

Application of large language models in medicine

Fenglin Liu^{1,18,19}✉, Hongjian Zhou^{1,18,19}, Boyang Gu^{2,18}, Xinyu Zou^{3,18}, Jinfa Huang^{4,18}, Jing Wu^{5,18}, Yiru Li⁶, Sam S. Chen⁷, Yining Hua⁸, Peilin Zhou⁹, Junling Liu¹⁰, Chengfeng Mao¹¹, Chenyu You¹², Xian Wu¹³, Yefeng Zheng^{13,14}, Lei Clifton¹⁵, Zheng Li^{16,19}, Jiebo Luo^{4,19} & David A. Clifton^{1,17,19}

Abstract

Large language models (LLMs), such as ChatGPT, have received great attention owing to their capabilities for understanding and generating human language. Despite a trend in researching the application of LLMs in supporting different medical tasks (such as enhancing clinical diagnostics and providing medical education), a comprehensive assessment of their development, practical applications and outcomes in the medical space is still missing. Therefore, this Review aims to provide an overview of the development and deployment of LLMs in medicine, including the challenges and opportunities they face. In terms of development, we discuss the principles of existing medical LLMs, including their basic model structures, number of parameters, and sources and scales of data used for model development. In terms of deployment, we compare different LLMs across various medical tasks and with state-of-the-art lightweight models.

Sections

Introduction

The principles of medical LLMs

Medical tasks

Clinical applications

Challenges

Outlook

¹Institute of Biomedical Engineering, Department of Engineering Science, University of Oxford, Oxford, UK.

²Department of Computing, Imperial College London, London, UK. ³Department of Systems Design Engineering, University of Waterloo, Waterloo, Ontario, Canada. ⁴Department of Computer Science, University of Rochester, Rochester, NY, USA. ⁵Institute of Health Informatics, University College London, London, UK. ⁶Western University, London, Ontario, Canada. ⁷Department of Kinesiology, University of Georgia, Athens, GA, USA. ⁸Harvard T.H. Chan School of Public Health, Boston, MA, USA. ⁹Hong Kong University of Science and Technology (Guangzhou), Guangzhou, China. ¹⁰School of ECE, Peking University, Shenzhen, China. ¹¹Massachusetts Institute of Technology, Cambridge, MA, USA. ¹²Stony Brook University, Stony Brook, NY, USA. ¹³Jarvis Research Center, Tencent YouTu Laboratory, Beijing, China. ¹⁴Medical Artificial Intelligence Laboratory, Westlake University, Hangzhou, China.

¹⁵Applied Digital Health (ADH), Nuffield Department of Primary Care Health Sciences, University of Oxford, Oxford, UK. ¹⁶Amazon, Palo Alto, CA, USA. ¹⁷Oxford-Suzhou Centre for Advanced Research, Suzhou, China.

¹⁸These authors contributed equally: Fenglin Liu, Hongjian Zhou, Boyang Gu, Xinyu Zou, Jinfa Huang, Jing Wu.

¹⁹These authors jointly supervised this work: Fenglin Liu, Hongjian Zhou, Zheng Li, Jiebo Luo, David A. Clifton.

✉e-mail: fenglin.liu@eng.ox.ac.uk

Key points

- Existing medical large language models (LLMs), ranging from 110 million to 520 billion parameters, are mainly developed through pre-training, fine-tuning and prompting methods using large-scale medical corpora from different sources.
- Their performance is mostly evaluated based on exam-style question-answering tasks. Combining different fine-tuning and prompting methods enables LLMs to achieve comparable or even better results than experts.
- LLMs perform poorly in non-question-answering tasks without pre-set options, thus requiring further improvements before integration into real clinical decision-making processes.
- Medical LLMs are being adapted to various clinical applications, but large-scale clinical trials are still missing.
- Mitigating hallucinations; establishing robust data, benchmarks and metrics; and addressing ethical, safety and regulatory concerns through interdisciplinary collaborations will help to accelerate the integration of LLMs into clinic practice.

Introduction

The recently emerged general large language models (LLMs)^{1,2}, such as PaLM³, LLaMA^{4,5}, GPT series^{6,7} and ChatGLM⁸, have advanced the state of the art in various natural language processing (NLP) tasks, including text generation, text summarization and question answering (QA), as well as to adapt them to the medical domain^{9,10} (Box 1). For example, PaLM³ and GPT-4 (ref. 7), MedPaLM-2 (ref. 10) and MedPrompt¹¹ have achieved a competitive accuracy of 86.5 and 90.2 compared with human experts (87.0)¹² in the United States Medical Licensing Examination (USMLE)¹³, respectively. Based on publicly available general LLMs (such as LLaMA)⁵ a wide range of medical LLMs, including ChatDoctor¹⁴, MedAlpaca¹⁵, PMC-LLaMA¹², BenTsao¹⁶ and Clinical Camel¹⁷, have been introduced to assist medical professionals to improve patient care^{18,19}. Despite these advances, limitations remain: to begin with, many of these models mainly focus on medical dialogue and QA tasks, but their utility in clinical practice (such as electronic health records (EHRs)²⁰, discharge summary generation¹⁹, health education²¹ and care planning¹⁰) is often overlooked¹⁸. Moreover, current LLMs often fail to provide practical guidelines and are tested on a small number of users.

In this Review, we begin by analysing the principles underpinning current medical LLMs, providing detailed descriptions of their architecture, parameter scales and the data sets used during their development (Box 1). Next, we evaluate their performance across ten biomedical NLP tasks, including both discriminative and generative tasks. We then explore the practical deployment of medical LLMs in clinical settings by providing guidelines tailored for seven distinct clinical application scenarios. Finally, we discuss challenges such as the risk of generating factually inaccurate yet plausible outputs (hallucination) and the ethical, legal and safety implications; we also propose promising research directions. We argue for a comprehensive evaluation framework that assesses the trustworthiness of medical LLMs to ensure their responsible and effective use in the health-care domain. We also maintain a regularly updated list of practical guides on medical LLMs at <https://github.com/Al-in-Health/MedLLMsPracticalGuide>.

The principles of medical LLMs

Existing medical LLMs are mainly pre-trained from scratch, fine-tuned from existing general LLMs or directly obtained through prompting to align the general LLMs to the medical domain (Table 1). In this section, we introduce the principles of medical LLMs in terms of pre-training, fine-tuning and prompting (Fig. 1).

Pre-training

Pre-training typically involves training an LLM on a large corpus of medical texts, including both structured and unstructured texts, such as EHRs²², clinical notes²⁰ and medical literature²³. PubMed, MIMIC-III clinical notes²⁴ and PubMed Central (PMC) literature are three widely used medical corpora for medical LLM pre-training²⁵ (Table 1). Pre-training medical LLMs typically involves refining the following training objectives: masked language modelling, next-sentence prediction and next-token prediction (Box 1). For example, BERT-series models (such as BioBERT²⁶, PubMedBERT²⁷, ClinicalBERT²³ and GatorTron²⁰, which are originally derived from the general domain BERT or RoBERTa models) mainly adopt the masked language modelling and the next-sentence prediction for pre-training, whereas GPT-series models (such as BioGPT²⁸ and GatorTronGPT²²) mainly adopt the next-token prediction for pre-training (Box 1). After pre-training, medical LLMs can learn rich medical knowledge that can be leveraged to achieve strong performance on different medical tasks.

Fine-tuning

Training a medical LLM from scratch is expensive and time-consuming; one solution is to fine-tune the general LLMs with medical data using methods such as supervised fine-tuning (SFT), instruction fine-tuning (IFT) and parameter-efficient fine-tuning (PEFT)^{10,15,17} (Table 1).

Supervised fine-tuning. SFT aims to leverage high-quality medical corpus, which can be physician–patient conversations¹⁴, medical QA¹⁵ and knowledge graphs^{16,29}. The developed SFT data serve as a continuation of the pre-training data to further pre-train the general LLMs with the same training objectives, such as next-token prediction. DoctorGLM³⁰ and ChatDoctor¹⁴ are obtained by fine-tuning ChatGLM⁸ and LLaMA⁴, respectively, on the physician–patient dialogue data. MedAlpaca¹⁵, based on the general LLM Alpaca, is fine-tuned using over 160,000 medical QA pairs sourced from diverse medical corpora. Clinical Camel¹⁷ combines physician–patient conversations, clinical literature and medical QA pairs to refine the LLaMA-2 model⁵. In particular, Qilin-Med²⁹ and Zhongjing³¹ are obtained by incorporating the knowledge graph to perform fine-tuning on the Baichuan³¹ and LLaMA⁴, respectively.

Instruction fine-tuning. IFT generates instruction-based training data sets^{32,33} that typically comprise instruction–input–output triples, such as instruction–QA.

