example, in the case of a logistic regression model, modifications can be made to include anti-discrimination elements by evaluating the covariance between the protected attribute and the signed distance from the sample’s feature vectors to the decision boundary. Another approach involves adjusting the parameters of the decision boundary to enhance fairness (by minimizing disparate impact or mistreatment), while still adhering to accuracy constraints [345]. Deep learning models, e.g., CNNs with adversarial-loss terms, render the internal feature representations invariant to variations in protected subgroups [346]. Tuning loss weights for different classes is also common for unbalanced label learning [347]. Furthermore, adjustments to the stochastic gradient descent technique can be implemented to incorporate fairness constraints within online

learning frameworks [348].

Post-processing involves methodologies that alter the output of a trained model, such as probability scores or decision thresholds, to adhere to group fairness criteria. In aiming for equalized odds, one approach is to establish appropriate thresholds for each group, ensuring that the model attains a consistent operating point across all groups. However, in scenarios where the receiver operating curves do not intersect or where the desired operating point does not align with an intersection point, implementing this strategy necessitates deliberately degrading performance for specific subgroups using a randomized decision rule. Essentially, this signifies that the model’s performance for certain groups may have to be intentionally diminished to fulfill the criteria of equalized odds.

For LLMs, bias mitigation methods are frequently studied in the context of instruction fine-tuning and prompt engineering [349]. The representative technique for instruction fine-tuning is RLHF. In the case of InstructGPT, GPT-3 is refined through a process involving RLHF, specifically aimed at adhering to human instructions. The procedure involves three sequential steps: firstly, gathering human-authored demonstration data to guide GPT-3’s learning; secondly, assembling comparative data consisting of model-generated outputs assessed by annotators to construct a reward model that predicts outputs preferred by humans; and lastly, fine-tuning policies based on this reward model using the Proximal Policy Optimization algorithm [350].

B. Accountability

LLMs possess the propensity to magnify the inherent social biases embedded within their training data because they can generate hallucinatory or counterfactual information and present a deficiency in robustness. These limits imply that LLMs are susceptible to perturbations and deviations from their intended performance, particularly when exposed to diverse inputs or scenarios. Thus, ensuring accountability emerges as a pivotal concern when integrating LLMs within the Healthcare domain.

The study [351] identify two primary factors that significantly contribute to the performance and instances of hallucination in generative Large Language Models (LLMs). Firstly, a major influence is attributed to the model’s memorization of the training data. Additionally, the authors demonstrate that named entity IDs serve as “indices” for accessing the memorized data. Secondly, the authors illustrate that LLMs employ an additional heuristic based on corpus-derived patterns involving word frequencies. They provide evidence that NLI test samples deviating from these patterns result in significantly lower scores compared to those adhering to them. Hallucinations are not a unique flaw of LLMs. They are also common in large vision language models (LVLMs). The study [352] further confirmed that objects that have a high frequency of occurrence in the visual instructions or co-occur with the objects present in the image are evidently more susceptible to hallucination by LVLMs.

Generated counterfactual speech presents an additional obstacle to accountable AI. In the evaluation conducted by

the study [353], ChatGPT was evaluated using fact-based question-answering datasets, revealing that its performance did not exhibit enhancements in comparison to earlier versions. Consequently, the reliability of ChatGPT in tasks necessitating faithfulness is called into question. For instance, its potential fabrication of references in the context of scientific article composition [354] and the invention of fictitious legal cases within the legal domain [355] accentuate the potential risks associated with its use in critical domains.

The research [356] uncovered a diverse range of viewpoints among Healthcare professionals regarding the impact of AI on their workload and decision-making processes. Certain practitioners believed that AI had the potential to alleviate their workload, enhance clinical decision-making, and ultimately enhance patient safety by assisting in diagnostics. Conversely, other practitioners voiced apprehensions about a heightened workload, encompassing the effort required to learn and manage the technology alongside patient care, as well as potential risks to patients if the AI system provided unsuitable recommendations. The study also underscored the participants' anxieties about the diminishing interpersonal connection following the integration of AI.

Numerous studies have noted that the apparent scientific style of language used by ChatGPT can mislead human observers regarding the reliability of its outputs [357], [358]. The study [15] contended that enabling users to access human-generated source references is crucial for enhancing the reliability of the model's responses. The study [359] advocated for the involvement of both AI developers and system safety engineers in evaluating the moral accountability concerning patient harm. Additionally, they recommend a transition from a static assurance model to a dynamic one, recognizing that ensuring safety is an ongoing process and cannot be entirely resolved during the initial design phase of the AI system prior to its deployment.

The study [356] proposed a solution to tackle the issue of accountability, advocating for the education and training of prospective AI users to discern the appropriateness of relying on AI recommendations. However, imparting this knowledge to practitioners demands a considerable investment of effort. Healthcare professionals frequently grapple with overwhelming workloads and burnout, making comprehensive training on AI a significant challenge. Moreover, not all Healthcare practitioners possess adequate statistical training to comprehend the underlying mechanics of AI algorithms. In addition to education, the study [356] recommended the establishment of policies and mechanisms to ensure the protection of both clinicians and AI within the Healthcare domain.

C. Transparency

The limited transparency of neural networks has been widely criticized, presenting significant obstacles to their application in the Healthcare domain. LLMs and PLMs are complex neural network models, which further exacerbate the challenges associated with interpretability. In recent years, there have been efforts to understand the inner workings of PLMs in Healthcare contexts. Probing PLMs have been extensively employed to uncover the underlying factors contributing

to their performance [360]. For example, [361] examined PLMs' disease knowledge, while [362] conducted in-depth analyses of attention in protein Transformer models, yielding valuable insights into their mechanisms.

In the general meaning learning domain, a transparent model is typically characterized by decision-making processes akin to those of white-box models, e.g., decision tree-based models or linear regression models. It often encompasses post hoc explanations [363], model-specific explanations [364] or model-agnostic explanations [365]. Sometimes, the explanation insights are derived from feature maps [366], generated natural language [367], factual and counterfactual examples [368], or decision-making evidence [369].

LLMs normally rely on Transformer structures. However, the conventional Transformer is not intrinsically explainable, as it is a stack of multiple layers of multi-head attention, skip connections, and non-linear transformation. The study [370] introduced an innovative approach for assessing relevancy in Transformer networks in the computer vision domain. Their method involves assigning local relevance using the Deep Taylor Decomposition principle and subsequently propagating these relevance scores through the network's layers. This propagation incorporates attention layers and skip connections, introducing a novel challenge compared to existing methods. The authors' solution is grounded in a unique formulation that has demonstrated the ability to preserve the overall relevancy across different layers of the network.

The study [363] introduced an innovative method accompanied by quantitative metrics aimed at mitigating the limitations observed in existing post hoc explanation approaches, as outlined in the literature. These drawbacks include reliance on human judgment, the necessity for retraining, and issues related to data distribution shifts during the occlusion of samples. The method proposed in this study allows for a quantitative assessment of interpretability methods without the need for retraining and effectively addresses distribution shifts between training and evaluation sets. Furthermore, the authors have developed novel metrics and indices to quantitatively evaluate time-series interpretability methods, offering a comprehensive evaluation of how closely an interpretability method aligns with the learned representation of the model in focus.

In the era of LLMs, CoT prompting [33] has emerged as a potential method for providing a certain level of interpretability by generating reasoning steps. The technique empowers LLMs to break down complex, multi-step problems into more manageable intermediate steps. This enables the allocation of additional computational resources to problems demanding deeper reasoning steps. Moreover, it offers a transparent view of the LLM's behavior, shedding light on its potential process of arriving at a specific answer and offering insights for identifying and rectifying errors in the reasoning path. Essentially, a chain of thought can be perceived as a systematic, step-by-step thought process leading to the derivation of an answer. However, this approach faces two primary challenges: the high cost of annotations required for CoT and the evaluation of interpretability. Acquiring demonstrations with annotated reasoning steps is an expensive task, particularly in professional fields such as Healthcare. Additionally, evaluating the

generated reasoning results as explainable justifications and ensuring their usability pose significant challenges.

D. Ethics

The ethical concerns about using LLMs in the Healthcare domain have been widely discussed. The study [29] argued that a primary concern relates to the potential perpetuation of misinformation and biases. Furthermore, the inaccurate results from LLMs inhibit their autonomous deployment, although their utilization in an assisting capacity could significantly enhance efficiency. Domain-specific fine-tuning may enhance their performance, as evidenced by variants like PubMedBERT and BioBERT derived from the BERT model. Addressing accountability issues involves ensuring that clinicians and researchers utilizing these tools take responsibility for the generated output. Lastly, evaluating clinical interventions utilizing LLMs should ideally involve randomized controlled trials to assess their impact on mortality and morbidity. However, determining the appropriate benchmark for such costly and risky trials remains an open question.

Healthcare LLMs typically possess a wide range of patient characteristics, including clinical measurements, molecular signatures, demographic information, and even behavioral and sensory tracking data. It is crucial to acknowledge that these models are susceptible to the problem of memorizing training data and simply reproducing it for users. As mentioned in Section IV-E3, EHRs serve as important training data for Healthcare LLMs, alongside public scientific literature and web data. However, it is worth noting that some EHRs remain private due to organizations' concerns about data exposure. For instance, clinical records may contain sensitive information such as patient visits and medical history, and exposing such data could lead to physical and mental harm to patients. It is important to recognize that de-identification techniques employed in EHR records (e.g., MIMIC III) may not always guarantee complete safety. Recent studies have shown that there can be instances of data leakage from PLMs in the general domain, allowing for the recovery of personal health information from models trained on such data sources [371], [372]. Additionally, approaches such as KART [373] have been proposed to assess the vulnerability of sensitive information in pre-trained biomedical language models using various attack strategies. The Federated Learning [374] is a promising technology to alleviate such a problem.

The study [15] brought attention to a pressing issue regarding the potential misuse of AI-generated content for training subsequent models. They emphasized that due to the probable presence of biases in content generated by LLMs, any sequential models trained using this content could inadvertently inherit and perpetuate those biases. This highlights a significant ethical concern, as the biases within AI-generated content could be inadvertently propagated, emphasizing the critical importance of addressing bias mitigation strategies during training processes involving LLM-generated data.

Complex software designed to aid in the diagnosis, prevention, monitoring, prediction, prognosis, treatment, or alleviation of diseases falls under the classification of a medical

device, necessitating adherence to regulatory controls. These controls encompass the development of these tools within a quality management system to ensure their reliability and safety. Within the medical domain, incorporating LLMs raises notable ethical considerations due to their intricate nature. The study [375] highlighted that in the European Union, stringent post-market surveillance and clinical follow-up are mandated, presenting specific challenges when applied to LLMs. Notably, LLMs lack inherent quality assurance directly from their developers, making their integration as external "plug-in" components of medical devices, for instance, through an API, impractical. This raises ethical concerns, as the use of LLMs in critical medical contexts necessitates stringent quality validation and assurance to ensure patient safety and uphold ethical standards within the Healthcare domain. The study [376] believed that LLMs are frequently trained using extensive data sources, which, given their opaque nature, raises apprehensions regarding potential violations of existing intellectual property rights. This concern, coupled with the limited transparency of current LLMs, has recently influenced amendments to the proposed Artificial Intelligence Act within the European Union. This amendment mandates that companies utilizing generative AI tools must provide disclosures concerning any copyrighted material employed in the development of their systems.

VII. FUTURE WORK AND CONCLUSION

We have discussed fairness, accountability, transparency, and ethics for Healthcare in Section VI. Further from technology aspects, we summarize four points that are the most significant for Healthcare LLMs, including medical knowledge enhancement, integration with the Healthcare process, effective interaction with patients and doctors, hallucinations, misunderstandings, and prompt brittleness.

A. Future Work

1) *Medical knowledge enhancement*: the integration of medical knowledge into PLMs has been a prominent topic of research for several years [160], [377]. In the era of PLMs, there have been attempts to consider PLMs as a form of soft knowledge base capable of capturing knowledge. Studies [378], [379] have explored explicit methods to inject knowledge into PLMs. In the knowledge-intensive Healthcare domain, models infused with medical knowledge hold tremendous potential for future applications. An example study [380] illustrates the integration of domain-specific knowledge (UMLS) into the pre-training process of PLMs. This integration is accomplished through a knowledge augmentation strategy that connects words sharing the same underlying "concept" in UMLS. Furthermore, the study leverages the semantic type knowledge available in UMLS to generate input embeddings with clinical significance. Another study [381] employs mention-neighbour hybrid attention to learn heterogeneous entity information. It infuses the semantic representations of entity types as external medical knowledge, enhancing the model's capabilities. While general PLMs like GPT-4 exhibit significant competence in answering

medical questions, explicit injection of medical knowledge remains a challenge, especially when dealing with smaller-sized PLMs [187], [234]. For both LLMs and PLMs, the most common approach to injecting medical knowledge is fine-tuning the model parameters using medical data. However, a major drawback of this method is that the knowledge remains fixed once training is complete, making it difficult to incorporate specific knowledge or update the overall knowledge without retraining. Retrieval-based LLMs [382] may offer a solution to these challenges, allowing for more flexible and updatable knowledge integration.

2) *Integration with Healthcare process*: is the application of artificial intelligence in the medical field just an “old myth”, or can it really change the status quo? Clearly, although current AI solutions are fragmented and mostly experimental without widespread adoption, there exist such problems because we believe they are mainly caused by the following three reasons based on the existing study [383]. First, it is difficult to integrate with existing hospital information technology (IT) systems. AI solutions require large amounts of data for training, and most of this data is currently stored in hospitals’ own information systems. Retrieving and integrating this data requires upgrades and modifications to existing systems, which will have an impact on hospitals’ daily operations. In addition, different hospitals use different data formats and standards, lack standardized interfaces, and have relatively complex workflows in the Healthcare domain. AI systems find it difficult to adapt to different interfaces, which also increases the difficulty of integration. Second, fragmentation of IT systems due to hospital consolidations. With the increase in hospital mergers and acquisitions, the original hospitals may use completely different IT systems. After consolidation, it is necessary to unify their respective clinical and management systems, which requires huge investment and a long transition period. Introducing new AI systems during this process will face great technical challenges. Third, regulations are unclear and challenging. Currently, laws and regulations for AI medical applications are incomplete. Key issues such as information security, privacy protection, and liability attribution lack clear provisions. In addition, regulations differ across countries and regions. These will bring uncertainties to the development and application of AI systems. At the same time, the application of AI in the medical industry involves complex ethical issues that are also difficult to resolve.