To ensure the high quality of training data and generalizability to different medical instructions and scenarios, MedPaLM-2 (ref. 10) invited qualified medical professionals for input, BenTsao¹⁶ developed knowledge-based instruction data from the knowledge graph (cMeKG)³⁴, whereas MedAlpaca¹⁵ incorporated both medical dialogues and QA pairs for IFT. Multimodal LLMs such as Med-Flamingo³⁵, LLaVA-Med³⁶ and Med-Gemini³⁷ have expanded the capabilities of LLMs to process complex and multimodal medical data; for example, Med-Flamingo³⁵ undergoes IFT on medical image–text data, learning to identify abnormalities and generate diagnostic reports;

Box 1 | Background of LLMs

The impressive performance of large language models (LLMs) can be attributed to Transformer-based language models, large-scale pre-training and scaling laws.

Language models

A language model^{184–186} is a probabilistic model that models the joint probability distribution of tokens (meaningful units of text, such as words or sub-words or morphemes) in a sequence, that is, the probabilities of how words and phrases are used in sequences. Therefore, it can predict the likelihood of a sequence of tokens given the previous tokens, which can be used to predict the next token in a sequence or to generate new sequences.

The Transformer architecture

The recurrent neural network^{185,187} has been widely used for language modelling by processing tokens sequentially and maintaining a vector named ‘hidden state’ that encodes the context of previous tokens. Nonetheless, sequential processing makes it unsuitable for parallel training and limits its ability to capture long-range dependencies, making it computationally expensive and hindering its learning ability for long sequences. The strength of the Transformer¹⁸⁸ lies in its fully attentive mechanism, which relies exclusively on the attention mechanism and eliminates the need for recurrence. When processing each token, the attention mechanism computes a weighted sum of the other input tokens, in which the weights are determined by the relevance between each input token and the current token. It enables the model to adaptively focus on different parts of the sequence to learn the joint probability distribution of tokens. Therefore, not only does Transformer enable modelling of long-text, but it also allows highly paralleled training¹⁸⁹, thus reducing training costs.

Large-scale pre-training

The LLMs are trained on massive corpora of unlabelled texts (for example, CommonCrawl, Wiki and Books) to learn rich linguistic knowledge and language patterns. The common training objectives are masked language modelling and next-token prediction. In masked language modelling, a portion of the input text is masked, and the model is tasked with predicting the masked text based on the remaining unmasked context, encouraging the model to capture the semantic and syntactic relationships between tokens¹⁸⁹. In next-token prediction, the model is required to predict the next token in a sequence, given the previous tokens⁶.

Scaling laws

LLMs are scaled-up versions of Transformer architecture¹⁸⁸ with increased numbers of Transformer layers, model parameters and volume of pre-training data. The ‘scaling laws’^{190,191} predict how much improvement can be expected in a model’s performance as its size increases (in terms of parameters, layers, data or the amount of training computed). The scaling laws proposed by OpenAI¹⁹⁰ show that to achieve optimal model performance, the budget allocation for model size should be larger than the data¹⁹².

The scaling laws proposed by Google DeepMind¹⁹¹ show that both model and data sizes should be increased in equal scales. The scaling

laws guide researchers to allocate resources and anticipate the benefits of scaling models.

General LLMs

Existing general LLMs can be divided into three categories based on their architecture (see the table).

Encoder-only LLMs consist of a stack of Transformer encoder layers; they utilize a bidirectional training strategy that enables them to integrate context from both the left and the right of a given token in the input sequence. This bidirectionality enables the models to achieve a deep understanding of the input sentences¹⁸⁹. Therefore, encoder-only LLMs are particularly suitable for language understanding tasks (such as sentiment analysis or document classification) for which the full context of the input is essential for accurate predictions. BERT¹⁸⁹ and DeBERTa¹⁹³ are the representative encoder-only LLMs.

Decoder-only LLMs use a stack of Transformer decoder layers and are characterized by their unidirectional (left-to-right) processing of text, enabling them to generate language sequentially. This architecture is trained unidirectionally using the next-token prediction training objective to predict the next token in a sequence, given all

Model structures	Models	Number of parameters	Pre-train data scale
Encoder-only	BERT ¹⁸⁹	110M/340M	3.3B tokens
	RoBERTa ⁸⁶	355M	161GB
	DeBERTa ¹⁹³	1.5B	160GB
Decoder-only	GPT-2 (ref. 194)	1.5B	40GB
	Vicuna ¹⁹⁵	7B/13B	LLaMA + 70k dialogues
	Alpaca	7B/13B	LLaMA + 52k IFT
	Mistral ¹⁹⁶	7B	–
	LLaMA ⁴	7B/13B/33B/65B	1.4T tokens
	LLaMA-2 (ref. 5)	7B/13B/34B/70B	2T tokens
	LLaMA-3	8B/70B	15T tokens
	GPT-3 (ref. 6)	6.7B/13B/175B	300B tokens
	Qwen ¹⁹⁷	1.8B/7B/14B/72B	3T tokens
	PaLM ³	8B/62B/540B	780B tokens
	FLAN-PaLM ⁷⁴	540B	–
	Gemini (Bard)	–	–
	GPT-3.5 (ref. 143)	–	–
	GPT-4 (ref. 7)	–	–
	Claude-3	–	–
Encoder-decoder	BART ¹⁹⁸	140M/400M	160GB
	ChatGLM ⁸	6.2B	1T tokens
	T5 (ref. 74)	11B	1T tokens
	FLAN-T5 (ref. 74)	3B/11B	780B tokens
	UL2 (ref. 199)	19.5B	1T tokens
	GLM ³	130B	400B tokens

B, billion; IFT, instruction fine-tuning; M, million; T, trillion.

(continued from previous page)

the previous tokens. After training, the decoder-only LLMs generate sequences autoregressively (that is, token by token). The examples are the GPT-series developed by OpenAI^{6,7}, the LLaMA-series developed by Meta^{4,5} and the PaLM developed by Google³.

Encoder–decoder LLMs are designed to simultaneously process input sequences and generate output sequences. They consist of a stack of bidirectional Transformer encoder layers followed by a

stack of unidirectional Transformer decoder layers. The encoder processes and understands the input sequences whereas the decoder generates the output sequences^{8,74}. Representative examples of encoder–decoder LLMs include Flan-T5 (ref. 67) and ChatGLM⁸. Specifically, ChatGLM has 6.2 billion parameters and is a conversational open-source LLM optimized for Chinese to support Chinese–English bilingual question answering.

LLaVA-Med's³⁶ two-stage IFT process involves aligning medical concepts across visual and textual modalities, followed by fine-tuning on different medical instructions; whereas Med-Gemini's³⁷ IFT uses a curated data set of medical instructions and multimodal data, enabling it to comprehend complex medical concepts, procedures and diagnostic reasoning. Conversely, MAIRA-1 (ref. 38) and RadFM³⁹ are two multimodal LLMs specifically designed for radiology applications. MAIRA-1 (ref. 38) undergoes IFT on a data set of radiology instructions and corresponding medical images to analyse radiological images and generate accurate diagnostic reports. By contrast, RadFM³⁹ leverages a pre-training approach on a large corpus of radiology-specific image–text data, followed by IFT on different sets of radiology instructions. These models' multimodal IFT approaches enable them to bridge the gap between visual and textual medical information, perform a wide range of medical tasks accurately, and generate context-aware responses to complex medical queries. The goal of IFT is to improve the model's ability to follow various human and task instructions, align their outputs with the medical domain and produce a specialized medical LLM. Therefore, the main difference between SFT and IFT is that the former focuses on injecting medical knowledge into a general LLM through continued pre-training, thus improving its ability to understand the medical text and accurately predict the next token. By contrast, IFT aims to improve the model's 'instruction-following' ability and adjust its outputs to match the given instructions rather than accurately predicting the next token as in SFT³². As a result, SFT emphasizes the 'quantity' of training data whereas IFT emphasizes their 'quality and diversity'. Combination of both fine-tuning methods have also been attempted^{29,31,40} to simultaneously improve LLM's ability to understand medical knowledge and follow human instructions, resulting in improved overall task performance. For example, Zhongjing³¹ first collected a large amount of real medical corpus for SFT and then incorporated multi-turn dialogue as instruction data to perform IFT.

Parameter-efficient fine-tuning. PEFT aims to substantially reduce computational and memory requirements for fine-tuning general LLMs. The main concept is to keep most of the parameters in pre-trained LLMs unchanged by fine-tuning only the smallest subset of parameters (or additional parameters). Commonly used PEFT techniques include low-rank adaptation (LoRA)⁴¹, prefix tuning⁴² and adapter tuning⁴³.

In contrast to fine-tuning full-rank weight matrices, LoRA preserves the parameters of the original LLMs and only adds trainable low-rank matrices into the self-attention module of each Transformer layer⁴¹. Therefore, LoRA can substantially reduce the number of trainable parameters and improve the efficiency of fine-tuning while preserving the ability of the fine-tuned LLM to capture the characteristics of the tasks. Prefix tuning instead adds a small set of continuous task-specific vectors (that is, 'prefixes') to the input of each Transformer layer^{1,42}.

These prefixes serve as the additional context to guide the generation of the model without changing the original pre-trained parameter weights. Adapter tuning introduces small neural network modules, known as adapters, into each Transformer layer of the pre-trained LLMs⁴⁴. These adapters are fine-tuned while the original model parameters are kept frozen⁴⁴ for flexible and efficient fine-tuning. Compared with full-rank fine-tuning, the number of trainable parameters introduced by adapters is relatively small, yet they enable the LLMs to adapt to clinical text classification and understanding tasks effectively, achieving similar performance with fewer trainable parameters. In general, PEFT is valuable for developing domain-specific LLMs (such as medical) owing to its ability to reduce computational demands while maintaining the model performance. For example, medical LLMs DoctorGLM³⁰, MedAlpaca⁴⁵, Baize-Healthcare⁴⁵, Zhongjing³¹, CPLLM⁴⁶ and Clinical Camel¹⁷ adopted the LoRA⁴¹ to align the general LLMs to the medical domain.

Prompting

Fine-tuning considerably reduces computational costs compared with pre-training, but it still requires model training and collections of high-quality data sets. By contrast, the 'prompting' methods align general LLMs to the medical domain without training any model parameters (such as MedPrompt¹¹ and MedPaLM⁹). Popular prompting methods include in-context learning (ICL), chain-of-thought (CoT) prompting, prompt tuning and retrieval-augmented generation (RAG).

In-context learning. ICL aims to give direct instructions to prompt the LLM to perform a task. ICL consists of four processes: task understanding, context learning, knowledge reasoning and answer generation. First, the model must understand the specific requirements and goals of the task. Second, the model learns to understand the contextual information related to the task with argument context. Third, it uses the model's internal knowledge and reasoning capabilities to understand the patterns and logic in the example. Fourth, the model generates the task-related answers. The advantage of ICL is that it does not require a large amount of labelled data for fine-tuning. Based on the type and number of input examples, ICL can be divided into three categories⁴⁷: one-shot prompting, in which only one example and task description can be entered; few-shot prompting, which allows the input of multiple instances and task descriptions; and zero-shot prompting, in which only task descriptions can be entered. ICL enables the LLMs to make task predictions based on contexts augmented with a few examples and task demonstrations. It allows the LLMs to learn from these examples or demonstrations to perform the task and follow the given examples to give corresponding answers⁶. For example, MedPaLM⁹ substantially improves the task performance by providing the general LLM, PaLM³, with a small number of task examples such as medical QA pairs.

Table 1 | Summary of medical-domain LLMs

Model development	Models	Number of parameters	Data scale	Data source
Pre-training	BioBERT ²⁶	110M	18B tokens	PubMed + PMC
	PubMedBERT ²⁷	110M/340M	3.2B tokens	PubMed + PMC
	SciBERT ¹⁶¹	110M	3.17B tokens	Literature
	NYUTron ⁷²	110M	7.25M notes, 4.1B tokens	NYU Notes ⁷²
	ClinicalBERT ²³	110M	112k clinical notes	MIMIC-III ²⁴
	BioM-ELECTRA ¹⁶²	110M/335M	–	PubMed
	BioMed-RoBERTa ¹⁶³	125M	7.55B tokens	S2ORC ¹⁶³
	BioLinkBERT ¹⁶⁴	110M/340M	21GB	PubMed
	BlueBERT ¹⁶⁵	110M/340M	>4.5B tokens	PubMed + MIMIC-III ²⁴
	SciFive ¹⁶⁶	220M/770M	–	PubMed + PMC
	ClinicalT5 (ref. 167)	220M/770M	2M clinical notes	MIMIC-III ²⁴
	MedCPT ¹⁶⁸	330M	255M articles	PubMed
	DRAGON ¹⁶⁹	360M	6GB	BookCorpus ¹⁶⁹
	BioGPT ²⁸	1.5B	15M articles	PubMed
	BioMedLM ¹⁷⁰	2.7B	110GB	Pile ¹⁷⁰
	OphGLM ¹⁷¹	6.2B	20k dialogues	MedDialog ¹⁷¹
	GatorTron ²⁰	8.9B	>82B tokens + 6B tokens	EHRs ²⁰ + PubMed
			2.5B tokens + 0.5B tokens	Wiki + MIMIC-III ²⁴
	GatorTronGPT ²²	5B/20B	277B tokens	EHRs ²²
Fine-tuning	DoctorGLM ³⁰	6.2B	323MB dialogues	CMD
	BianQue ¹⁷²	6.2B	2.4M dialogues	BianQueCorpus ¹⁷²
	ClinicalGPT ¹⁷³	7B	96k EHRs + 100k dialogues	MD-EHR ¹⁷³ + MedDialog ¹⁷³
			192 medical QA	VariousMedQA ¹⁷³
	Qilin-Med ²⁹	7B	3GB	ChiMed ²⁹
	ChatDoctor ¹⁴	7B	110k dialogues	HealthCareMagic + iCliniq
	BenTsao ¹⁶	7B	8k instructions	CMeKG-8K ³⁴
	HuatuoGPT ¹⁷⁴	7B	226k instructions and dialogues	Hybrid SFT ¹⁷⁴
	Baize-healthcare ⁴⁵	7B	101K dialogues	Quora + MedQuAD ⁴⁵
	BioMedGPT ¹⁷⁵	10B	>26B tokens	S2ORC ¹⁷⁵
	MedAlpaca ¹⁵	7B/13B	160k medical QA	Medical Meadow ¹⁵
	AlpaCare ⁴⁰	7B/13B	52k instructions	MedInstruct-52k ⁴⁰
	Zhongjing ³¹	13B	70k dialogues	CMtMedQA ³¹
	PMC-LLaMA ¹²	13B	79.2B tokens	Books + Literature ¹² + MedC-I ¹²
	CPLLM ⁴⁶	13B	109k EHRs	eICU-CRD ⁴⁶ + MIMIC-IV
	Med42 (ref. 176)	7B/70B	250M tokens	PubMed + MedQA ¹³ + OpenOrca
	MEDITRON ¹⁷⁷	7B/70B	48.1B tokens	PubMed + Guidelines ¹⁷⁷
	OpenBioLLM	8B/70B	–	–
	MedLlama3-v20	8B/70B	–	–
	Clinical Camel ¹⁷	13B/70B	70k dialogues + 100k articles	ShareGPT + PubMed
			4k medical QA	MedQA ¹³
	MedPaLM-2 (ref. 10)	340B	193k medical QA	MultiMedQA ¹⁰
	Med-Flamingo ³⁵	–	600k pairs	Multimodal Textbook ³⁵ + PMC-OA ³⁵
				VQA-RAD ³⁵ + PathVQA ¹⁷⁷

Table 1 (continued) | Summary of medical-domain LLMs

Model development	Models	Number of parameters	Data scale	Data source
Fine-tuning (continued)	LLaVA-Med ³⁶	–	660k pairs	PMC-15M ³⁶ + VQA-RAD ³⁶ SLAKE ³⁶ + PathVQA ¹⁷⁷
	MAIRA-1 (ref. 38)	–	337k pairs	MIMIC-CXR ¹⁷⁸
	RadFM ³⁹	–	32M pairs	MedMD ³⁹
	Med-Gemini ¹⁷⁹	–	–	MedQA-R&RS ³⁷ + MultiMedQA ¹⁰ MIMIC-III ²⁴ + MultiMedBench ¹⁷⁹
Prompting	CodeX ¹⁸⁰	GPT-3.5/LLaMA-2	CoT ⁴⁸	–
	DeID-GPT ⁴⁹	ChatGPT/GPT-4	CoT ⁴⁸	–
	ChatCAD ⁸⁰	ChatGPT	ICL	–
	Dr. Knows ⁷³	ChatGPT	ICL	UMLS
	MedPaLM ⁹	PaLM (540B)	CoT and ICL	MultiMedQA ¹⁰
	MedPrompt ¹¹	GPT-4	CoT and ICL ⁴⁸	–
	Chat-Orthopedist ⁶²	ChatGPT	RAG	PubMed + Guidelines + UpToDate + Dymene
	QA-RAG ⁶¹	ChatGPT	RAG	FDA QA ⁶¹
	Almanac ⁶⁰	ChatGPT	RAG and CoT	Clinical QA ⁶⁰
	Oncology-GPT-4 (ref. 176)	GPT-4	RAG and ICL	Oncology guidelines from ASCO and ESMO

ASCO, American Society of Clinical Oncology; B, billion; CoT, chain-of-thought; EHR, electronic health record; ESMO, European Society for Medical Oncology; ICL, in-context learning; LLM, large language model; M, million; PMC, PubMed Central; QA, question answering; RAG, retrieval-augmented generation; SFT, supervised fine-tuning; UMLS, Unified Medical Language System.

Chain-of-thought prompting. CoT further improves the accuracy and logic of model output compared with ICT. Specifically, through prompting words, CoT aims to prompt the model to generate intermediate steps or paths of reasoning when dealing with downstream (complex) problems⁴⁸. Moreover, CoT can be combined with few-shot prompting by giving reasoning examples, thus enabling medical LLMs to give reasoning processes when generating responses. CoT improves model performance for tasks involving complex reasoning (such as medical QA)^{9,10}. Medical LLMs such as DeID-GPT⁴⁹, MedPaLM⁹ and MedPrompt¹¹ use CoT prompting to assist them in simulating a diagnostic thought process, thus providing more transparent and interpretable predictions or diagnoses. In particular, MedPrompt¹¹ directly prompts a general LLM, GPT-4 (ref. 7), to outperform the fine-tuned medical LLMs on medical QA without training any model parameters.