3) *Effective Interaction with Patients and Doctors*: despite the existing fluency of LLMs in human communication, the unique nature of the medical domain necessitates specific requirements for the interaction between LLMs and their users, namely doctors and patients. These requirements include the ability of LLMs to proactively inquire about symptoms, pose targeted questions, and effectively manage the pace and flow of conversations. Additionally, it is desirable for LLMs to perceive and appropriately address patient emotions such as anxiety and fear, thereby providing suitable emotional support. Moreover, an augmentation to the dialogue system could involve incorporating a virtual human design. This design would enable the model to portray a doctor’s image, encompassing elements such as tone, speech speed, and facial expressions,

with the intention of enhancing rapport in communication. Simultaneously, we aim to establish a continuous learning mechanism that enables the model to learn from doctor-patient dialogues and continually enhance its communication capabilities. This entails automatic pragmatic learning, ensuring the model optimizes its ability to effectively communicate with doctors and patients, ultimately leading to more successful interactions.

4) *Hallucinations, Misunderstandings and Prompt Brittleness*: hallucinations, misunderstandings, and prompt brittleness are three fundamental challenges encountered by both general LLMs and Healthcare LLMs. Hallucinations refer to instances where LLMs generate responses that lack coherence or relevance to the given input. These “hallucinations” can pose significant issues, particularly when users are unfamiliar with the discussed concepts, as they may struggle to identify the inaccuracies in the model’s output. Misunderstandings represent a misalignment problem where the behavior of LLMs fails to align with human values, objectives, and expectations. In other words, LLMs may provide incorrect actions or responses despite receiving proper instructions. Prompt brittleness signifies that even minor modifications to the input prompt can yield dramatically different outputs, as first observed in the study by [384]. In the Healthcare context, these issues could lead to unacceptable consequences. While additional instructions or reinforcement learning from human feedback can partially mitigate these challenges, they do not fully satisfy the stability requirements within the Healthcare domain. Regarding prompt brittleness, the current state of prompt engineering heavily relies on extensive experimentation, with a limited theoretical understanding of why a specific phrasing or formulation of a task is more sensible beyond achieving improved empirical results. Consequently, the development of LLMs that exhibit robustness to different prompt styles and formats remains an unsolved problem.

B. Conclusion

Recently, there has been a growing interest in LLMs and their potential applications across various fields. In this study, we aim to provide a comprehensive survey specifically focusing on Healthcare LLMs. Our survey encompasses an extensive examination of data, technologies, applications, fairness, accountability, transparency, ethics, and limitations associated with Healthcare LLMs. A noteworthy transformation has been observed from Discriminative AI to Generative AI, as well as from model-centered to data-centered approaches, marking a significant shift from PLMs to LLMs. This transition has enabled Healthcare LLMs to support more advanced applications beyond conventional NLP-based fundamental tasks. Consequently, the emergence of these advanced applications has inspired numerous related studies.

To facilitate the development of Healthcare LLMs, various instruction datasets and training and inference technologies have been proposed. These resources have played a crucial role in accelerating the progress of LLMs, particularly within the Healthcare domain. Our objective is to summarize these existing resources, providing valuable support to researchers

intending to embark on the development of their own Healthcare LLMs.

However, despite the opportunities presented by Healthcare LLMs, several significant challenges persist, impeding their implementation in Healthcare settings. Issues pertaining to interpretability, privacy protection, medical knowledge enhancement, integration with Healthcare processes, and effective interaction with patients and doctors pose substantial obstacles. These challenges hinder the translation of innovative LLMs into practical adoption within the Healthcare field. Consequently, physicians and other Healthcare professionals must carefully consider the potential benefits and limitations associated with LLMs as they navigate the selection and integration of these models into their medical practice.

ACKNOWLEDGMENTS

This work has been supported by the National Research Foundation Singapore under AI Singapore Programme (Award Number: AISG-GC-2019-001-2A and AISG2-TC-2022-004); The RIE2025 Industry Alignment Fund (I2101E0002 – CISCONUS Accelerated Digital Economy Corporate Laboratory).

REFERENCES

- [1] M. Ge, R. Mao, and E. Cambria, "A survey on computational metaphor processing techniques: From identification, interpretation, generation to application," *Artificial Intelligence Review*, 2023.
- [2] R. Liu, R. Mao, A. T. Luu, and E. Cambria, "A brief survey on advances in coreference resolution," *Artificial Intelligence Review*, 2023.
- [3] R. Mao, K. He, X. Zhang, G. Chen, J. Ni, Z. Yang, and E. Cambria, "A survey on semantic processing techniques," *Information Fusion*, p. 101988, 2023.
- [4] X. Zhang, R. Mao, and E. Cambria, "A survey on syntactic processing techniques," *Artificial Intelligence Review*, vol. 56, p. 5645–5728, 2023. [Online]. Available: <https://link.springer.com/article/10.1007/s10462-022-10300-7>
- [5] J. Li *et al.*, "Recent advances in end-to-end automatic speech recognition," *APSIPA Transactions on Signal and Information Processing*, vol. 11, no. 1, 2022.
- [6] Y. Peng, S. Dalmia, I. Lane, and S. Watanabe, "Branchformer: Parallel mlp-attention architectures to capture local and global context for speech recognition and understanding," in *International Conference on Machine Learning*. PMLR, 2022, pp. 17 627–17 643.
- [7] R. Mao, X. Li, K. He, M. Ge, and E. Cambria, "Metapro online: A computational metaphor processing online system," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, 2023, pp. 127–135.
- [8] H. Bao, K. He, X. Yin, X. Li, X. Bao, H. Zhang, J. Wu, and Z. Gao, "Bert-based meta-learning approach with looking back for sentiment analysis of literary book reviews," in *Natural Language Processing and Chinese Computing: 10th CCF International Conference, NLPCC 2021, Qingdao, China, October 13–17, 2021, Proceedings, Part II 10*. Springer, 2021, pp. 235–247.
- [9] R. Mao and X. Li, "Bridging towers of multi-task learning with a gating mechanism for aspect-based sentiment analysis and sequential metaphor identification," *Proceedings of the AAI Conference on Artificial Intelligence*, vol. 35, no. 15, pp. 13 534–13 542, 2021.
- [10] K. He, Y. Huang, R. Mao, T. Gong, C. Li, and E. Cambria, "Virtual prompt pre-training for prototype-based few-shot relation extraction," *Expert Systems with Applications*, vol. 213, p. 118927, 2023.
- [11] K. He, R. Mao, T. Gong, E. Cambria, and C. Li, "Jcbie: a joint continual learning neural network for biomedical information extraction," *BMC bioinformatics*, vol. 23, no. 1, pp. 1–20, 2022.
- [12] Y. Huang, K. He, Y. Wang, X. Zhang, T. Gong, R. Mao, and C. Li, "Copner: Contrastive learning with prompt guiding for few-shot named entity recognition," in *Proceedings of the 29th International conference on computational linguistics*, 2022, pp. 2515–2527.
- [13] S. Ranathunga, E.-S. A. Lee, M. Prifti Skenduli, R. Shekhar, M. Alam, and R. Kaur, "Neural machine translation for low-resource languages: A survey," *ACM Computing Surveys*, vol. 55, no. 11, pp. 1–37, 2023.
- [14] S. Agrawal, C. Zhou, M. Lewis, L. Zettlemoyer, and M. Ghazvininejad, "In-context examples selection for machine translation," *arXiv preprint arXiv:2212.02437*, 2022.
- [15] R. Mao, G. Chen, X. Zhang, F. Guerin, and E. Cambria, "GPTEval: A survey on assessments of ChatGPT and GPT-4," *arXiv preprint arXiv:2308.12488*, 2023.
- [16] K. Singhal, T. Tu, J. Gottweis, R. Sayres, E. Wulczyn, L. Hou, K. Clark, S. Pfohl, H. Cole-Lewis, D. Neal *et al.*, "Towards expert-level medical question answering with large language models," *arXiv preprint arXiv:2305.09617*, 2023.
- [17] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [18] R. Schaeffer, B. Miranda, and S. Koyejo, "Are emergent abilities of large language models a mirage?" *arXiv preprint arXiv:2304.15004*, 2023.
- [19] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei, "Scaling laws for neural language models," *arXiv preprint arXiv:2001.08361*, 2020.
- [20] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, "Deep contextualized word representations," in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. New Orleans, Louisiana: Association for Computational Linguistics, Jun. 2018, pp. 2227–2237. [Online]. Available: <https://aclanthology.org/N18-1202>
- [21] J. D. M.-W. C. Kenton and L. K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of NAACL-HLT*, 2019, pp. 4171–4186.
- [22] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized bert pretraining approach," 2019.
- [23] K. He, N. Hong, S. Lalpalme-Remis, Y. Lan, M. Huang, C. Li, and L. Yao, "Understanding the patient perspective of epilepsy treatment through text mining of online patient support groups," *Epilepsy & Behavior*, vol. 94, pp. 65–71, 2019.
- [24] Y. Li, X. Ma, X. Zhou, P. Cheng, K. He, and C. Li, "Knowledge enhanced lstm for coreference resolution on biomedical texts," *Bioinformatics*, vol. 37, no. 17, pp. 2699–2705, 2021.
- [25] K. He, L. Yao, J. Zhang, Y. Li, and C. Li, "Construction of genealogical knowledge graphs from obituaries: Multitask neural network extraction system," *Journal of Medical Internet Research*, vol. 23, no. 8, p. e25670, 2021.
- [26] B. Mao, C. Jia, Y. Huang, K. He, J. Wu, T. Gong, and C. Li, "Uncertainty-guided mutual consistency training for semi-supervised biomedical relation extraction," in *2022 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. IEEE, 2022, pp. 2318–2325.
- [27] J. Wu, K. He, R. Mao, C. Li, and E. Cambria, "Megacare: Knowledge-guided multi-view hypergraph predictive framework for healthcare," *Information Fusion*, p. 101939, 2023.
- [28] C. Li, X. Xu, G. Zhou, K. He, T. Qi, W. Zhang, F. Tian, Q. Zheng, J. Hu *et al.*, "Implementation of national health informatization in china: survey about the status quo," *JMIR medical informatics*, vol. 7, no. 1, p. e12238, 2019.
- [29] A. J. Thirunavukarasu, D. S. J. Ting, K. Elangovan, L. Gutierrez, T. F. Tan, and D. S. W. Ting, "Large language models in medicine," *Nature medicine*, pp. 1–11, 2023.
- [30] K. S. Kalyan, A. Rajasekharan, and S. Sangeetha, "Ammu: a survey of transformer-based biomedical pretrained language models," *Journal of biomedical informatics*, vol. 126, p. 103982, 2022.
- [31] M. Moor, O. Banerjee, Z. S. H. Abad, H. M. Krumholz, J. Leskovec, E. J. Topol, and P. Rajpurkar, "Foundation models for generalist medical artificial intelligence," *Nature*, vol. 616, no. 7956, pp. 259–265, 2023.
- [32] K. He, J. Wu, X. Ma, C. Zhang, M. Huang, C. Li, and L. Yao, "Extracting kinship from obituary to enhance electronic health records for genetic research," in *Proceedings of the Fourth social media mining for health applications (# SMM4H) workshop & shared task*, 2019, pp. 1–10.
- [33] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 24 824–24 837, 2022.