Prompt tuning. Prompt tuning aims to improve the model performance by implementing both prompting and fine-tuning techniques⁵⁰. The prompt tuning method introduces learnable prompts (that is, trainable continuous vectors) that can be optimized or adjusted during the fine-tuning process to better adapt to different medical scenarios and tasks. Therefore, they provide a more flexible way of prompting LLMs than the ‘prompting alone’ methods that use discrete and fixed prompts. In contrast to traditional fine-tuning methods that train all model parameters, prompt tuning only tunes a very small set of parameters (that is, less than 3% of the total model parameters) associated with the prompts themselves, instead of extensively training the model parameters.

MedPaLM⁹ and MedPaLM-2 (ref. 10) have combined all the above prompting methods resulting in an ‘instruction prompt tuning’ to improve performances on various medical QA data sets. Using the MedQA data set for the USMLE, MedPaLM-2 (ref. 10) achieves a

competitive overall accuracy of 86.5% compared with human experts (87.0%), surpassing MedPaLM⁹ by a large margin (19%).

Retrieval-augmented generation. RAG improves the performance of LLMs by integrating external knowledge into the generation process to minimize hallucinations, obscure reasoning processes and reliance on outdated information⁵¹. RAG consists of three main components: retrieval, augmentation and generation. The retrieval component uses various indexing strategies and input query processing techniques to search and rank relevant information from an external knowledge base. The retrieved external data are then augmented into the LLM’s prompt, providing additional context and grounding for the generated response. By directly updating the external knowledge base, RAG mitigates the risk of amnesia (that is, catastrophic forgetting⁵²) associated with model weight modifications, making it particularly suitable for domains with low error tolerance and rapidly evolving information, such as the medical field. In contrast to traditional fine-tuning methods, RAG incorporates new medical information without compromising the model’s previously acquired knowledge. MIRAGE⁵³ is the first benchmark to be proposed based on medical information RAG, including 7,663 questions from five medical QA data sets.

In RAG, retrieval can be achieved by calculating the similarity between the textual embeddings of the question and document chunks. Textual embeddings are vector representations of text, which are encoded by embedding models (such as AngIE⁵⁴, Voyage⁵⁵ and BGE⁵⁶) into numerical vectors that capture the semantic information of the text. These embedding models serve as the foundation for ensuring that the retrieval step accurately captures the semantic meaning of the question, retrieves accurate documents and ultimately enhances the quality of the generated responses. In addition to embedding, the retrieval process can be optimized by adaptive retrieval, recursive

retrieval and iterative retrieval, among others^{57–59}. Adaptive retrieval dynamically adjusts retrieval strategies (for example, prioritizing specific document chunks) based on the complexity or specificity of the question, improving relevance and efficiency. Recursive retrieval refines retrieval results in multiple rounds, during which the retrieved documents are re-evaluated and refined in successive rounds to progressively narrow down the results. Iterative retrieval incorporates prior retrieval results to improve subsequent retrieval attempts. For example, Almanac⁶⁰ is a large language framework augmented with retrieval capabilities for medical guidelines and treatment recommendations, surpassing ChatGPT on clinical scenario evaluations, particularly in terms of completeness and safety. Another example is QA-RAG⁶¹, which utilizes RAG with LLM for pharmaceutical regulatory tasks, in which the model searches for relevant guideline documents and provides answers based on the retrieved guidelines. Chat-Orthopedist⁶², a retrieval-augmented LLM, assists people with adolescent idiopathic scoliosis to prepare for discussions with clinicians by providing accurate and comprehensible responses to patient inquiries, leveraging the adolescent idiopathic scoliosis domain knowledge from PubMed clinical papers, Scoliosis Research Society's practice guidelines and UpToDate.com.

Medical tasks

In this section, we introduce two popular types of medical machine-learning tasks: generative and discriminative tasks, including ten representative tasks that further build up clinical applications (Fig. 2). A detailed definition of the task and performance comparisons across different LLMs can be found in the Supplementary information.

Discriminative tasks

Discriminative tasks categorize or differentiate data (for example, structured or unstructured text) into specific classes or categories based on given input data. The representative tasks include QA, entity extraction, relation extraction, text classification, natural language inference, semantic textual similarity and information retrieval (Supplementary Information). The typical input for discriminative

tasks can be medical questions, clinical notes, medical documents, research papers and patient EHRs. The output can be labels, categories, extracted entities, relationships or answers to specific questions, which are often structured and categorized information derived from the input text. In existing LLMs, the discriminative tasks are widely studied and used to make predictions and extract information from input text. Entity extraction can automatically identify and categorize critical information (that is, entities) such as symptoms, medications, diseases, diagnoses and lab results from patient EHRs, thus assisting in organizing and managing patient data. Entity linking is a subsequent step to entity extraction, in which the identified entities are mapped to entries in a structured knowledge base or a standardized terminology system. For example, an extracted entity such as 'diabetes' can be linked to a specific concept in SNOMED CT, Unified Medical Language System (UMLS) or International Classification of Diseases (ICD) codes. This process ensures that the extracted entities are standardized, enabling interoperability across different systems and facilitating more precise medical analysis and data management.

Generative tasks

Different from discriminative tasks that focus on understanding and categorizing the input text, generative tasks require a model to generate new, fluent and accurate text based on given inputs. These tasks include medical text summarization^{63,64}, medical text generation²⁸ and medical text simplification⁶⁵.

For medical text summarization, the input and output are typically a long and detailed medical text (for example, 'Findings' in radiology reports) and a concise summarized text (such as the 'Impression' in radiology reports), respectively. In medical text generation (for example, discharge instruction generation⁶⁶), the input can be medical conditions, symptoms, patient demographics or even a set of medical notes or test results. The output can be a diagnosis recommendation of a medical condition, personalized instructional information or health advice for the patient to manage their condition outside the hospital.

Medical text simplification⁶⁵ aims to generate a simplified version of the complex medical text by, for example, clarifying and explaining

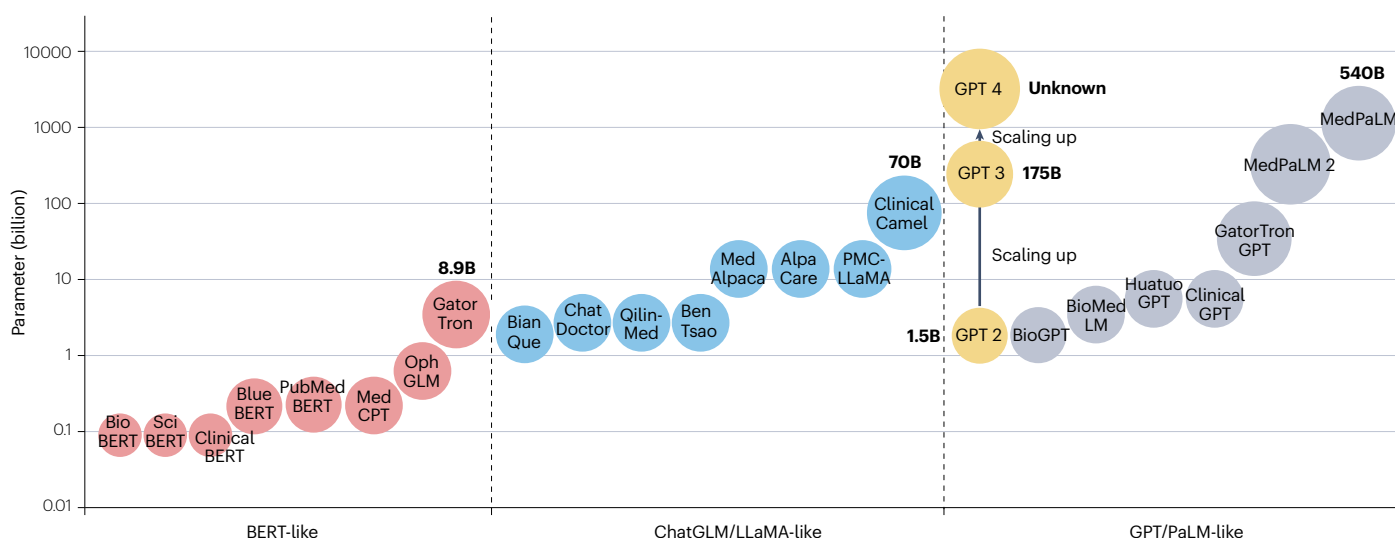
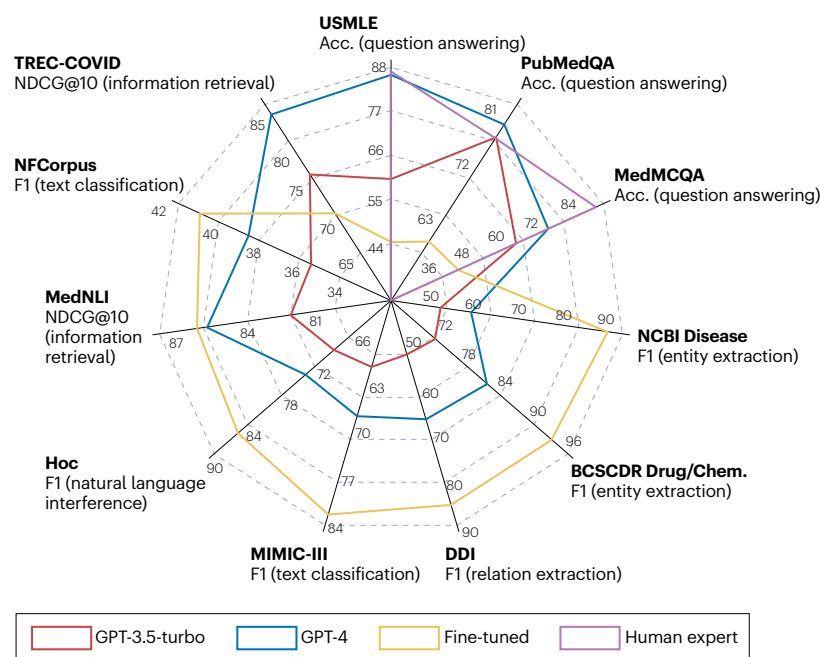


Fig. 1 | Model size for different medical LLMs. Data from Table 1 were used to illustrate the development of model sizes for medical large language models (LLMs) in different model architectures, that is, BERT, ChatGLM/LLaMA and GPT/PaLM. B, billion.

a Medical large language model for medical discriminative tasks



b Medical large language model for medical generative tasks

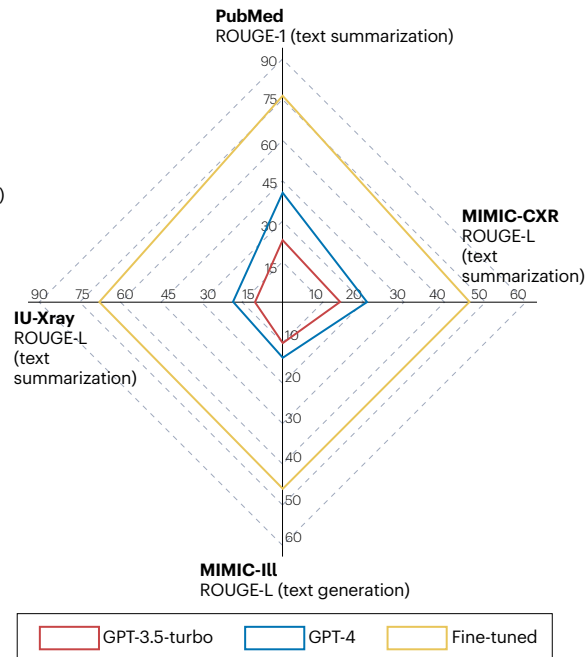


Fig. 2 | Performance comparison. Performance (Dataset-Metric (task)) comparison between the GPT-3.5 turbo, GPT-4, state-of-the-art task-specific lightweight models (fine-tuned), and human experts, on seven medical tasks across 11 data sets. All data presented in our figures originate from published and peer-reviewed literature (Supplementary Information). NDCG@10 represents the normalized discounted cumulative gain at rank 10, evaluating ranking quality

by considering the relevance and position of retrieved items. Recall-oriented understudy for gisting evaluation (ROUGE-1) is a score based on unigram (that is, 1-gram, word-level) overlap between predicted and reference texts. ROUGE-L is a ROUGE score considering the longest (L) common subsequence between predicted and reference texts. Acc., accuracy; F1, harmonic mean of precision and recall.

medical terms and therefore improve readability. Specifically, complicated or opaque words are replaced, complex syntactic structures improved and rare concepts explained⁶⁷.

Performance comparisons

General LLMs such as GPT-3.5-turbo and GPT-4 (ref. 7) have achieved strong performance on existing medical machine-learning tasks, and they even outperform existing task-specific fine-tuned models and human experts for QA tasks (MedQA (USMLE)¹³, PubMedQA⁶⁸ and MedMCQA⁶⁹) (Figs. 2 and 3). However, existing general LLMs perform worse than the task-specific fine-tuned models on non-QA discriminative tasks (Fig. 2). For example, task-specific fine-tuned model BioBERT²⁶ achieves an F1 score of 89.36 compared with 56.73 achieved by GPT-4 on the entity extraction task. A higher F1 score indicates better performance in balancing precision (correctly identified entities) and recall (completeness of identified entities). Notably, GPT-4 underperforms against task-specific lightweight models on all data sets in generative tasks. We hypothesize that general LLMs are strong in QA because these are close-ended tasks, that is, the correct answer is already provided by multiple candidates. By contrast, most non-QA tasks are open-ended, in which the model must predict the correct answer from a large pool of possible candidates, or even without any candidates being provided.

Real-world clinical practice often involves answering open-ended questions without pre-set options and is therefore different from the structured nature of exam-taking. Therefore, medical LLMs should be evaluated not only on medical QA tasks but also non-QA ones.

Clinical applications

Current medical LLMs are still in the research-and-development stage, with few clinical trials ongoing⁷⁰. In this section, we discuss the application of LLMs in medicine in detail (Table 2 and Fig. 4).

Medical decision-making

Medical decision-making, including diagnosis, prognosis, treatment suggestion, risk prediction and clinical trial matching, relies heavily on the synthesis and interpretation of vast amounts of information from different sources, such as patient medical histories, clinical data and the latest medical literature. LLMs can rapidly process and understand such data and potentially assist health-care professionals in making more informed and legally sound decisions across a wide range of clinical scenarios^{18,71}. For example, in medical diagnosis, LLMs can assist practitioners in analysing medical data from tests and self-described symptoms to conclude the most likely health problem⁷¹. Similarly, LLMs can support treatment planning by providing personalized recommendations based on the latest clinical evidence and patient-specific information¹⁸. Furthermore, LLMs can contribute to prognosis and risk prediction by identifying patterns and risk factors from large-scale patient data, enabling more accurate and timelier interventions⁷².

For example, Dr. Knows⁷³ can integrate knowledge graphs from UMLS to improve diagnosis prediction and provide treatment suggestions. This approach involves fine-tuning an encoder-decoder LLM T5 (ref. 74) with extracted diagnoses as prompts and using zero-shot prompting for LLMs like ChatGPT. Alternatively, models like DDx

PaLM-2 (ref. 75), based on IFT general LLMs (such as Google's PaLM-2) with extensive medical data sets MultiMedQA¹⁰ and MIMIC-III²⁴, enables interactive diagnosis assistance, in which humans and LLMs can iteratively communicate to arrive at a final diagnosis. NYUTron⁷² is pretrained and supervised fine-tuned on various NYU hospitals and can perform three clinical tasks (in-patient mortality prediction, comorbidity index prediction and readmission prediction) and two operational tasks (insurance claim denial prediction and inpatient length of stay prediction). Similarly, Foresight⁷⁶ is trained on UK hospital patient data and can be used for forecasting the risk of disorders, performing differential diagnoses and suggesting medications. For clinical trial matching, TrialGPT⁷⁷ predicts criterion-level eligibility with faithful explanations, reducing screening time for human experts. Ongoing clinical trials (NCT06002425)⁷⁸ across Germany, Italy, China and the USA are examining the accuracy and efficiency of LLMs in clinical decision-making and treatment recommendations for gastrointestinal cancers, or whether medical laypeople make better decisions when using LLMs (DRKS00033775)⁷⁹.

Evaluating LLM-based medical diagnosis systems requires task-specific approaches. For general diagnostic accuracy, metrics like the area under the curve (AUC), which measures classification performance, as well as precision, recall and F1 score are used with annotated data sets^{75–77}. Diagnostic information can also be evaluated using text summarization and medical concept extraction performance⁷³. One limitation of using LLMs is their heavy reliance on subjective text inputs from patients; moreover, because LLMs are text-based, they cannot analyse medical images, an essential component of diagnostic assessments⁷⁸. However, they can help with diagnosis as a logical

reasoning tool for improving accuracy in other vision-based models; for example, in ChatCAD⁸⁰, images are first fed into an existing computer-aided diagnosis model to obtain tensor outputs. These outputs are translated into natural language, which is then fed into ChatCAD to summarize results and formulate diagnoses. ChatCAD achieves a recall score of 0.781, substantially higher than that of the state-of-the-art task-specific model (0.382).