- [34] Y. Hu, I. Ameer, X. Zuo, X. Peng, Y. Zhou, Z. Li, Y. Li, J. Li, X. Jiang, and H. Xu, "Zero-shot clinical entity recognition using chatgpt," *arXiv preprint arXiv:2303.16416*, 2023.
- [35] V. A. Korthikanti, J. Casper, S. Lym, L. McAfee, M. Andersch, M. Shoeybi, and B. Catanzaro, "Reducing activation recomputation in large transformer models," *Proceedings of Machine Learning and Systems*, vol. 5, 2023.
- [36] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, T. Mihaylov, M. Ott, S. Shleifer, K. Shuster, D. Simig, P. S. Koura, A. Sridhar, T. Wang, and L. Zettlemoyer, "Opt: Open pre-trained transformer language models," 2022.
- [37] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. de Las Casas, L. A. Hendricks, J. Welbl, A. Clark, T. Hennigan, E. Noland, K. Millican, G. van den Driessche, B. Damoc, A. Guy, S. Osindero, K. Simonyan, E. Elsen, J. W. Rae, O. Vinyals, and L. Sifre, "Training compute-optimal large language models," 2022.
- [38] R. Taylor, M. Kardas, G. Cucurull, T. Scialom, A. Hartshorn, E. Saravia, A. Poulton, V. Kerkez, and R. Stojnic, "Galactica: A large language model for science," 2022.
- [39] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample, "Llama: Open and efficient foundation language models," 2023.
- [40] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann *et al.*, "Palm: Scaling language modeling with pathways," *arXiv preprint arXiv:2204.02311*, 2022.
- [41] H. Yang, X.-Y. Liu, and C. D. Wang, "Fingpt: Open-source financial large language models," 2023.
- [42] S. Milano, J. A. McGrane, and S. Leonelli, "Large language models challenge the future of higher education," *Nature Machine Intelligence*, vol. 5, no. 4, pp. 333–334, 2023.
- [43] A. Arora and A. Arora, "The promise of large language models in health care," *The Lancet*, vol. 401, no. 10377, p. 641, 2023.
- [44] H. Zhang, J. Chen, F. Jiang, F. Yu, Z. Chen, J. Li, G. Chen, X. Wu, Z. Zhang, Q. Xiao *et al.*, "Huatuogpt, towards taming language model to be a doctor," *arXiv preprint arXiv:2305.15075*, 2023.
- [45] S. Chang, C. Baian, L. Fangyu, F. Zihao, S. Ehsan, and N. Collier, "Visual med-alpaca: A parameter-efficient biomedical llm with visual capabilities," 2023. [Online]. Available: <https://cambridgeidl.github.io/visual-med-alpaca/>
- [46] J. A. Omiye, H. Gui, S. J. Rezaei, J. Zou, and R. Daneshjou, "Large language models in medicine: the potentials and pitfalls," 2023.
- [47] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong *et al.*, "A survey of large language models," *arXiv preprint arXiv:2303.18223*, 2023.
- [48] B. Min, H. Ross, E. Sulem, A. P. B. Veyseh, T. H. Nguyen, O. Sainz, E. Agirre, I. Heintz, and D. Roth, "Recent advances in natural language processing via large pre-trained language models: A survey," *ACM Computing Surveys*, 2021.
- [49] P. Lavanya and E. Sasikala, "Deep learning techniques on text classification using natural language processing (nlp) in social healthcare network: A comprehensive survey," in *2021 3rd international conference on signal processing and communication (ICSPSC)*. IEEE, 2021, pp. 603–609.
- [50] N. Fei, Z. Lu, Y. Gao, G. Yang, Y. Huo, J. Wen, H. Lu, R. Song, X. Gao, T. Xiang *et al.*, "Towards artificial general intelligence via a multimodal foundation model," *Nature Communications*, vol. 13, no. 1, p. 3094, 2022.
- [51] W. Chen, Z. Li, H. Fang, Q. Yao, C. Zhong, J. Hao, Q. Zhang, X. Huang, J. Peng, and Z. Wei, "A benchmark for automatic medical consultation system: frameworks, tasks and datasets," *Bioinformatics*, vol. 39, no. 1, p. btac817, 2023.
- [52] X. Shi, Z. Liu, C. Wang, H. Leng, K. Xue, X. Zhang, and S. Zhang, "Midmed: Towards mixed-type dialogues for medical consultation," *arXiv preprint arXiv:2306.02923*, 2023.
- [53] J. H. Moon, H. Lee, W. Shin, Y.-H. Kim, and E. Choi, "Multimodal understanding and generation for medical images and text via vision-language pre-training," *IEEE Journal of Biomedical and Health Informatics*, vol. 26, no. 12, pp. 6070–6080, 2022.
- [54] G. Yang, Q. Ye, and J. Xia, "Unbox the black-box for the medical explainable ai via multi-modal and multi-centre data fusion: A mini-review, two showcases and beyond," *Information Fusion*, vol. 77, pp. 29–52, 2022.
- [55] D. Demner-Fushman, K. W. Fung, P. Do, R. D. Boyce, and T. R. Goodwin, "Overview of the tac 2018 drug-drug interaction extraction from drug labels track," in *TAC*, 2018.
- [56] Y. Deng, X. Xu, Y. Qiu, J. Xia, W. Zhang, and S. Liu, "A multimodal deep learning framework for predicting drug–drug interaction events," *Bioinformatics*, vol. 36, no. 15, pp. 4316–4322, 2020.
- [57] B.-W. Zhao, L. Hu, Z.-H. You, L. Wang, and X.-R. Su, "Hingrl: predicting drug–disease associations with graph representation learning on heterogeneous information networks," *Briefings in bioinformatics*, vol. 23, no. 1, p. bbab515, 2022.
- [58] M. Krallinger, O. Rabal, A. Lourenco, J. Oyarzabal, and A. Valencia, "Information retrieval and text mining technologies for chemistry," *Chemical reviews*, vol. 117, no. 12, pp. 7673–7761, 2017.
- [59] K. He, B. Mao, X. Zhou, Y. Li, T. Gong, C. Li, and J. Wu, "Knowledge enhanced coreference resolution via gated attention," in *2022 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. IEEE, 2022, pp. 2287–2293.
- [60] A. Nesterov and D. Umerenkov, "Distantly supervised end-to-end medical entity extraction from electronic health records with human-level quality," *arXiv preprint arXiv:2201.10463*, 2022.
- [61] R. Mythili, N. Parthiban, and V. Kavitha, "Construction of heterogeneous medical knowledge graph from electronic health records," *Journal of Discrete Mathematical Sciences and Cryptography*, vol. 25, no. 4, pp. 921–930, 2022.
- [62] L. Xiong, Y. Guo, Y. Chen, and S. Liang, "How can entities improve the quality of medical dialogue generation?" in *2023 2nd International Conference on Big Data, Information and Computer Network (BDICN)*. IEEE, 2023, pp. 225–229.
- [63] X. Lin, X. He, Q. Chen, H. Tou, Z. Wei, and T. Chen, "Enhancing dialogue symptom diagnosis with global attention and symptom graph," 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)*, 2019, pp. 5033–5042.
- [64] D. S. Wishart, C. Knox, A. C. Guo, S. Shrivastava, M. Hassanali, P. Stothard, Z. Chang, and J. Woolsey, "Drugbank: a comprehensive resource for in silico drug discovery and exploration," *Nucleic acids research*, vol. 34, no. suppl_1, pp. D668–D672, 2006.
- [65] D. A. Lindberg, B. L. Humphreys, and A. T. McCray, "The unified medical language system," *Yearbook of medical informatics*, vol. 2, no. 01, pp. 41–51, 1993.
- [66] G. Wang, A. Badal, X. Jia, J. S. Maltz, K. Mueller, K. J. Myers, C. Niu, M. Vannier, P. Yan, Z. Yu *et al.*, "Development of metaverse for intelligent healthcare," *Nature Machine Intelligence*, vol. 4, no. 11, pp. 922–929, 2022.
- [67] X. Yu, W. Hu, S. Lu, X. Sun, and Z. Yuan, "Biobert based named entity recognition in electronic medical record," in *2019 10th international conference on information technology in medicine and education (ITME)*. IEEE, 2019, pp. 49–52.
- [68] M. Chen, F. Du, G. Lan, and V. S. Lobanov, "Using pre-trained transformer deep learning models to identify named entities and syntactic relations for clinical protocol analysis," in *AAAI Spring Symposium: Combining Machine Learning with Knowledge Engineering (1)*, 2020, pp. 1–8.
- [69] Z. Miftahutdinov, I. Alimova, and E. Tutubalina, "On biomedical named entity recognition: experiments in interlingual transfer for clinical and social media texts," in *European Conference on Information Retrieval*. Springer, 2020, pp. 281–288.
- [70] Q. Wei, Z. Ji, Y. Si, J. Du, J. Wang, F. Tiryaki, S. Wu, C. Tao, K. Roberts, and H. Xu, "Relation extraction from clinical narratives using pre-trained language models," in *AMIA annual symposium proceedings*, vol. 2019. American Medical Informatics Association, 2019, p. 1236.
- [71] A. Dunn, J. Dagdelen, N. Walker, S. Lee, A. S. Rosen, G. Ceder, K. Persson, and A. Jain, "Structured information extraction from complex scientific text with fine-tuned large language models," *arXiv preprint arXiv:2212.05238*, 2022.
- [72] M. Agrawal, S. Heggelmann, H. Lang, Y. Kim, and D. Sontag, "Large language models are few-shot clinical information extractors," in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2022, pp. 1998–2022.
- [73] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, "Training language models to follow instructions with human feedback," *Advances in Neural Information Processing Systems*, vol. 35, pp. 27 730–27 744, 2022.
- [74] B. Nye, J. J. Li, R. Patel, Y. Yang, I. J. Marshall, A. Nenkova, and B. C. Wallace, "A corpus with multi-level annotations of patients, interventions and outcomes to support language processing for medical literature," in *Proceedings of the conference. Association for Compu-*

- tational Linguistics. Meeting*, vol. 2018. NIH Public Access, 2018, p. 197.
- [75] S. Moon, S. Pakhomov, N. Liu, J. O. Ryan, and G. B. Melton, "A sense inventory for clinical abbreviations and acronyms created using clinical notes and medical dictionary resources," *Journal of the American Medical Informatics Association*, vol. 21, no. 2, pp. 299–307, 2014.
- [76] S. K. Prabhakar and D.-O. Won, "Medical text classification using hybrid deep learning models with multihead attention," *Computational intelligence and neuroscience*, vol. 2021, 2021.
- [77] S. Baker, A. Korhonen, and S. Pyysalo, "Cancer hallmark text classification using convolutional neural networks," 2017.
- [78] M. A. Al-Garadi, Y.-C. Yang, H. Cai, Y. Ruan, K. O'Connor, G.-H. Graciela, J. Perrone, and A. Sarker, "Text classification models for the automatic detection of nonmedical prescription medication use from social media," *BMC medical informatics and decision making*, vol. 21, no. 1, pp. 1–13, 2021.
- [79] X. Sun, X. Li, J. Li, F. Wu, S. Guo, T. Zhang, and G. Wang, "Text classification via large language models," *arXiv preprint arXiv:2305.08377*, 2023.
- [80] H. Wang, C. Xu, and J. McAuley, "Automatic multi-label prompting: Simple and interpretable few-shot classification," in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2022, pp. 5483–5492.
- [81] T. Schick, H. Schmid, and H. Schütze, "Automatically identifying words that can serve as labels for few-shot text classification," *arXiv preprint arXiv:2010.13641*, 2020.
- [82] M. Rastegar-Mojarad, S. Liu, Y. Wang, N. Afzal, L. Wang, F. Shen, S. Fu, and H. Liu, "Biocreative/ohnlp challenge 2018," in *ACM-BCB 2018 - Proceedings of the 2018 ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*, ser. ACM-BCB 2018 - Proceedings of the 2018 ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics. Association for Computing Machinery, Inc, Aug. 2018, p. 575.
- [83] T. Botsis, G. Hartvigsen, F. Chen, and C. Weng, "Secondary use of ehr: data quality issues and informatics opportunities," *Summit on translational bioinformatics*, vol. 2010, p. 1, 2010.
- [84] D. Mahajan, A. Poddar, J. J. Liang, Y.-T. Lin, J. M. Prager, P. Suryanarayanan, P. Raghavan, C.-H. Tsou *et al.*, "Identification of semantically similar sentences in clinical notes: Iterative intermediate training using multi-task learning," *JMIR medical informatics*, vol. 8, no. 11, p. e22508, 2020.
- [85] Y. Wang, F. Liu, K. Verspoor, and T. Baldwin, "Evaluating the utility of model configurations and data augmentation on clinical semantic textual similarity," in *Proceedings of the 19th SIGBioMed workshop on biomedical language processing*, 2020, pp. 105–111.
- [86] X. Yang, X. He, H. Zhang, Y. Ma, J. Bian, Y. Wu *et al.*, "Measurement of semantic textual similarity in clinical texts: comparison of transformer-based models," *JMIR medical informatics*, vol. 8, no. 11, p. e19735, 2020.
- [87] X. Yang, A. Chen, N. PourNejatian, H. C. Shin, K. E. Smith, C. Parisien, C. Compas, C. Martin, A. B. Costa, M. G. Flores *et al.*, "A large language model for electronic health records," *NPJ Digital Medicine*, vol. 5, no. 1, p. 194, 2022.
- [88] S. Fox and M. Duggan, "Health online 2013," 2012.
- [89] S. Nerella, S. Bandyopadhyay, J. Zhang, M. Contreras, S. Siegel, A. Bumin, B. Silva, J. Sena, B. Shickel, A. Bihorac *et al.*, "Transformers in healthcare: A survey," *arXiv preprint arXiv:2307.00067*, 2023.
- [90] G. Pergola, E. Kochkina, L. Gui, M. Liakata, and Y. He, "Boosting low-resource biomedical QA via entity-aware masking strategies," in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Online: Association for Computational Linguistics, Apr. 2021, pp. 1977–1985. [Online]. Available: <https://aclanthology.org/2021.eacl-main.169>
- [91] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, "Biobert: a pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020.
- [92] Z. Chen, G. Li, and X. Wan, "Align, reason and learn: Enhancing medical vision-and-language pre-training with knowledge," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 5152–5161.
- [93] Q. Liu, R. Mao, X. Geng, and E. Cambria, "Semantic matching in machine reading comprehension: An empirical study," *Information Processing & Management*, vol. 60, no. 2, p. 103145, 2023.
- [94] K. He, R. Mao, Y. Huang, T. Gong, C. Li, and E. Cambria, "Template-free prompting for few-shot named entity recognition via semantic-enhanced contrastive learning," in *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [95] T. Gao, J. Fang, H. Liu, Z. Liu, C. Liu, P. Liu, Y. Bao, and W. Yan, "Lego-absa: A prompt-based task assemblable unified generative framework for multi-task aspect-based sentiment analysis," in *Proceedings of the 29th international conference on computational linguistics*, 2022, pp. 7002–7012.
- [96] K. He, R. Mao, T. Gong, C. Li, and E. Cambria, "Meta-based self-training and re-weighting for aspect-based sentiment analysis," *IEEE Transactions on Affective Computing*, 2022.
- [97] C. Li, F. Gao, J. Bu, L. Xu, X. Chen, Y. Gu, Z. Shao, Q. Zheng, N. Zhang, Y. Wang *et al.*, "Sentiprompt: Sentiment knowledge enhanced prompt-tuning for aspect-based sentiment analysis," *arXiv preprint arXiv:2109.08306*, 2021.
- [98] M. M. Amin, R. Mao, E. Cambria, and B. W. Schuller, "A wide evaluation of ChatGPT on affective computing tasks," *arXiv preprint arXiv:2308.13911*, 2023.
- [99] K. Singhal, S. Azizi, T. Tu, S. S. Mahdavi, J. Wei, H. W. Chung, N. Scales, A. Tanwani, H. Cole-Lewis, S. Pfohl *et al.*, "Large language models encode clinical knowledge," *arXiv preprint arXiv:2212.13138*, 2022.
- [100] A. Pal, L. K. Umapathi, and M. Sankarasubbu, "Medmqc: A large-scale multi-subject multi-choice dataset for medical domain question answering," in *Conference on Health, Inference, and Learning*. PMLR, 2022, pp. 248–260.
- [101] Q. Jin, B. Dhingra, Z. Liu, W. W. Cohen, and X. Lu, "Pubmedqa: A dataset for biomedical research question answering," *arXiv preprint arXiv:1909.06146*, 2019.
- [102] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt, "Measuring massive multitask language understanding," *arXiv preprint arXiv:2009.03300*, 2020.
- [103] A. Hamidi and K. Roberts, "Evaluation of ai chatbots for patient-specific ehr questions," *arXiv preprint arXiv:2306.02549*, 2023.
- [104] Q. Guo, S. Cao, and Z. Yi, "A medical question answering system using large language models and knowledge graphs," *International Journal of Intelligent Systems*, vol. 37, no. 11, pp. 8548–8564, 2022.
- [105] T. Kowatsch, M. Nißen, C.-H. I. Shih, D. Rügger, D. Volland, A. Filler, F. Künzler, F. Barata, S. Hung, D. Büchter *et al.*, "Text-based healthcare chatbots supporting patient and health professional teams: preliminary results of a randomized controlled trial on childhood obesity," *Persuasive Embodied Agents for Behavior Change (PEACH2017)*, 2017.
- [106] B. Chaix, J.-E. Bibault, A. Pienkowski, G. Delamon, A. Guillemassé, P. Nectoux, B. Brouard *et al.*, "When chatbots meet patients: one-year prospective study of conversations with patients with breast cancer and a chatbot," *JMIR cancer*, vol. 5, no. 1, p. e12856, 2019.
- [107] S. Ji, T. Zhang, K. Yang, S. Ananiadou, E. Cambria, and J. Tiedemann, "Domain-specific continued pretraining of language models for capturing long context in mental health," *arXiv preprint arXiv:2304.10447*, 2023.
- [108] L. Reis, C. Maier, J. Mattke, and T. Weitzel, "Chatbots in healthcare: Status quo, application scenarios for physicians and patients and future directions," 2020.
- [109] J. Ni, T. Young, V. Pandelea, F. Xue, and E. Cambria, "Recent advances in deep learning based dialogue systems: A systematic survey," *Artificial intelligence review*, vol. 56, no. 4, pp. 3055–3155, 2023.
- [110] A. A. Abd-Alrazaq, M. Alajlani, N. Ali, K. Denecke, B. M. Bewick, and M. Househ, "Perceptions and opinions of patients about mental health chatbots: scoping review," *Journal of medical Internet research*, vol. 23, no. 1, p. e17828, 2021.
- [111] S. Ji, T. Zhang, L. Ansari, J. Fu, P. Tiwari, and E. Cambria, "Mentalbert: Publicly available pretrained language models for mental healthcare," 2021.
- [112] R. L. Pande, M. Morris, A. Peters, C. M. Spettell, R. Feifer, and W. Gillis, "Leveraging remote behavioral health interventions to improve medical outcomes and reduce costs," *Am J Manag Care*, vol. 21, no. 2, pp. e141–e151, 2015.
- [113] D. Milward and M. Beveridge, "Ontology-based dialogue systems," in *Proc. 3rd Workshop on Knowledge and Reasoning in practical dialogue systems (IJCAI03)*, 2003, pp. 9–18.
- [114] L. Xu, Q. Zhou, K. Gong, X. Liang, J. Tang, and L. Lin, "End-to-end knowledge-routed relational dialogue system for automatic diagnosis," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 7346–7353.

- [115] W. Qin, Z. Chen, L. Wang, Y. Lan, W. Ren, and R. Hong, "Read, diagnose and chat: Towards explainable and interactive llms-augmented depression detection in social media," *arXiv preprint arXiv:2305.05138*, 2023.
- [116] L. Yunxiang, L. Zihan, Z. Kai, D. Ruilong, and Z. You, "Chatdoctor: A medical chat model fine-tuned on llama model using medical domain knowledge," *arXiv preprint arXiv:2303.14070*, 2023.
- [117] B. Jing, P. Xie, and E. Xing, "On the automatic generation of medical imaging reports," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2018, pp. 2577–2586.
- [118] Y. Xue, T. Xu, L. Rodney Long, Z. Xue, S. Antani, G. R. Thoma, and X. Huang, "Multimodal recurrent model with attention for automated radiology report generation," in *Medical Image Computing and Computer Assisted Intervention—MICCAI 2018: 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part I*. Springer, 2018, pp. 457–466.
- [119] J. Chen, H. Guo, K. Yi, B. Li, and M. Elhoseiny, "Visualgpt: Data-efficient adaptation of pretrained language models for image captioning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 18 030–18 040.
- [120] S. Wang, Z. Zhao, X. Ouyang, Q. Wang, and D. Shen, "Chatcad: Interactive computer-aided diagnosis on medical image using large language models," *arXiv preprint arXiv:2302.07257*, 2023.
- [121] Z. Chen, Y. Song, T.-H. Chang, and X. Wan, "Generating radiology reports via memory-driven transformer," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 1439–1449.
- [122] A. Nicolson, J. Dowling, and B. Koopman, "Improving chest x-ray report generation by leveraging warm-starting," *arXiv preprint arXiv:2201.09405*, 2022.
- [123] Z. Zhao, S. Wang, J. Gu, Y. Zhu, L. Mei, Z. Zhuang, Z. Cui, Q. Wang, and D. Shen, "Chatcad+: Towards a universal and reliable interactive cad using llms," *arXiv preprint arXiv:2305.15964*, 2023.
- [124] Z. Gao, B. Hong, X. Zhang, Y. Li, C. Jia, J. Wu, C. Wang, D. Meng, and C. Li, "Instance-based vision transformer for subtyping of papillary renal cell carcinoma in histopathological image," in *Medical Image Computing and Computer Assisted Intervention—MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part VIII 24*. Springer, 2021, pp. 299–308.
- [125] Z. Gao, A. Mao, K. Wu, Y. Li, L. Zhao, X. Zhang, J. Wu, L. Yu, C. Xing, T. Gong *et al.*, "Childhood leukemia classification via information bottleneck enhanced hierarchical multi-instance learning," *IEEE Transactions on Medical Imaging*, 2023.
- [126] J. Shi, L. Tang, Y. Li, X. Zhang, Z. Gao, Y. Zheng, C. Wang, T. Gong, and C. Li, "A structure-aware hierarchical graph-based multiple instance learning framework for pt staging in histopathological image," *IEEE Transactions on Medical Imaging*, 2023.
- [127] S. Liu, Z. Zhu, Q. Qu, and C. You, "Robust training under label noise by over-parameterization," in *International Conference on Machine Learning*. PMLR, 2022, pp. 14 153–14 172.
- [128] C. Zhou, P. Liu, P. Xu, S. Iyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu *et al.*, "Lima: Less is more for alignment," *arXiv preprint arXiv:2305.11206*, 2023.
- [129] X. Li, P. Yu, C. Zhou, T. Schick, L. Zettlemoyer, O. Levy, J. Weston, and M. Lewis, "Self-alignment with instruction backtranslation," 2023.
- [130] Y. Peng, S. Yan, and Z. Lu, "Transfer learning in biomedical natural language processing: an evaluation of bert and elmo on ten benchmarking datasets," *arXiv preprint arXiv:1906.05474*, 2019.
- [131] Z. Kraljevic, A. Shek, D. Bean, R. Bendayan, J. Teo, and R. Dobson, "Medgpt: Medical concept prediction from clinical narratives," *arXiv preprint arXiv:2107.03134*, 2021.
- [132] S. Sharma and R. Daniel Jr, "Bioflair: Pretrained pooled contextualized embeddings for biomedical sequence labeling tasks," *arXiv preprint arXiv:1908.05760*, 2019.
- [133] I. B. Ozyurt, "On the effectiveness of small, discriminatively pretrained language representation models for biomedical text mining," in *Proceedings of the First Workshop on Scholarly Document Processing*, 2020, pp. 104–112.
- [134] Y.-P. Chen, Y.-Y. Chen, J.-J. Lin, C.-H. Huang, F. Lai *et al.*, "Modified bidirectional encoder representations from transformers extractive summarization model for hospital information systems based on character-level tokens (alphabert): development and performance evaluation," *JMIR medical informatics*, vol. 8, no. 4, p. e17787, 2020.
- [135] L. Akhtyamova, "Named entity recognition in spanish biomedical literature: Short review and bert model," in *2020 26th Conference of Open Innovations Association (FRUCT)*. IEEE, 2020, pp. 1–7.
- [136] N. Poerner, U. Waltinger, and H. Schütze, "Inexpensive domain adaptation of pretrained language models: Case studies on biomedical ner and covid-19 qa," *arXiv preprint arXiv:2004.03354*, 2020.
- [137] Y. Li, S. Rao, J. R. A. Solares, A. Hassaine, R. Ramakrishnan, D. Canoy, Y. Zhu, K. Rahimi, and G. Salimi-Khorshidi, "Behrt: transformer for electronic health records," *Scientific reports*, vol. 10, no. 1, p. 7155, 2020.
- [138] S. Gururangan, A. Marasović, S. Swayamdipta, K. Lo, I. Beltagy, D. Downey, and N. A. Smith, "Don't stop pretraining: Adapt language models to domains and tasks," *arXiv preprint arXiv:2004.10964*, 2020.
- [139] X. Meng, C. H. Ganoe, R. T. Sieberg, Y. Y. Cheung, and S. Hassanpour, "Self-supervised contextual language representation of radiology reports to improve the identification of communication urgency," *AMIA Summits on Translational Science Proceedings*, vol. 2020, p. 413, 2020.
- [140] M. Müller, M. Salathé, and P. E. Kummervold, "Covid-twitter-bert: A natural language processing model to analyse covid-19 content on twitter," *Frontiers in Artificial Intelligence*, vol. 6, p. 1023281, 2023.
- [141] J. Copara, J. Knafou, N. Naderi, C. Moro, P. Ruch, and D. Teodoro, "Contextualized French language models for biomedical named entity recognition," in *Actes de la 6e conférence conjointe Journées d'Études sur la Parole (JEP, 33e édition), Traitement Automatique des Langues Naturelles (TALN, 27e édition), Rencontre des Étudiants Chercheurs en Informatique pour le Traitement Automatique des Langues (RÉCITAL, 22e édition), Atelier DÉfi Fouille de Textes*. Nancy, France: ATALA et AFCEP, 6 2020, pp. 36–48. [Online]. Available: <https://aclanthology.org/2020.jepalnrecital-deft.4>
- [142] K. K. Bressemer, L. C. Adams, R. A. Gaudin, D. Tröltzsch, B. Hamm, M. R. Makowski, C.-Y. Schüle, J. L. Vahldiek, and S. M. Niehues, "Highly accurate classification of chest radiographic reports using a deep learning natural language model pre-trained on 3.8 million text reports," *Bioinformatics*, vol. 36, no. 21, pp. 5255–5261, 2020.
- [143] Y. Kawazoe, D. Shibata, E. Shinohara, E. Aramaki, and K. Ohe, "A clinical specific bert developed with huge size of japanese clinical narrative," *medRxiv*, pp. 2020–07, 2020.
- [144] N. Zhang, Q. Jia, K. Yin, L. Dong, F. Gao, and N. Hua, "Conceptualized representation learning for chinese biomedical text mining," *arXiv preprint arXiv:2008.10813*, 2020.
- [145] U. Naseem, M. Khushi, V. Reddy, S. Rajendran, I. Razzak, and J. Kim, "Bioalbert: A simple and effective pre-trained language model for biomedical named entity recognition," in *2021 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2021, pp. 1–7.
- [146] H.-C. Shin, Y. Zhang, E. Bakhturina, R. Puri, M. Patwary, M. Shoeybi, and R. Mani, "Biomegatron: Larger biomedical domain language model," 2020.
- [147] K. Huang, J. Altosaar, and R. Ranganath, "Clinicalbert: Modeling clinical notes and predicting hospital readmission," *arXiv preprint arXiv:1904.05342*, 2019.
- [148] K. Huang, A. Singh, S. Chen, E. Moseley, C.-Y. Deng, N. George, and C. Lindvall, "Clinical xlnet: Modeling sequential clinical notes and predicting prolonged mechanical ventilation," in *Proceedings of the 3rd Clinical Natural Language Processing Workshop*, 2020, pp. 94–100.
- [149] P. Lewis, M. Ott, J. Du, and V. Stoyanov, "Pretrained language models for biomedical and clinical tasks: understanding and extending the state-of-the-art," in *Proceedings of the 3rd Clinical Natural Language Processing Workshop*, 2020, pp. 146–157.
- [150] E. T. R. Schneider, J. V. A. de Souza, J. Knafou, L. E. S. e. Oliveira, J. Copara, Y. B. Gumiel, L. F. A. d. Oliveira, E. C. Paraiso, D. Teodoro, and C. M. C. M. Barra, "BioBERTpt - a Portuguese neural language model for clinical named entity recognition," in *Proceedings of the 3rd Clinical Natural Language Processing Workshop*. Online: Association for Computational Linguistics, Nov. 2020, pp. 65–72. [Online]. Available: <https://aclanthology.org/2020.clinicalnlp-1.7>
- [151] X. Yang, J. Bian, W. R. Hogan, and Y. Wu, "Clinical concept extraction using transformers," *Journal of the American Medical Informatics Association*, vol. 27, no. 12, pp. 1935–1942, 2020.
- [152] B. Hao, H. Zhu, and I. C. Paschalidis, "Enhancing clinical bert embedding using a biomedical knowledge base," in *28th International Conference on Computational Linguistics (COLING 2020)*, 2020.
- [153] J. Wang, G. Zhang, W. Wang, K. Zhang, and Y. Sheng, "Cloud-based intelligent self-diagnosis and department recommendation service using chinese medical bert," *Journal of Cloud Computing*, vol. 10, pp. 1–12, 2021.
- [154] Y. Gu, R. Tinn, H. Cheng, M. Lucas, N. Usuyama, X. Liu, T. Naumann, J. Gao, and H. Poon, "Domain-specific language model pretraining for biomedical natural language processing," *ACM Transactions on Computing for Healthcare (HEALTH)*, vol. 3, no. 1, pp. 1–23, 2021.