Clinical coding

Clinical coding, such as the ICD, medication coding and procedure coding help to standardize diagnostic, procedural and treatment information. These codes are essential for tracking health metrics, treatment outcomes, and billing and reimbursement processes; however, their manual entry is time-consuming and prone to errors. LLMs can automate this process by extracting relevant medical terms from clinical notes and assigning corresponding codes, including ICD codes^{81–84}, medication codes (such as the [National Drug Code Directory](#)) and procedure codes (for example, Current Procedural Terminology)⁸⁵. For example, PLM-ICD⁸¹, which builds upon the RoBERTa model⁸⁶, an optimized version of BERT, when fine-tuned for ICD coding, can understand medical terms, and it achieves strong coding performance on 70,539 notes from the MIMIC-II and MIMIC-III data sets²⁴. Other examples include DRG-LLaMA⁸², which leverages the LLaMA model and applies PEFT techniques such as LoRA to adapt the model to this task. ChatICD⁸³ and LLM-codex⁸⁴ both use ChatGPT with prompts for ICD coding, with LLM-codex⁸⁴ taking a step further by training a language model on top of the ChatGPT responses, demonstrating its strong coding performance in MIMIC-III data set²⁴.

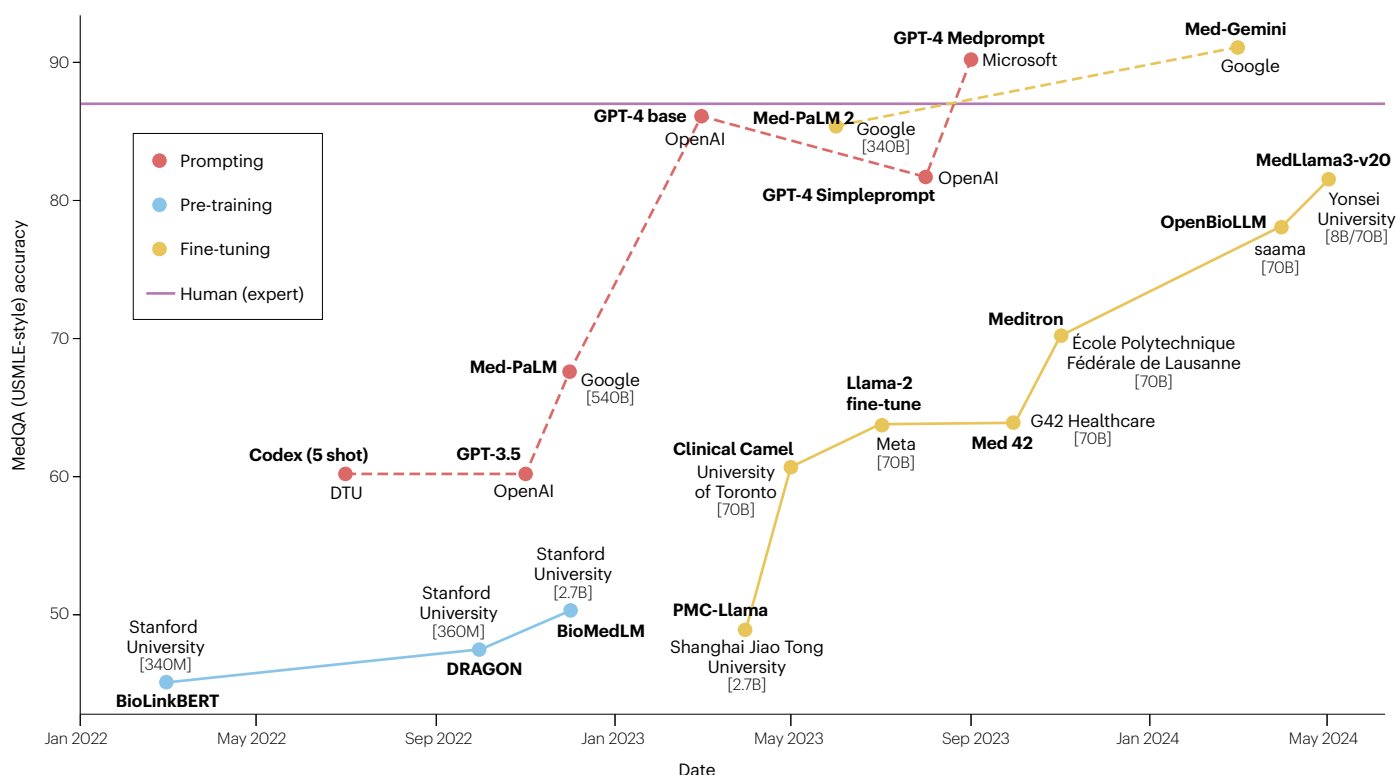


Fig. 3 | Development of medical LLMs over time. Illustration of the development of different medical large language models (LLMs) over time by assessing the scores of the United States Medical Licensing Examination (USMLE)

from the MedQA data set. Solid and dashed lines represent open-source and closed-source models, respectively. The number in square brackets indicates the number of model parameters. B, billion.

Table 2 | Summary of medical LLMs tailored to various clinical applications

Application	Model	Architecture	Model development	Number of parameters	Data scale	Data source	Evaluation (task: score)
Medical decision-making	Dr. Knows ⁷³	GPT-3.5	ICL	154B	5,820 notes	MIMIC-III ²⁴ + IN-HOUSE ⁷³	Diagnosis summarization: 30.72 ROUGE-L
	DDx PaLM-2 (ref. 75)	PaLM-2	FT and ICL	340B	–	MultiMedQA ¹⁰ + MIMIC-III ²⁴	Differential diagnosis: 0.591 top-10 accuracy
	NYUTron ⁷²	BERT	PT and FT	110M	7.25M notes, 4.1B tokens	NYU Notes ⁷²	Readmission prediction: 0.799 AUC
							In-hospital mortality pPrediction: 0.949 AUC
							Comorbidity index prediction: 0.894 AUC
							Length of stay prediction: 0.787 AUC
							Insurance denial prediction: 0.872 AUC
	Foresight ⁷⁶	GPT-2	PT and FT	1.5B	35M notes	King's College Hospital, MIMIC-III	Next biomedical concept forecast: 0.913 F1
						South London and Maudsley Hospital	–
	TrialGPT ⁷⁷	GPT-4	–	–	184 patients	2016 SIGIR ⁷⁷ , 2021 and 2022 TREC ⁷⁷	Ranking clinical trials: 0.733 Precision@10, 0.817 NDCG@10
Clinical coding							Excluding clinical trials: 0.775 AUROC
	PLM-ICD ⁸¹	RoBERTa	FT	355M	70,539 notes	MIMIC-II + MIMIC-III ²⁴	ICD code prediction: 0.926 AUC, 0.104 F1
	DRG-LLaMA ⁸²	LLaMA-7B	FT	7B	25k pairs	MIMIC-IV ¹⁷⁸	Diagnosis-related group prediction: 0.327 F1
	ChatICD ⁸³	ChatGPT	ICL	–	10k pairs	MIMIC-III ²⁴	ICD code prediction: 0.920 AUC, 0.681 F1
Clinical report generation	LLM-codex ⁸⁴	ChatGPT + LSTM	ICL	–	–	MIMIC-III ²⁴	ICD code prediction: 0.834 AUC, 0.468 F1
	ImpressionGPT ⁸⁹	ChatGPT	ICL and RAG	110M	184k reports	MIMIC-CXR ¹⁷⁸ + IU X-ray	Report summarization: 47.93 ROUGE-L
	RadAdapt ⁹⁰	T5	FT	223M, 738M	80k reports	MIMIC-III ²⁴	Report summarization: 36.8 ROUGE-L
	ChatCAD ⁸⁰	GPT-3	ICL	175B	300 reports	MIMIC-CXR ¹⁷⁸	Report generation: 0.605 F1
	MAIRA-1 (ref. 38)	ViT + Vicuna-7B	FT	8B	337k pairs	MIMIC-CXR ¹⁷⁸	Report generation: 28.9 ROUGE-L
Medical robotics	RadFM ³⁹	ViT + LLaMA-13B	PT and FT	14B	32M pairs	MedMD ³⁹	Report generation: 18.22 ROUGE-L
	SuFIA ⁹⁹	GPT-4	ICL	–	4 tasks	ORBIT-Surgical ⁹⁹	Surgical tasks: 100 success rate
	UltrasoundGPT ⁹⁷	GPT-4	ICL	–	522 tasks	–	Task completion: 80 success rate
Medical language translation	Robotic X-ray ¹⁰⁰	GPT-4	ICL	–	–	–	X-ray surgery: 7.6/10 human rating
	Medical mT5 (ref. 102)	T5	PT	738M, 3B	4.5B pairs	PubMed + EMEA ¹⁰²	(Multi-task) sequence labelling: 0.767 F1
						ClinicalTrials ¹⁰² , among others	Augment mining 0.733 F1
	Apollo ¹⁰³	Qwen	PT and FT	1.8B-7B	2.5B pairs	ApolloCorpora ¹⁰³	QA: 0.588 accuracy
	BiMedix ¹⁰⁴	Mistral	FT	13B	1.3M pairs	BiMed1.3M ¹⁰⁴	QA: 0.654 accuracy
	Biomed-sum ¹⁰⁵	BART	FT	406M	27k papers	BioCiteDB ¹⁰⁵	Abstractive summarization: 32.33 ROUGE-L
	RALL ¹⁰⁶	BART	FT and RAG	406M	63k pairs	CELLS ¹⁰⁵	Lay language generation: NA

Table 2 (continued) | Summary of medical LLMs tailored to various clinical applications

Application	Model	Architecture	Model development	Number of parameters	Data scale	Data source	Evaluation (task: score)
Medical education	ChatGPT ¹¹¹	GPT-3.5/GPT-4	ICL	–	–	–	Curriculum generation, learning planning
	Med-Gemini ³⁷	Gemini	FT and CoT	–	–	MedQA-R/RS ³⁷ + MultiMedQA ¹⁰	Text-based QA: 0.911 accuracy
						MIMIC-III ²⁴ + MultiMedBench ³⁷	Multimodal QA: 0.935 accuracy
Mental health support	PsyChat ¹²⁰	ChatGLM	FT	6B	350k pairs	Xingling ¹²⁰ + Smilechat ¹²⁰	Text generation: 27.6 ROUGE-L
	ChatCounselor ¹¹⁷	Vicuna	FT	7B	8k instructions	Psych8K ¹¹⁷	Question answering: evaluated by ChatGPT
	Mental-LLM ¹²²	Alpaca, FLAN-T5	FT and ICL	7B, 11B	31k pairs	Dreaddit + DepSeverity + SDCNL	Mental health prediction: 0.741 accuracy
						CSSRS-Suicide + Red-Sam	
Medical inquiry and response	AMIE ¹²⁶	PaLM2	FT	340B	>2M pairs	MedQA ¹³ + MultiMedBench ³⁷ + MIMIC-III ²⁴ + Dialogue ¹²⁶	Diagnostic Accuracy: 0.920 top-10 accuracy
	Healthcare Copilot ¹²⁵	ChatGPT	ICL	–	–	MedDialog ¹²⁶	Inquiry capability: 4.62/5 (ChatGPT)
							Conversational fluency: 4.06/5 (ChatGPT)
							Response accuracy: 4.56/5 (ChatGPT)
							Response safety: 3.88/5 (ChatGPT)
	Conversational Diagnosis ¹⁸¹	GPT-4/LLaMA	ICL	–	40k pairs	MIMIC-IV	Disease screening: 0.770 top-10 hit rate Differential diagnosis: 0.910 accuracy

AMIE, Articulate Medical Intelligence Explorer; AUC, area under the curve; AUROC, area under the receiver operating characteristic curve; B, billion; CoT, chain-of-thought prompting; FT, fine tuning; F1, harmonic mean of precision and recall; ICD, International Classification of Diseases; ICL, in-context learning; LLM, large language model; M, million; NA, not applicable; Precision@10, proportion of relevant items in the top 10 ranked/retrieved results; PT, pre-training; QA, question answering; RAG: retrieval-augmented generation.

ICD coding is typically formulated as a multi-label classification task using the MIMIC-III data set for training and evaluation. Models are assessed based on their F1 score, AUC and Precision@k (which measures the proportion of relevant items among the top k retrieved items), considering either the top k most frequent labels or the full label set. One challenge of deploying LLMs for clinical coding is the potential for biases and hallucinations; in particular, traditional multi-label classification models can easily constrain their outputs to a predefined list of (usually >1,000) candidate codes through a classification neural network. By contrast, generative LLMs could suffer from major hallucinations because the input text is lengthy. As a result, the LLM may assign a code that is not in the candidate list or that is a non-existent clinical code to the input text. It is therefore essential to establish a proactive mechanism to detect and correct errors before patient EHRs are entered. Most LLMs for clinical coding focus on ICD coding, but there is a growing need to expand to other types of clinical coding, such as medication and procedure coding, which are equally important to accurately capture patient information, facilitate billing and reimbursement processes, and support clinical decision-making.

Clinical report generation. Clinical reports, such as radiology reports⁸⁷, discharge summaries and patient clinic letters, refer to

standardized documentation that health-care workers complete after each patient visit⁸⁸. Clinical report generation usually involves text generation/summarization and information retrieval, a large portion of which often consists of medical diagnostic results. It is typically tedious for overworked clinicians to write clinical reports, and therefore they are often incomplete or error-prone. LLMs can act as an assistant tool to improve efficiency and reduce errors in lengthy reports^{89,90}. Another popular approach involves incorporating a vision-based model to provide complementary information^{38,39,80}; the vision model analyses the input medical image and generates an annotation, which serves as a supplementary input to the LLM, alongside additional text prompts.

General medical vision–language models such as Med-Gemini³⁷, LLaVA-Med³⁶ and Med-Flamingo³⁵ can serve as foundation models for broad medical domains, including radiology and pathology, with other models trained specifically on radiographs, such as ChatCAD⁸⁰, MAIRA-1 (ref. 38) and RadFM³⁹, with superior performance in specific subdomains.

LLMs can also leverage textual data for report summarization to generate radiology reports. This can be achieved using either unimodal LLMs, which input a long report and generate a summary, or multi-modal LLMs, which input both the long report and the corresponding image to generate a summary. Vision–language models can also be

Review article

developed for report summarization. For example, ImpressionGPT⁸⁹ is a unimodal LLM that uses dynamic prompts and iterative optimization to generate report summaries. RadAdapt⁹⁰ systematically evaluates

various language models and lightweight adaptation methods, achieving optimal report generation performance with a 36.8 ROUGE⁹¹ score, through pre-training on clinical text and parameter-efficient

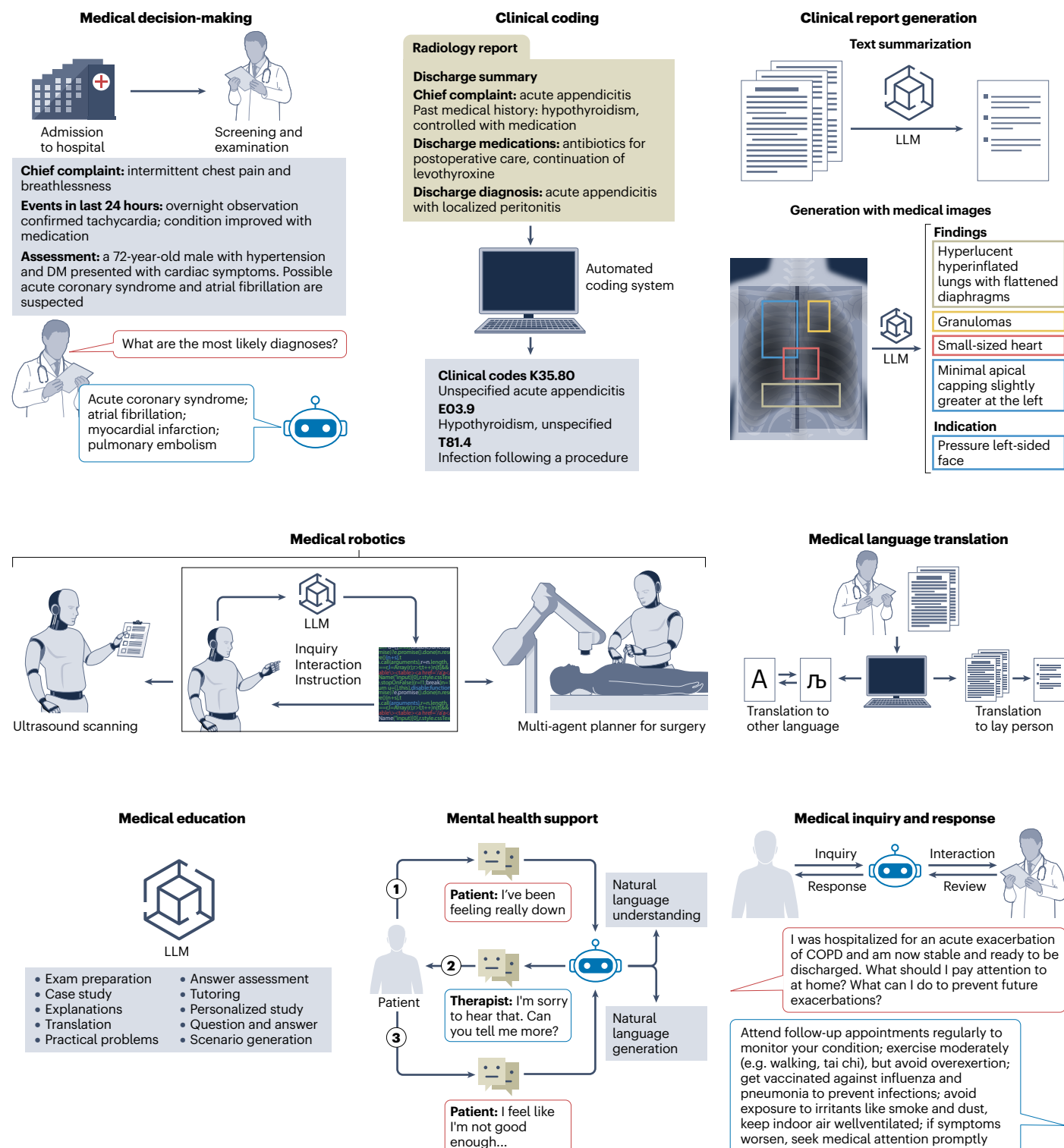


Fig. 4 | Application of LLMs in medicine. Integrated overview of applications^{73,109,115,182,183} of large language models (LLMs) in medicine. COPD, chronic obstructive pulmonary disease; DM, diabetes mellitus.

fine-tuning with LoRA, while also investigating the impact of few-shot prompting.