- [155] S. Wada, T. Takeda, S. Manabe, S. Konishi, J. Kamohara, and Y. Matsumura, "Pre-training technique to localize medical bert and enhance biomedical bert," *arXiv preprint arXiv:2005.07202*, 2020.
- [156] Y. Meng, W. Speier, M. K. Ong, and C. W. Arnold, "Bidirectional representation learning from transformers using multimodal electronic health record data to predict depression," *IEEE Journal of Biomedical and Health Informatics*, vol. 25, no. 8, pp. 3121–3129, 2021.
- [157] W. Antoun, F. Baly, and H. Hajji, "Arabert: Transformer-based model for arabic language understanding," in *Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection*, 2020, pp. 9–15.
- [158] N. Boudjellal, H. Zhang, A. Khan, A. Ahmad, R. Naseem, J. Shang, and L. Dai, "Abioner: a bert-based model for arabic biomedical named-entity recognition," *Complexity*, vol. 2021, pp. 1–6, 2021.
- [159] G. Miolo, G. Mantoan, and C. Orsenigo, "Electrmed: a new pre-trained language representation model for biomedical nlp," *arXiv preprint arXiv:2104.09585*, 2021.
- [160] Z. Yuan, Y. Liu, C. Tan, S. Huang, and F. Huang, "Improving biomedical pretrained language models with knowledge," *arXiv preprint arXiv:2104.10344*, 2021.
- [161] N. Taghizadeh, E. Doostmohammadi, E. Seifossadat, H. R. Rabiee, and M. S. Tahaei, "Sina-bert: a pre-trained language model for analysis of medical texts in persian," *arXiv preprint arXiv:2104.07613*, 2021.
- [162] L. Rasmy, Y. Xiang, Z. Xie, C. Tao, and D. Zhi, "Med-bert: pretrained contextualized embeddings on large-scale structured electronic health records for disease prediction," *NPJ digital medicine*, vol. 4, no. 1, p. 86, 2021.
- [163] G. López-García, J. M. Jerez, N. Ribelles, E. Alba, and F. J. Veredas, "Transformers for clinical coding in spanish," *IEEE Access*, vol. 9, pp. 72 387–72 397, 2021.
- [164] L. N. Phan, J. T. Anibal, H. Tran, S. Chanana, E. Bahadroglu, A. Peltekian, and G. Altan-Bonnet, "Scifive: a text-to-text transformer model for biomedical literature," *arXiv preprint arXiv:2106.03598*, 2021.
- [165] K. r. Kanakarajan, B. Kundumani, and M. Sankarasubbu, "BioELECTRA:pretrained biomedical text encoder using discriminators," in *Proceedings of the 20th Workshop on Biomedical Language Processing*. Online: Association for Computational Linguistics, Jun. 2021, pp. 143–154. [Online]. Available: <https://aclanthology.org/2021.bionlp-1.16>
- [166] Z. Yuan, Z. Zhao, H. Sun, J. Li, F. Wang, and S. Yu, "Coder: Knowledge-infused cross-lingual medical term embedding for term normalization," *Journal of biomedical informatics*, vol. 126, p. 103983, 2022.
- [167] M. Yasunaga, J. Leskovec, and P. Liang, "Linkbert: Pretraining language models with document links," *arXiv preprint arXiv:2203.15827*, 2022.
- [168] U. Naseem, A. G. Dunn, M. Khushi, and J. Kim, "Benchmarking for biomedical natural language processing tasks with a domain specific albert," *BMC bioinformatics*, vol. 23, no. 1, pp. 1–15, 2022.
- [169] H. Yuan, Z. Yuan, R. Gan, J. Zhang, Y. Xie, and S. Yu, "Biobart: Pretraining and evaluation of a biomedical generative language model," *BioNLP 2022 @ ACL 2022*, p. 97, 2022.
- [170] F. Liu, E. Shareghi, Z. Meng, M. Basaldella, and N. Collier, "Self-alignment pretraining for biomedical entity representations," *arXiv preprint arXiv:2010.11784*, 2020.
- [171] X. Zhang, C. Wu, Y. Zhang, Y. Wang, and W. Xie, "Knowledge-enhanced visual-language pre-training on chest radiology images," 2023.
- [172] Y. Choi, C. Y.-I. Chiu, and D. Sontag, "Learning low-dimensional representations of medical concepts," *AMIA Summits on Translational Science Proceedings*, vol. 2016, p. 41, 2016.
- [173] Z. Zeng, C. Xiao, Y. Yao, R. Xie, Z. Liu, F. Lin, L. Lin, and M. Sun, "Knowledge transfer via pre-training for recommendation: A review and prospect," *Frontiers in big Data*, vol. 4, p. 602071, 2021.
- [174] X. Qiu, T. Sun, Y. Xu, Y. Shao, N. Dai, and X. Huang, "Pre-trained models for natural language processing: A survey," *Science China Technological Sciences*, vol. 63, no. 10, pp. 1872–1897, 2020.
- [175] F. Shamshad, S. Khan, S. W. Zamir, M. H. Khan, M. Hayat, F. S. Khan, and H. Fu, "Transformers in medical imaging: A survey," *Medical Image Analysis*, p. 102802, 2023.
- [176] H. Bao, L. Dong, F. Wei, W. Wang, N. Yang, X. Liu, Y. Wang, J. Gao, S. Piao, M. Zhou *et al.*, "Unilmv2: Pseudo-masked language models for unified language model pre-training," in *International conference on machine learning*. PMLR, 2020, pp. 642–652.
- [177] L. Dong, N. Yang, W. Wang, F. Wei, X. Liu, Y. Wang, J. Gao, M. Zhou, and H.-W. Hon, "Unified language model pre-training for natural language understanding and generation," *Advances in neural information processing systems*, vol. 32, 2019.
- [178] Y. Tay, M. Dehghani, V. Q. Tran, X. Garcia, D. Bahri, T. Schuster, H. S. Zheng, N. Houshy, and D. Metzler, "Unifying language learning paradigms," *arXiv preprint arXiv:2205.05131*, 2022.
- [179] Y. Sun, S. Wang, Y. Li, S. Feng, H. Tian, H. Wu, and H. Wang, "Ernie 2.0: A continual pre-training framework for language understanding," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 05, 2020, pp. 8968–8975.
- [180] Y. Sun, S. Wang, S. Feng, S. Ding, C. Pang, J. Shang, J. Liu, X. Chen, Y. Zhao, Y. Lu *et al.*, "Ernie 3.0: Large-scale knowledge enhanced pre-training for language understanding and generation," *arXiv preprint arXiv:2107.02137*, 2021.
- [181] X. Yang, N. PourNejatian, H. C. Shin, K. E. Smith, C. Parisien, C. Compas, C. Martin, M. G. Flores, Y. Zhang, T. Magoc *et al.*, "Gatortron: A large language model for clinical natural language processing," *medRxiv*, pp. 2022–02, 2022.
- [182] V. Liévin, C. E. Hother, and O. Winther, "Can large language models reason about medical questions?" *arXiv preprint arXiv:2207.08143*, 2022.
- [183] H. Nori, N. King, S. M. McKinney, D. Carignan, and E. Horvitz, "Capabilities of gpt-4 on medical challenge problems," *arXiv preprint arXiv:2303.13375*, 2023.
- [184] Z. Liu, X. Yu, L. Zhang, Z. Wu, C. Cao, H. Dai, L. Zhao, W. Liu, D. Shen, Q. Li *et al.*, "Deid-gpt: Zero-shot medical text de-identification by gpt-4," *arXiv preprint arXiv:2303.11032*, 2023.
- [185] H. Xiong, S. Wang, Y. Zhu, Z. Zhao, Y. Liu, Q. Wang, and D. Shen, "Doctorglm: Fine-tuning your chinese doctor is not a herculean task," *arXiv preprint arXiv:2304.01097*, 2023.
- [186] T. Han, L. C. Adams, J.-M. Papaioannou, P. Grundmann, T. Oberhauser, A. Löser, D. Truhn, and K. K. Bressen, "Medalpaca—an open-source collection of medical conversational ai models and training data," *arXiv preprint arXiv:2304.08247*, 2023.
- [187] H. Wang, C. Liu, N. Xi, Z. Qiang, S. Zhao, B. Qin, and T. Liu, "Huatuo: Tuning llama model with chinese medical knowledge," 2023.
- [188] C. Wu, X. Zhang, Y. Zhang, Y. Wang, and W. Xie, "Pmc-llama: Further finetuning llama on medical papers," *arXiv preprint arXiv:2304.14454*, 2023.
- [189] C. Yirong, W. Zhenyu, X. Xiaofen, X. Zhipei, F. Kai, L. Sihang, W. Junhong, and X. Xiangmin, "Bianque-1.0: Improving the "question" ability of medical chat model through finetuning with hybrid instructions and multi-turn doctor qa datasets," 2023. [Online]. Available: <https://github.com/scutcyr/BianQue>
- [190] C. Peng, X. Yang, A. Chen, K. E. Smith, N. PourNejatian, A. B. Costa, C. Martin, M. G. Flores, Y. Zhang, T. Magoc *et al.*, "A study of generative large language model for medical research and healthcare," *arXiv preprint arXiv:2305.13523*, 2023.
- [191] G. Wang, G. Yang, Z. Du, L. Fan, and X. Li, "Clinicalgpt: Large language models finetuned with diverse medical data and comprehensive evaluation," *arXiv preprint arXiv:2306.09968*, 2023.
- [192] J. Zhou, X. Chen, and X. Gao, "Path to medical agi: Unify domain-specific medical llms with the lowest cost," *arXiv preprint arXiv:2306.10765*, 2023.
- [193] C. Li, C. Wong, S. Zhang, N. Usuyama, H. Liu, J. Yang, T. Naumann, H. Poon, and J. Gao, "Llava-med: Training a large language-and-vision assistant for biomedicine in one day," *arXiv preprint arXiv:2306.00890*, 2023.
- [194] W. Gao, Z. Deng, Z. Niu, F. Rong, C. Chen, Z. Gong, W. Zhang, D. Xiao, F. Li, Z. Cao *et al.*, "Ophglm: Training an ophthalmology large language-and-vision assistant based on instructions and dialogue," *arXiv preprint arXiv:2306.12174*, 2023.
- [195] C. Yirong, X. Xiaofen, W. Zhenyu, and X. Xiangmin, "Soulchat: The "empathy" ability of the large model is improved by mixing and fine-tuning the data set of long text consultation instructions and multiple rounds of empathy dialogue," 6 2023. [Online]. Available: <https://github.com/scutcyr/SoulChat>
- [196] M. Moor, Q. Huang, S. Wu, M. Yasunaga, C. Zakka, Y. Dalmia, E. P. Reis, P. Rajpurkar, and J. Leskovec, "Med-flamingo: a multimodal medical few-shot learner," *arXiv preprint arXiv:2307.15189*, 2023.
- [197] Ö. Uzuner, B. R. South, S. Shen, and S. L. DuVall, "2010 i2b2/va challenge on concepts, assertions, and relations in clinical text," *Journal of the American Medical Informatics Association*, vol. 18, no. 5, pp. 552–556, 2011.
- [198] W. Sun, A. Rumshisky, and O. Uzuner, "Evaluating temporal relations in clinical text: 2012 i2b2 challenge," *Journal of the American Medical Informatics Association*, vol. 20, no. 5, pp. 806–813, 2013.

- [199] X. Yang, J. Bian, R. Fang, R. I. Bjarnadottir, W. R. Hogan, and Y. Wu, "Identifying relations of medications with adverse drug events using recurrent convolutional neural networks and gradient boosting," *Journal of the American Medical Informatics Association*, vol. 27, no. 1, pp. 65–72, 2020.
- [200] Y. Wang, S. Fu, F. Shen, S. Henry, O. Uzuner, H. Liu *et al.*, "The 2019 n2c2/ohnlp track on clinical semantic textual similarity: overview," *JMIR medical informatics*, vol. 8, no. 11, p. e23375, 2020.
- [201] C. Shivade, "Mednli-a natural language inference dataset for the clinical domain (version 1.0.0). physionet," 2019.
- [202] A. Pampari, P. Raghavan, J. Liang, and J. Peng, "emrqa: A large corpus for question answering on electronic medical records," *arXiv preprint arXiv:1809.00732*, 2018.
- [203] S. Smith, M. Patwary, B. Norick, P. LeGresley, S. Rajbhandari, J. Casper, Z. Liu, S. Prabhunoye, G. Zerveas, V. Korthikanti *et al.*, "Using deepspeed and megatron to train megatron-turing nlq 530b, a large-scale generative language model," *arXiv preprint arXiv:2201.11990*, 2022.
- [204] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. d. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman *et al.*, "Evaluating large language models trained on code," *arXiv preprint arXiv:2107.03374*, 2021.
- [205] D. Jin, E. Pan, N. Oufattole, W.-H. Weng, H. Fang, and P. Szolovits, "What disease does this patient have? a large-scale open domain question answering dataset from medical exams," *Applied Sciences*, vol. 11, no. 14, p. 6421, 2021.
- [206] D. Hendrycks and K. Gimpel, "Gaussian error linear units (gelu)," *arXiv preprint arXiv:1606.08415*, 2016.
- [207] O. Press, N. A. Smith, and M. Lewis, "Train short, test long: Attention with linear biases enables input length extrapolation," *arXiv preprint arXiv:2108.12409*, 2021.
- [208] A. Nentidis, G. Katsimpras, E. Vandrova, A. Krithara, L. Gasco, M. Krallinger, and G. Paliouras, "Overview of bioasq 2021: The ninth bioasq challenge on large-scale biomedical semantic indexing and question answering," in *Experimental IR Meets Multilinguality, Multimodality, and Interaction: 12th International Conference of the CLEF Association, CLEF 2021, Virtual Event, September 21–24, 2021, Proceedings 12*. Springer, 2021, pp. 239–263.
- [209] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 2002, pp. 311–318.
- [210] A. B. Abacha, E. Agichtein, Y. Pinter, and D. Demner-Fushman, "Overview of the medical question answering task at trec 2017 liveqa," in *TREC*, 2017, pp. 1–12.
- [211] A. B. Abacha, Y. Mrabet, M. Sharp, T. R. Goodwin, S. E. Shooshan, and D. Demner-Fushman, "Bridging the gap between consumers' medication questions and trusted answers," in *MedInfo*, 2019, pp. 25–29.
- [212] P. Barham, A. Chowdhery, J. Dean, S. Ghemawat, S. Hand, D. Hurt, M. Isard, H. Lim, R. Pang, S. Roy *et al.*, "Pathways: Asynchronous distributed dataflow for ml," *Proceedings of Machine Learning and Systems*, vol. 4, pp. 430–449, 2022.
- [213] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, E. Li, X. Wang, M. Dehghani, S. Brahma *et al.*, "Scaling instruction-finetuned language models," *arXiv preprint arXiv:2210.11416*, 2022.
- [214] R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, and T. B. Hashimoto, "Stanford alpaca: An instruction-following llama model," https://github.com/tatsu-lab/stanford_alpaca, 2023.
- [215] Z. Du, Y. Qian, X. Liu, M. Ding, J. Qiu, Z. Yang, and J. Tang, "Glm: General language model pretraining with autoregressive blank infilling," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022, pp. 320–335.
- [216] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," *arXiv preprint arXiv:2106.09685*, 2021.
- [217] T. Dettmers, M. Lewis, Y. Belkada, and L. Zettlemoyer, "Llm.int8(): 8-bit matrix multiplication for transformers at scale," *arXiv preprint arXiv:2208.07339*, 2022.
- [218] T. Dettmers, M. Lewis, S. Shleifer, and L. Zettlemoyer, "8-bit optimizers via block-wise quantization," *arXiv preprint arXiv:2110.02861*, 2021.
- [219] K. Lo, L. L. Wang, M. Neumann, R. Kinney, and D. S. Weld, "S2orc: The semantic scholar open research corpus," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 4969–4983.
- [220] J. Fries, L. Weber, N. Seelam, G. Altay, D. Datta, S. Garda, S. Kang, R. Su, W. Kusa, S. Cahyawijaya *et al.*, "Bigbio: a framework for data-centric biomedical natural language processing," *Advances in Neural Information Processing Systems*, vol. 35, pp. 25 792–25 806, 2022.
- [221] L. X. Xuanwei Zhang and K. Zhao, "Chatyuan: A large language model for dialogue in chinese and english," Dec. 2022. [Online]. Available: <https://github.com/clue-ai/ChatYuan>
- [222] S. Chen, Z. Ju, X. Dong, H. Fang, S. Wang, Y. Yang, J. Zeng, R. Zhang, R. Zhang, M. Zhou, P. Zhu, and P. Xie, "Meddialog: a large-scale medical dialogue dataset," *arXiv preprint arXiv:2004.03329*, 2020.
- [223] W. Chen, Z. Li, H. Fang, Q. Yao, C. Zhong, J. Hao, Q. Zhang, X. Huang, J. Peng, and Z. Wei, "A Benchmark for Automatic Medical Consultation System: Frameworks, Tasks and Datasets," *Bioinformatics*, 12 2022, btac817. [Online]. Available: <https://doi.org/10.1093/bioinformatics/btac817>
- [224] N. Zhang, M. Chen, Z. Bi, X. Liang, L. Li, X. Shang, K. Yin, C. Tan, J. Xu, F. Huang, L. Si, Y. Ni, G. Xie, Z. Sui, B. Chang, H. Zong, Z. Yuan, L. Li, J. Yan, H. Zan, K. Zhang, B. Tang, and Q. Chen, "CBLUE: A Chinese biomedical language understanding evaluation benchmark," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 7888–7915. [Online]. Available: <https://aclanthology.org/2022.acl-long.544>
- [225] W. Chen, Z. Li, H. Fang, Q. Yao, C. Zhong, J. Hao, Q. Zhang, X. Huang, J. Peng, and Z. Wei, "A Benchmark for Automatic Medical Consultation System: Frameworks, Tasks and Datasets," *Bioinformatics*, 12 2022, btac817.
- [226] S. Zhang, X. Zhang, H. Wang, L. Guo, and S. Liu, "Multi-scale attentive interaction networks for chinese medical question answer selection," *IEEE Access*, vol. 6, pp. 74 061–74 071, 2018.
- [227] R. Anil, A. M. Dai, O. Firat, M. Johnson, D. Lepikhin, A. Passos, S. Shakeri, E. Taropa, P. Bailey, Z. Chen *et al.*, "Palm 2 technical report," *arXiv preprint arXiv:2305.10403*, 2023.
- [228] L. Gao, S. Biderman, S. Black, L. Golding, T. Hoppe, C. Foster, J. Phang, H. He, A. Thite, N. Nabeshima *et al.*, "The pile: An 800gb dataset of diverse text for language modeling," *arXiv preprint arXiv:2101.00027*, 2020.
- [229] M. Herrero-Zazo, I. Segura-Bedmar, P. Martínez, and T. Declerck, "The ddi corpus: An annotated corpus with pharmacological substances and drug–drug interactions," *Journal of biomedical informatics*, vol. 46, no. 5, pp. 914–920, 2013.
- [230] J. Li, Y. Sun, R. J. Johnson, D. Sciaky, C.-H. Wei, R. Leaman, A. P. Davis, C. J. Mattingly, T. C. Wieggers, and Z. Lu, "Biocreative v cdr task corpus: a resource for chemical disease relation extraction," *Database*, vol. 2016, 2016.
- [231] Y. Hou, Y. Xia, L. Wu, S. Xie, Y. Fan, J. Zhu, T. Qin, and T.-Y. Liu, "Discovering drug–target interaction knowledge from biomedical literature," *Bioinformatics*, vol. 38, no. 22, pp. 5100–5107, 2022.
- [232] S. Zhang, X. Zhang, H. Wang, L. Guo, and S. Liu, "Multi-scale attentive interaction networks for chinese medical question answer selection," *IEEE Access*, vol. 6, pp. 74 061–74 071, 2018.
- [233] J. He, M. Fu, and M. Tu, "Applying deep matching networks to chinese medical question answering: a study and a dataset," *BMC medical informatics and decision making*, vol. 19, no. 2, pp. 91–100, 2019.
- [234] J. Li, X. Wang, X. Wu, Z. Zhang, X. Xu, J. Fu, P. Tiwari, X. Wan, and B. Wang, "Huatuo-26m, a large-scale chinese medical qa dataset," *arXiv preprint arXiv:2305.01526*, 2023.
- [235] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *The Journal of Machine Learning Research*, vol. 21, no. 1, pp. 5485–5551, 2020.
- [236] X. He, S. Chen, Z. Ju, X. Dong, H. Fang, S. Wang, Y. Yang, J. Zeng, R. Zhang, R. Zhang *et al.*, "Meddialog: Two large-scale medical dialogue datasets," *arXiv preprint arXiv:2004.03329*, 2020.
- [237] H. Liu, C. Li, Q. Wu, and Y. J. Lee, "Visual instruction tuning," *arXiv preprint arXiv:2304.08485*, 2023.
- [238] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [239] S. Zhang, Y. Xu, N. Usuyama, J. Bagga, R. Tinn, S. Preston, R. Rao, M. Wei, N. Valluri, C. Wong *et al.*, "Large-scale domain-specific pretraining for biomedical vision-language processing," *arXiv preprint arXiv:2303.00915*, 2023.

- [240] J. J. Lau, S. Gayen, A. Ben Abacha, and D. Demner-Fushman, "A dataset of clinically generated visual questions and answers about radiology images," *Scientific data*, vol. 5, no. 1, pp. 1–10, 2018.
- [241] B. Liu, L.-M. Zhan, L. Xu, L. Ma, Y. Yang, and X.-M. Wu, "Slake: A semantically-labeled knowledge-enhanced dataset for medical visual question answering," in *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2021, pp. 1650–1654.
- [242] X. He, Z. Cai, W. Wei, Y. Zhang, L. Mou, E. Xing, and P. Xie, "Pathological visual question answering," *arXiv preprint arXiv:2010.12435*, 2020.
- [243] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds *et al.*, "Flamingo: a visual language model for few-shot learning," *Advances in Neural Information Processing Systems*, vol. 35, pp. 23 716–23 736, 2022.
- [244] W. Lin, Z. Zhao, X. Zhang, C. Wu, Y. Zhang, Y. Wang, and W. Xie, "Pmc-clip: Contrastive language-image pre-training using biomedical documents," *arXiv preprint arXiv:2303.07240*, 2023.
- [245] X. He, Y. Zhang, L. Mou, E. Xing, and P. Xie, "Pathvqa: 30000+ questions for medical visual question answering," 2020.
- [246] R. Rafailov, A. Sharma, E. Mitchell, S. Ermon, C. D. Manning, and C. Finn, "Direct preference optimization: Your language model is secretly a reward model," *arXiv preprint arXiv:2305.18290*, 2023.
- [247] S. Rajbhandari, J. Rasley, O. Ruwase, and Y. He, "Zero: Memory optimizations toward training trillion parameter models," in *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 2020, pp. 1–16.
- [248] A. Webson and E. Pavlick, "Do prompt-based models really understand the meaning of their prompts?" in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2022, pp. 2300–2344.
- [249] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [250] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, "Albert: A lite bert for self-supervised learning of language representations," *arXiv preprint arXiv:1909.11942*, 2019.
- [251] Z. Zhang, X. Han, Z. Liu, X. Jiang, M. Sun, and Q. Liu, "Ernie: Enhanced language representation with informative entities," *arXiv preprint arXiv:1905.07129*, 2019.
- [252] P. He, X. Liu, J. Gao, and W. Chen, "Deberta: Decoding-enhanced bert with disentangled attention," *arXiv preprint arXiv:2006.03654*, 2020.
- [253] K. Clark, M.-T. Luong, Q. V. Le, and C. D. Manning, "Electra: Pre-training text encoders as discriminators rather than generators," *arXiv preprint arXiv:2003.10555*, 2020.
- [254] M. Joshi, D. Chen, Y. Liu, D. S. Weld, L. Zettlemoyer, and O. Levy, "Spanbert: Improving pre-training by representing and predicting spans," *Transactions of the association for computational linguistics*, vol. 8, pp. 64–77, 2020.
- [255] J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le, "Finetuned language models are zero-shot learners," *arXiv preprint arXiv:2109.01652*, 2021.
- [256] G. Penedo, Q. Malartic, D. Hesslow, R. Cojocar, A. Cappelli, H. Alobeidli, B. Pannier, E. Almazrouei, and J. Launay, "The RefinedWeb dataset for Falcon LLM: outperforming curated corpora with web data, and web data only," *arXiv preprint arXiv:2306.01116*, 2023. [Online]. Available: <https://arxiv.org/abs/2306.01116>
- [257] S. Casper, X. Davies, C. Shi, T. K. Gilbert, J. Scheurer, J. Rando, R. Freedman, T. Korbak, D. Lindner, P. Freire *et al.*, "Open problems and fundamental limitations of reinforcement learning from human feedback," *arXiv preprint arXiv:2307.15217*, 2023.
- [258] Y. Wang, S. Mishra, P. Alipoormolabashi, Y. Kordi, A. Mirzaei, A. Arunkumar, A. Ashok, A. S. Dhanasekaran, A. Naik, D. Stap *et al.*, "Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks," *arXiv preprint arXiv:2204.07705*, 2022.
- [259] O. Honovich, T. Scialom, O. Levy, and T. Schick, "Unnatural instructions: Tuning language models with (almost) no human labor," *arXiv preprint arXiv:2212.09689*, 2022.
- [260] Y. Wang, Y. Kordi, S. Mishra, A. Liu, N. A. Smith, D. Khashabi, and H. Hajishirzi, "Self-instruct: Aligning language model with self generated instructions," *arXiv preprint arXiv:2212.10560*, 2022.
- [261] C. Xu, D. Guo, N. Duan, and J. McAuley, "Baize: An open-source chat model with parameter-efficient tuning on self-chat data," *arXiv preprint arXiv:2304.01196*, 2023.
- [262] O. Byambasuren, Y. Yang, Z. Sui, D. Dai, B. Chang, S. Li, and H. Zan, "Preliminary study on the construction of chinese medical knowledge graph," *Journal of Chinese Information Processing*, vol. 33, no. 10, pp. 1–9, 2019.
- [263] W.-L. Chiang, Z. Li, Z. Lin, Y. Sheng, Z. Wu, H. Zhang, L. Zheng, S. Zhuang, Y. Zhuang, J. E. Gonzalez *et al.*, "Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality," *See <https://vicuna.lmsys.org> (accessed 14 April 2023)*, 2023.
- [264] L. Pan, M. Saxon, W. Xu, D. Nathani, X. Wang, and W. Y. Wang, "Automatically correcting large language models: Surveying the landscape of diverse self-correction strategies," *arXiv preprint arXiv:2308.03188*, 2023.
- [265] S. Min, X. Lyu, A. Holtzman, M. Artetxe, M. Lewis, H. Hajishirzi, and L. Zettlemoyer, "Rethinking the role of demonstrations: What makes in-context learning work?" *arXiv preprint arXiv:2202.12837*, 2022.
- [266] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer, "Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 7871–7880.
- [267] Y. Bengio *et al.*, "From system 1 deep learning to system 2 deep learning," in *Neural Information Processing Systems*, 2019.
- [268] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, "Large language models are zero-shot reasoners," 2023.
- [269] S. Roy and D. Roth, "Solving general arithmetic word problems," *arXiv preprint arXiv:1608.01413*, 2016.
- [270] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano *et al.*, "Training verifiers to solve math word problems," *arXiv preprint arXiv:2110.14168*, 2021.
- [271] L. Weng, "Llm-powered autonomous agents," *lilianweng.github.io*, Jun 2023. [Online]. Available: <https://lilianweng.github.io/posts/2023-06-23-agent/>
- [272] X. Liu, H. Yu, H. Zhang, Y. Xu, X. Lei, H. Lai, Y. Gu, H. Ding, K. Men, K. Yang *et al.*, "Agentbench: Evaluating llms as agents," *arXiv preprint arXiv:2308.03688*, 2023.
- [273] J. He, C. Zhou, X. Ma, T. Berg-Kirkpatrick, and G. Neubig, "Towards a unified view of parameter-efficient transfer learning," *arXiv preprint arXiv:2110.04366*, 2021.
- [274] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, and S. Gelly, "Parameter-efficient transfer learning for nlp," in *International Conference on Machine Learning*. PMLR, 2019, pp. 2790–2799.
- [275] Z. Hu, Y. Lan, L. Wang, W. Xu, E.-P. Lim, R. K.-W. Lee, L. Bing, and S. Poria, "Llm-adapters: An adapter family for parameter-efficient finetuning of large language models," *arXiv preprint arXiv:2304.01933*, 2023.
- [276] X. L. Li and P. Liang, "Prefix-tuning: Optimizing continuous prompts for generation," *arXiv preprint arXiv:2101.00190*, 2021.
- [277] Y. Huang, Y. Cheng, A. Babna, O. Firat, D. Chen, M. Chen, H. Lee, J. Ngiam, Q. V. Le, Y. Wu *et al.*, "Gpipe: Efficient training of giant neural networks using pipeline parallelism," *Advances in neural information processing systems*, vol. 32, 2019.
- [278] A. Harlap, D. Narayanan, A. Phanishayee, V. Seshadri, N. Devanur, G. Ganger, and P. Gibbons, "Pipedream: Fast and efficient pipeline parallel dnn training," *arXiv preprint arXiv:1806.03377*, 2018.
- [279] J. Ren, S. Rajbhandari, R. Y. Aminabadi, O. Ruwase, S. Yang, M. Zhang, D. Li, and Y. He, "{ZeRO-Offload}: Democratizing {Billion-Scale} model training," in *2021 USENIX Annual Technical Conference (USENIX ATC 21)*, 2021, pp. 551–564.
- [280] S. Rajbhandari, O. Ruwase, J. Rasley, S. Smith, and Y. He, "Zero-infinity: Breaking the gpu memory wall for extreme scale deep learning," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, 2021, pp. 1–14.
- [281] E. Herrett, A. M. Gallagher, K. Bhaskaran, H. Forbes, R. Mathur, T. Van Staa, and L. Smeeth, "Data resource profile: clinical practice research datalink (cprd)," *International journal of epidemiology*, vol. 44, no. 3, pp. 827–836, 2015.
- [282] B. C. Wallace, S. Saha, F. Soboczenski, and I. J. Marshall, "Generating (Factual?) Narrative Summaries of RCTs: Experiments with Neural Multi-Document Summarization," in *Proceedings of AMIA Informatics Summit*, 2021.
- [283] J. DeYoung, I. Beltagy, M. van Zuylen, B. Kuehl, and L. L. Wang, "MS²: Multi-document summarization of medical studies," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 7494–7513. [Online]. Available: <https://aclanthology.org/2021.emnlp-main.594>
- [284] Y. Guo, W. Qiu, Y. Wang, and T. Cohen, "Automated lay language summarization of biomedical scientific reviews," *arXiv preprint arXiv:2012.12573*, 2020.

- [285] V. Gupta, P. Bharti, P. Nokhiz, and H. Karnick, "Sumpubmed: Summarization dataset of pubmed scientific article," in *Proceedings of the 2021 Conference of the Association for Computational Linguistics: Student Research Workshop*. Association for Computational Linguistics, 2021.
- [286] J. Bishop, Q. Xie, and S. Ananiadou, "Gencomparesum: a hybrid unsupervised summarization method using salience," in *Proceedings of the 21st workshop on biomedical language processing*, 2022, pp. 220–240.
- [287] L. L. Wang, K. Lo, Y. Chandrasekhar, R. Reas, J. Yang, D. Burdick, D. Eide, K. Funk, Y. Katsis, R. M. Kinney, Y. Li, Z. Liu, W. Merrill, P. Mooney, D. A. Murdick, D. Rishi, J. Sheehan, Z. Shen, B. Stilson, A. D. Wade, K. Wang, N. X. R. Wang, C. Wilhelm, B. Xie, D. M. Raymond, D. S. Weld, O. Etzioni, and S. Kohlmeier, "CORD-19: The COVID-19 open research dataset," in *Proceedings of the 1st Workshop on NLP for COVID-19 at ACL 2020*. Online: Association for Computational Linguistics, Jul. 2020. [Online]. Available: <https://www.aclweb.org/anthology/2020.nlpCOVID19-acl.1>
- [288] A. Ben Abacha and D. Demner-Fushman, "On the summarization of consumer health questions," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28th - August 2, 2019*.
- [289] S. Yadav, D. Gupta, and D. Demner-Fushman, "Chq-summ: A dataset for consumer healthcare question summarization," *arXiv preprint arXiv:2206.06581*, 2022.
- [290] M. Basaldella, F. Liu, E. Shareghi, and N. Collier, "Cometa: A corpus for medical entity linking in the social media," *arXiv preprint arXiv:2010.03295*, 2020.
- [291] G. Zeng, W. Yang, Z. Ju, Y. Yang, S. Wang, R. Zhang, M. Zhou, J. Zeng, X. Dong, R. Zhang *et al.*, "Meddialog: Large-scale medical dialogue datasets," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 9241–9250.
- [292] Z. Ju, S. Chakravorty, X. He, S. Chen, X. Yang, and P. Xie, "Covid-dialog: Medical dialogue datasets about covid-19," 2020.
- [293] M. Savery, A. B. Abacha, S. Gayen, and D. Demner-Fushman, "Question-driven summarization of answers to consumer health questions," *Scientific Data*, vol. 7, no. 1, p. 322, 2020.
- [294] B. Yu, Y. Li, and J. Wang, "Detecting causal language use in science findings," 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)*. Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 4664–4674. [Online]. Available: <https://aclanthology.org/D19-1473>
- [295] Y. Li, Z. Li, K. Zhang, R. Dan, S. Jiang, and Y. Zhang, "Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge," *Cureus*, vol. 15, no. 6, 2023.
- [296] D. Demner-Fushman, M. D. Kohli, M. B. Rosenman, S. E. Shooshan, L. Rodriguez, S. Antani, G. R. Thoma, and C. J. McDonald, "Preparing a collection of radiology examinations for distribution and retrieval," *Journal of the American Medical Informatics Association*, vol. 23, no. 2, pp. 304–310, 2016.
- [297] O. Pelka, S. Koitka, J. Rückert, F. Nensa, and C. M. Friedrich, "Radiology objects in context (roco): a multimodal image dataset," in *Intravascular Imaging and Computer Assisted Stenting and Large-Scale Annotation of Biomedical Data and Expert Label Synthesis: 7th Joint International Workshop, CVII-STENT 2018 and Third International Workshop, LABELS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Proceedings 3*. Springer, 2018, pp. 180–189.
- [298] S. Subramanian, L. L. Wang, S. Mehta, B. Bogin, M. van Zuylen, S. Parasa, S. Singh, M. Gardner, and H. Hajishirzi, "Medicat: A dataset of medical images, captions, and textual references," *arXiv preprint arXiv:2010.06000*, 2020.
- [299] J. Irvin, P. Rajpurkar, M. Ko, Y. Yu, S. Ciurea-Ilicus, C. Chute, H. Marklund, B. Haghgoo, R. Ball, K. Shpanskaya *et al.*, "Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 590–597.
- [300] A. Bustos, A. Pertusa, J.-M. Salinas, and M. De La Iglesia-Vaya, "Padchest: A large chest x-ray image dataset with multi-label annotated reports," *Medical image analysis*, vol. 66, p. 101797, 2020.
- [301] H. Zhi, B. Federico, Y. Mert, M. Thomas, J., and Z. James, "A visual-language foundation model for pathology image analysis using medical twitter," *Nature Medicine*, 2023.
- [302] Z. Yao, R. Y. Aminabadi, O. Ruwase, S. Rajbhandari, X. Wu, A. A. Awan, J. Rasley, M. Zhang, C. Li, C. Holmes, Z. Zhou, M. Wyatt, M. Smith, L. Kurilenko, H. Qin, M. Tanaka, S. Che, S. L. Song, and Y. He, "DeepSpeed-Chat: Easy, Fast and Affordable RLHF Training of ChatGPT-like Models at All Scales," *arXiv preprint arXiv:2308.01320*, 2023.
- [303] K. raj Kanakarajan, B. Kundumani, and M. Sankarasubbu, "Bioelectra: pretrained biomedical text encoder using discriminators," in *Proceedings of the 20th Workshop on Biomedical Language Processing*, 2021, pp. 143–154.
- [304] M. Basaldella, F. Liu, E. Shareghi, and N. Collier, "COMETA: A corpus for medical entity linking in the social media," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Online: Association for Computational Linguistics, Nov. 2020, pp. 3122–3137. [Online]. Available: <https://www.aclweb.org/anthology/2020.emnlp-main.253>
- [305] T. Chavan, S. Patankar, A. Kane, O. Gokhale, and R. Joshi, "A twitter bert approach for offensive language detection in marathi," *arXiv preprint arXiv:2212.10039*, 2022.
- [306] X. Zhang, Y. Malkov, O. Florez, S. Park, B. McWilliams, J. Han, and A. El-Kishky, "Twhin-bert: A socially-enriched pre-trained language model for multilingual tweet representations," *arXiv preprint arXiv:2209.07562*, 2022.
- [307] D. L. Wheeler, T. Barrett, D. A. Benson, S. H. Bryant, K. Canese, V. Chetvernin, D. M. Church, M. DiCuccio, R. Edgar, S. Federhen *et al.*, "Database resources of the national center for biotechnology information," *Nucleic Acids Research*, vol. 36, no. suppl_1, pp. D13–D21, 2007.
- [308] S. Liu, H. Yang, J. Li, and S. Kolmanič, "Preliminary study on the knowledge graph construction of chinese ancient history and culture," *Information*, vol. 11, no. 4, p. 186, 2020.
- [309] C. Li, M. Donizelli, N. Rodriguez, H. Dharuri, L. Endler, V. Chelliah, L. Li, E. He, A. Henry, M. I. Stefan *et al.*, "Biomodels database: An enhanced, curated and annotated resource for published quantitative kinetic models," *BMC systems biology*, vol. 4, no. 1, pp. 1–14, 2010.
- [310] D. S. Wishart, Y. D. Feunang, A. C. Guo, E. J. Lo, A. Marcu, J. R. Grant, T. Sajed, D. Johnson, C. Li, Z. Sayeeda *et al.*, "DrugBank 5.0: a major update to the DrugBank database for 2018," *Nucleic Acids Research*, vol. 46, no. D1, pp. D1074–D1082, 2018.
- [311] Y. Fang, X. Liang, N. Zhang, K. Liu, R. Huang, Z. Chen, X. Fan, and H. Chen, "Mol-instructions: A large-scale biomolecular instruction dataset for large language models," *arXiv preprint arXiv:2306.08018*, 2023.
- [312] J. W. Rae, S. Borgeaud, T. Cai, K. Millican, J. Hoffmann, F. Song, J. Aslanides, S. Henderson, R. Ring, S. Young *et al.*, "Scaling language models: Methods, analysis & insights from training gopher," *arXiv preprint arXiv:2112.11446*, 2021.
- [313] N. Bian, X. Han, L. Sun, H. Lin, Y. Lu, and B. He, "Chatgpt is a knowledgeable but inexperienced solver: An investigation of commonsense problem in large language models," *arXiv preprint arXiv:2303.16421*, 2023.
- [314] M. M. Amin, E. Cambria, and B. W. Schuller, "Will affective computing emerge from foundation models and general artificial intelligence? a first evaluation of chatgpt," *IEEE Intelligent Systems*, vol. 38, no. 2, pp. 15–23, 2023.
- [315] Y. Bang, S. Cahyawijaya, N. Lee, W. Dai, D. Su, B. Willie, H. Lovenia, Z. Ji, T. Yu, W. Chung *et al.*, "A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity," *arXiv preprint arXiv:2302.04023*, 2023.
- [316] Y.-T. Lin and Y.-N. Chen, "Llm-eval: Unified multi-dimensional automatic evaluation for open-domain conversations with large language models," *arXiv preprint arXiv:2305.13711*, 2023.
- [317] L. Wang, C. Lyu, T. Ji, Z. Zhang, D. Yu, S. Shi, and Z. Tu, "Document-level machine translation with large language models," *arXiv preprint arXiv:2304.02210*, 2023.
- [318] J. Liu, C. S. Xia, Y. Wang, and L. Zhang, "Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation," *arXiv preprint arXiv:2305.01210*, 2023.
- [319] Y. Zhuang, Q. Liu, Y. Ning, W. Huang, R. Lv, Z. Huang, G. Zhao, Z. Zhang, Q. Mao, S. Wang *et al.*, "Efficiently measuring the cognitive ability of llms: An adaptive testing perspective," *arXiv preprint arXiv:2306.10512*, 2023.
- [320] F. Xu, Q. Lin, J. Han, T. Zhao, J. Liu, and E. Cambria, "Are large language models really good logical reasoners? a comprehensive evaluation and beyond," *arXiv preprint arXiv:2306.09841*, 2023.
- [321] A. J. Thirunavukarasu, R. Hassan, S. Mahmood, R. Sanghera, K. Barzangi, M. El Mukashfi, and S. Shah, "Trialling a large language model (chatgpt) in general practice with the applied knowledge test:

- observational study demonstrating opportunities and limitations in primary care,” *JMIR Medical Education*, vol. 9, no. 1, p. e46599, 2023.
- [322] A. Gilson, C. W. Safranek, T. Huang, V. Socrates, L. Chi, R. A. Taylor, D. Chartash *et al.*, “How does chatgpt perform on the united states medical licensing examination? the implications of large language models for medical education and knowledge assessment,” *JMIR Medical Education*, vol. 9, no. 1, p. e45312, 2023.
- [323] T. H. Kung, M. Cheatham, A. Medenilla, C. Sillos, L. De Leon, C. Elepaño, M. Madriaga, R. Aggabao, G. Diaz-Candido, J. Maningo *et al.*, “Performance of chatgpt on usmle: Potential for ai-assisted medical education using large language models,” *PLoS digital health*, vol. 2, no. 2, p. e0000198, 2023.
- [324] D. Johnson, R. Goodman, J. Patrinely, C. Stone, E. Zimmerman, R. Donald, S. Chang, S. Berkowitz, A. Finn, E. Jahangir *et al.*, “Assessing the accuracy and reliability of ai-generated medical responses: an evaluation of the chat-gpt model,” 2023.
- [325] J. Holmes, Z. Liu, L. Zhang, Y. Ding, T. T. Sio, L. A. McGee, J. B. Ashman, X. Li, T. Liu, J. Shen *et al.*, “Evaluating large language models on a highly-specialized topic, radiation oncology physics,” *arXiv preprint arXiv:2304.01938*, 2023.
- [326] J. S. Samaan, Y. H. Yeo, N. Rajeev, L. Hawley, S. Abel, W. H. Ng, N. Srinivasan, J. Park, M. Burch, R. Watson *et al.*, “Assessing the accuracy of responses by the language model chatgpt to questions regarding bariatric surgery,” *Obesity surgery*, pp. 1–7, 2023.
- [327] D. Duong and B. D. Solomon, “Analysis of large-language model versus human performance for genetics questions,” *European Journal of Human Genetics*, pp. 1–3, 2023.
- [328] J. Chervenak, H. Lieman, M. Blanco-Breindel, and S. Jindal, “The promise and peril of using a large language model to obtain clinical information: Chatgpt performs strongly as a fertility counseling tool with limitations,” *Fertility and Sterility*, 2023.
- [329] N. Oh, G.-S. Choi, and W. Y. Lee, “Chatgpt goes to the operating room: evaluating gpt-4 performance and its potential in surgical education and training in the era of large language models,” *Annals of Surgical Treatment and Research*, vol. 104, no. 5, p. 269, 2023.
- [330] Z. Wang, R. Li, B. Dong, J. Wang, X. Li, N. Liu, C. Mao, W. Zhang, L. Dong, J. Gao *et al.*, “Can llms like gpt-4 outperform traditional ai tools in dementia diagnosis? maybe, but not today?” *arXiv preprint arXiv:2306.01499*, 2023.
- [331] A. Lahat, E. Shachar, B. Avidan, Z. Shatz, B. S. Glicksberg, and E. Klang, “Evaluating the use of large language model in identifying top research questions in gastroenterology,” *Scientific reports*, vol. 13, no. 1, p. 4164, 2023.
- [332] Q. Lyu, J. Tan, M. E. Zapadka, J. Ponnaturam, C. Niu, G. Wang, and C. T. Whitlow, “Translating radiology reports into plain language using chatgpt and gpt-4 with prompt learning: Promising results, limitations, and potential,” *arXiv preprint arXiv:2303.09038*, 2023.
- [333] I. Jahan, M. T. R. Laskar, C. Peng, and J. Huang, “Evaluation of chatgpt on biomedical tasks: A zero-shot comparison with fine-tuned generative transformers,” *arXiv preprint arXiv:2306.04504*, 2023.
- [334] M. Cascella, J. Montomoli, V. Bellini, and E. Bignami, “Evaluating the feasibility of chatgpt in healthcare: an analysis of multiple clinical and research scenarios,” *Journal of Medical Systems*, vol. 47, no. 1, p. 33, 2023.
- [335] T. Y. Zhuo, Z. Li, Y. Huang, Y.-F. Li, W. Wang, G. Haffari, and F. Shiri, “On robustness of prompt-based semantic parsing with large pre-trained language model: An empirical study on codex,” *arXiv preprint arXiv:2301.12868*, 2023.
- [336] Y. Zhao, T. Pang, C. Du, X. Yang, C. Li, N.-M. Cheung, and M. Lin, “On evaluating adversarial robustness of large vision-language models,” *arXiv preprint arXiv:2305.16934*, 2023.
- [337] C. Zhang, Y. Xie, H. Bai, B. Yu, W. Li, and Y. Gao, “A survey on federated learning,” *Knowledge-Based Systems*, vol. 216, p. 106775, 2021.
- [338] C. Wagner, M. Strohmaier, A. Olteanu, E. Kiciman, N. Contractor, and T. Eliassi-Rad, “Measuring algorithmically infused societies,” *Nature*, vol. 595, no. 7866, pp. 197–204, 2021.
- [339] Z.-X. Yong, R. Zhang, J. Zosa Forde, S. Wang, S. Cahyawijaya, H. Lovenia, G. Indra Winata, L. Sutawika, J. C. Blaise Cruz, L. Phan *et al.*, “Prompting multilingual large language models to generate code-mixed texts: The case of south east asian languages,” *arXiv e-prints*, pp. arXiv-2303, 2023.
- [340] R. Mao, Q. Liu, K. He, W. Li, and E. Cambria, “The biases of pre-trained language models: An empirical study on prompt-based sentiment analysis and emotion detection,” *IEEE Transactions on Affective Computing*, 2022.
- [341] R. J. Chen, J. J. Wang, D. F. Williamson, T. Y. Chen, J. Lipkova, M. Y. Lu, S. Sahai, and F. Mahmood, “Algorithmic fairness in artificial intelligence for medicine and healthcare,” *Nature Biomedical Engineering*, vol. 7, no. 6, pp. 719–742, 2023.
- [342] L. Seyyed-Kalantari, H. Zhang, M. B. McDermott, I. Y. Chen, and M. Ghassemi, “Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations,” *Nature medicine*, vol. 27, no. 12, pp. 2176–2182, 2021.
- [343] B. Delahunt, L. Egevad, H. Samaratunga, G. Martignoni, J. N. Nacey, and J. R. Strigley, “Gleason and fuhrman no longer make the grade,” *Histopathology*, vol. 68, no. 4, pp. 475–481, 2016.
- [344] A. Loupy, M. Mengel, and M. Haas, “Thirty years of the international banff classification for allograft pathology: the past, present, and future of kidney transplant diagnostics,” *Kidney International*, vol. 101, no. 4, pp. 678–691, 2022.
- [345] M. B. Zafar, I. Valera, M. G. Rogriguez, and K. P. Gummadi, “Fairness constraints: Mechanisms for fair classification,” in *Artificial intelligence and statistics*. PMLR, 2017, pp. 962–970.
- [346] R. Zemel, Y. Wu, K. Swersky, T. Pitassi, and C. Dwork, “Learning fair representations,” in *International conference on machine learning*. PMLR, 2013, pp. 325–333.
- [347] R. Mao, C. Lin, and F. Guerin, “End-to-end sequential metaphor identification inspired by linguistic theories,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 3888–3898.
- [348] M. Kim, O. Reingold, and G. Rothblum, “Fairness through computationally-bounded awareness,” *Advances in neural information processing systems*, vol. 31, 2018.
- [349] Y. Li, M. Du, R. Song, X. Wang, and Y. Wang, “A survey on fairness in large language models,” *arXiv preprint arXiv:2308.10149*, 2023.
- [350] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 2017.
- [351] N. McKenna, T. Li, L. Cheng, M. J. Hosseini, M. Johnson, and M. Steedman, “Sources of hallucination by large language models on inference tasks,” *arXiv preprint arXiv:2305.14552*, 2023.
- [352] Y. Li, Y. Du, K. Zhou, J. Wang, W. X. Zhao, and J.-R. Wen, “Evaluating object hallucination in large vision-language models,” *arXiv preprint arXiv:2305.10355*, 2023.
- [353] T. Y. Zhuo, Y. Huang, C. Chen, and Z. Xing, “Red teaming ChatGPT via jailbreaking: Bias, robustness, reliability and toxicity,” *arXiv preprint arXiv:2301.12867*, 2023.
- [354] S. A. Athaluri, S. V. Manthena, V. K. M. Kesapragada, V. Yarlagadda, T. Dave, and R. T. S. Duddumpudi, “Exploring the boundaries of reality: Investigating the phenomenon of artificial intelligence hallucination in scientific writing through ChatGPT references,” *Cureus*, vol. 15, no. 4, 2023.
- [355] A. Dero, K. Ghosh, and S. Ghosh, “How ready are pre-trained abstractive models and LLMs for legal case judgement summarization?” *arXiv preprint arXiv:2306.01248*, 2023.
- [356] A. Choudhury and O. Asan, “Impact of accountability, training, and human factors on the use of artificial intelligence in healthcare: Exploring the perceptions of healthcare practitioners in the us,” *Human Factors in Healthcare*, vol. 2, p. 100021, 2022.
- [357] A. Gudibande, E. Wallace, C. Snell, X. Geng, H. Liu, P. Abbeel, S. Levine, and D. Song, “The false promise of imitating proprietary LLMs,” *arXiv preprint arXiv:2305.15717*, 2023.
- [358] M. N. Dahlkemper, S. Z. Lahme, and P. Klein, “How do physics students evaluate artificial intelligence responses on comprehension questions? A study on the perceived scientific accuracy and linguistic quality,” *arXiv preprint arXiv:2304.05906*, 2023.
- [359] I. Habli, T. Lawton, and Z. Porter, “Artificial intelligence in health care: accountability and safety,” *Bulletin of the World Health Organization*, vol. 98, no. 4, p. 251, 2020.
- [360] F. Koto, J. H. Lau, and T. Baldwin, “Discourse probing of pretrained language models,” *arXiv preprint arXiv:2104.05882*, 2021.
- [361] L. Akhtyamova, P. Martínez, K. Verspoor, and J. Cardiff, “testing contextualized word embeddings to improve ner in spanish clinical case narratives,” *IEEE Access*, vol. 8, pp. 164 717–164 726, 2020.
- [362] H. Tan and M. Bansal, “Lxmert: Learning cross-modality encoder representations from transformers,” *arXiv preprint arXiv:1908.07490*, 2019.
- [363] H. Turbé, M. Bjelogrić, C. Lovis, and G. Mengaldo, “Evaluation of post-hoc interpretability methods in time-series classification,” *Nature Machine Intelligence*, vol. 5, no. 3, pp. 250–260, 2023.

- [364] S. Han, R. Mao, and E. Cambria, "Hierarchical attention network for explainable depression detection on Twitter aided by metaphor concept mappings," in *Proceedings of the 29th International Conference on Computational Linguistics (COLING)*. Gyeongju, Republic of Korea: International Committee on Computational Linguistics, 2022, pp. 94–104.
- [365] M. T. Ribeiro, S. Singh, and C. Guestrin, "“why should I trust you?” explaining the predictions of any classifier," in *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 2016, pp. 1135–1144.
- [366] R. Mao, C. Lin, and F. Guerin, "Word embedding and WordNet based metaphor identification and interpretation," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL)*, vol. 1. Melbourne, Australia: Association for Computational Linguistics, 2018, pp. 1222–1231. [Online]. Available: <https://aclanthology.org/P18-1113>
- [367] M. Ge, R. Mao, and E. Cambria, "Explainable metaphor identification inspired by conceptual metaphor theory," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 10, pp. 10681–10689, 2022.
- [368] S. H. Cho and K.-s. Shin, "Feature-weighted counterfactual-based explanation for bankruptcy prediction," *Expert Systems with Applications*, vol. 216, p. 119390, 2023.
- [369] W. Li, L. Zhu, R. Mao, and E. Cambria, "SKIER: A symbolic knowledge integrated model for conversational emotion recognition," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 11, pp. 13121–13129, 2023.
- [370] H. Chefer, S. Gur, and L. Wolf, "Transformer interpretability beyond attention visualization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 782–791.
- [371] J. Huang, H. Shao, and K. C.-C. Chang, "Are large pre-trained language models leaking your personal information?" 2022.
- [372] E. Lehman, S. Jain, K. Pichotta, Y. Goldberg, and B. C. Wallace, "Does bert pretrained on clinical notes reveal sensitive data?" *arXiv preprint arXiv:2104.07762*, 2021.
- [373] Y. Nakamura, S. Hanaoka, Y. Nomura, N. Hayashi, O. Abe, S. Yada, S. Wakamiya, and E. Aramaki, "Kart: Parameterization of privacy leakage scenarios from pre-trained language models," *arXiv preprint arXiv:2101.00036*, 2020.
- [374] L. Li, Y. Fan, M. Tse, and K.-Y. Lin, "A review of applications in federated learning," *Computers & Industrial Engineering*, vol. 149, p. 106854, 2020.
- [375] S. Gilbert, H. Harvey, T. Melvin, E. Vollebregt, and P. Wicks, "Large language model AI chatbots require approval as medical devices," *Nature Medicine*, pp. 1–3, 2023.
- [376] T. Minssen, E. Vayena, and I. G. Cohen, "The challenges for regulating medical use of ChatGPT and other large language models." *JAMA: Journal of the American Medical Association*, vol. 330, no. 4, 2023.
- [377] I. Alghanmi, L. E. Anke, and S. Schockaert, "Probing pre-trained language models for disease knowledge," in *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, 2021, pp. 3023–3033.
- [378] F. Petroni, T. Rocktäschel, P. Lewis, A. Bakhtin, Y. Wu, A. H. Miller, and S. Riedel, "Language models as knowledge bases?" *arXiv preprint arXiv:1909.01066*, 2019.
- [379] Z. Zhang, Z. Zeng, Y. Lin, H. Wang, D. Ye, C. Xiao, X. Han, Z. Liu, P. Li, M. Sun *et al.*, "Plug-and-play knowledge injection for pre-trained language models," *arXiv preprint arXiv:2305.17691*, 2023.
- [380] G. Michalopoulos, Y. Wang, H. Kaka, H. Chen, and A. Wong, "UmlsBERT: Clinical domain knowledge augmentation of contextual embeddings using the Unified Medical Language System Metathesaurus," in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Online: Association for Computational Linguistics, Jun. 2021, pp. 1744–1753. [Online]. Available: <https://aclanthology.org/2021.naacl-main.139>
- [381] T. Zhang, Z. Cai, C. Wang, M. Qiu, B. Yang, and X. He, "Smedbert: A knowledge-enhanced pre-trained language model with structured semantics for medical text mining," *arXiv preprint arXiv:2108.08983*, 2021.
- [382] A. Asai, S. Min, Z. Zhong, and D. Chen, "Retrieval-based language models and applications," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 6: Tutorial Abstracts)*. Toronto, Canada: Association for Computational Linguistics, Jul. 2023, pp. 41–46. [Online]. Available: <https://aclanthology.org/2023.acl-tutorials.6>
- [383] A. S. Rao, M. Pang, J. Kim, M. Kamineni, W. Lie, A. K. Prasad, A. Landman, K. Dryer, and M. D. Succi, "Assessing the utility of chatgpt throughout the entire clinical workflow," *medRxiv*, pp. 2023–02, 2023.
- [384] X. Liu, Y. Zheng, Z. Du, M. Ding, Y. Qian, Z. Yang, and J. Tang, "Gpt understands, too," 2021.