Evaluation of the performance of LLM-based radiology report generation models relies on the MIMIC-III or MIMIC-IV data sets, as they are the largest publicly available free-text EHRs. Common automatic evaluation metrics include BLEU⁹² and ROUGE⁹¹. Additionally, radiology-specific metrics such as RadCliQ⁹³ have been developed to better assess the quality and accuracy of the generated reports in the context of radiology. A clinical trial (NCT06263855)⁹⁴ in the USA is currently assessing whether using an LLM would improve the clarity and efficiency of discharge summaries.

LLM-generated reports tend to be less concise than human-written ones. Evaluation of these models is challenging because of the specialized nature of the content and the generative nature of the task. Current automatic evaluation methods focus on lexical metrics, which can lead to biased and inaccurate assessments of the contextual information present in the reports⁹⁵. For example, consider two sentences with similar meanings but different wordings: “The patient’s blood glucose level is within normal limits” and “The patient does not exhibit signs of hyperglycaemia”. Although both convey the absence of hyperglycaemia, lexical evaluation metrics might struggle to capture their semantic equivalence, as they rely on direct word-level comparisons. This discrepancy highlights the need for more sophisticated evaluation techniques that go beyond surface-level similarities, consider the underlying medical context, and can account for the nuances and variations in expressing clinical information.

Medical robotics

Medical robotics is revolutionizing health care by improving surgical procedures and medical imaging⁹⁶. LLMs can complement robotic technology by augmenting their decision-making, communication, interaction and control abilities. For example, surgical robots assisted with LLMs enable minimally invasive procedures with higher accuracy and shorter patient recovery times than traditional surgical robots without LLMs^{97,98}. Multi-agent planning systems designed with LLMs involve the coordination of multiple robotic units to perform collaborative tasks, improving surgical accuracy and operational efficiency⁹⁸. Similarly, SuFIA⁹⁹ combines the advanced reasoning capabilities of LLMs, specifically GPT-4 Turbo, with perception modules to implement high-level planning and low-level control of surgical robots for tasks such as instrument navigation and tissue manipulation.

In the field of medical imaging, UltrasoundGPT⁹⁷ equips ultrasound robots with LLMs and domain-specific knowledge by using an ultrasound operation knowledge database to enable precise motion planning. UltrasoundGPT utilizes a dynamic scanning strategy based on prompt engineering, which enables LLMs to adjust motion planning during procedures. This system demonstrates faster scan completion times and higher image quality than conventional ultrasound systems. Similarly, a simplified set of standardized commands and instructions enabled GPT-4 to control a robotic X-ray system named the Brainlab Loop-X¹⁰⁰ device.

The complexity of medical procedures, ethical considerations and patient safety concerns make it difficult to evaluate these systems in real health-care environments. Therefore, most current evaluations rely on simulated data and controlled laboratory settings; for example, SuFIA and Robotic X-ray’s performance are assessed using a combination of simulated surgical scenarios and expert human evaluation^{99,100}. Similarly, UltrasoundGPT is tested through the assessment of task completion⁹⁷. Moreover, the complex and dynamic nature of

shared human–robot workspaces might lead to LLM-powered medical robots misjudging human intentions or making inappropriate decisions, posing safety risks. Future research could explore safety features such as using sensing technologies and physical design constraints to minimize errors¹⁰¹.

Medical language translation

There are two main areas of medical language translation; the translation of medical terminology from one language to another^{102–104} and the translation of medical dialogue for ease of interpretation by non-professional personnel^{105,106}. Effective medical language translation is essential for providing high-quality health care to diverse patient populations.

Multilingual LLMs such as Medical mT5 (ref. 102), Apollo¹⁰³ and BiMediX¹⁰⁴, which are trained on extensive medical data sets in multiple languages, can be further fine-tuned to translate medical terminology between languages such as English, French, Spanish, Chinese and Arabic. When translating medical dialogue for non-professional understanding, it is crucial to fine-tune LLMs on data sets that encompass both technical medical conversations and their corresponding lay-language explanations. This training approach allows the models to learn the mapping between complex medical jargon and more accessible language. Techniques such as retrieval augmentation, which involves retrieving relevant lay-language explanations from external knowledge sources, can further enhance the quality and clarity of the translated dialogue^{105,106}. By integrating domain-specific knowledge from various sources, LLMs can generate more accurate and informative translations that cater to the needs of non-professional audiences.

Evaluating the performance of multilingual LLMs in medical language translation requires a multifaceted approach. Some models, such as Apollo¹⁰³ and BiMediX¹⁰⁴, use multiple-choice QA test data with the calculation of accuracy score^{103,104}. For generative benchmarks such as summarization^{105,106}, quantitative metrics like BLEU⁹² and ROUGE⁹¹ are commonly used to assess translation quality, but they should be supplemented with domain-specific evaluation criteria. For medical translations, accuracy of terminology, preservation of clinical meaning and consistency across languages are crucial factors. Human evaluation by bilingual medical experts is essential to validate the nuanced understanding of medical concepts across languages. For patient-oriented translations, comprehension tests with lay individuals can assess the effectiveness of jargon simplification.

In both translation and simplification tasks, misinterpretation is a common occurrence that can have damaging consequences. In developing and deploying medical translation and simplification platforms, developers should prioritize professional data sets, such as textbooks and peer-reviewed journals for medical knowledge recall. This way, it will be less likely for misinformation from unreliable sources to skew the output¹⁰⁷. Another ethical consideration of using LLMs to perform medical translation is the potential for discriminatory wording to be inserted inadvertently into the output. Such wording is difficult to prevent due to potential biases in training data and the multi-step processing of input through various model components¹⁰⁸.

Medical education

LLMs can be incorporated into the medical education system by facilitating study through explanations, language translation, answering questions, assisting with medical exam preparation and providing Socratic-style tutoring¹⁰⁹. Therefore, medical education could involve text generation, text simplification, semantic textual similarity and

information retrieval, among others. Medical education can be augmented by generating scenarios, problems and corresponding answers by an LLM. Moreover, students can gain a richer educational experience through personalized study modules and case-based assessments, including a wider array of challenges and scenarios beyond those found in standard textbooks¹⁰⁸. LLMs can also generate feedback on student responses to practical problems, informing them about their areas of weakness in real time¹¹⁰. LLMs can also be used to educate the public; medical dialogues are often complex and difficult to understand for the average patient. LLMs can tune the textual output of prompts to use varying degrees of medical terminology for different audiences¹⁰⁸.

Integrating LLMs into medical education can start with existing pre-trained models such as ChatGPT and Med-Gemini³⁷. For example, ChatGPT¹¹¹ can provide explanations and clarifications on complex medical concepts, whereas Med-Gemini³⁷, a multimodal model, can analyse medical images and generate detailed reports. Institutions, such as Second Xiangya Hospital of Central South University in China¹¹² and Carleton University in Canada¹¹³, are exploring the integration of these language models into curricula, leveraging their strengths while ensuring proper oversight and ethical considerations.

To evaluate the effectiveness of integrating LLMs into medical education, a combination of quantitative and qualitative methods should be used. Current research focuses on the QA-based evaluation³⁷; quantitative metrics can include student performance on assessments, such as exam scores and clinical skills evaluations, comparing outcomes before and after the introduction of LLM-based tools. Qualitative methods, such as surveys and focus groups, can gather feedback from students and educators on the perceived benefits, challenges and areas for improvement in using LLMs for learning and teaching. Additionally, longitudinal studies can track the long-term impact of LLM integration on student learning outcomes, clinical competence and career readiness.

Potential downsides of using LLMs in medical education include the current lack of ethical training and the presence of biases in training data sets²¹. These biases, if not addressed, can propagate through the generated outputs, reinforcing stereotypes and potentially leading to discrimination in medical education. The lack of explicit ethical training during LLM development may also result in the generation of content that does not align with the ethical principles and guidelines of the medical profession, such as promoting unethical practices or violating patient privacy. Moreover, LLMs can generate plausible-sounding but factually incorrect information, which can mislead students and health-care professionals if relied upon without proper verification. This can lead to the propagation of misconceptions, inappropriate treatment strategies or misdiagnosis¹¹⁴. To mitigate these risks, it is essential to establish rigorous fact-checking and validation processes and emphasize the importance of critical thinking, evidence-based practice and the verification of information from multiple reliable sources in medical education.

Mental health support

Mental health support involves both diagnosis and treatment; for example, depression is treated through a variety of psychotherapies, including cognitive behaviour therapy, interpersonal psychotherapy and psychodynamic therapy, among others¹¹⁵. Many of these techniques are based on patient–doctor conversations, with lengthy and expensive treatment plans. The ability of LLMs to serve as conversation partners could lower the barrier to entry for patients with financial or

physical constraints¹¹⁶, increasing the accessibility to mental health treatments¹¹⁷.

The willingness and level of self-disclosure has a strong influence on the effectiveness of mental health diagnosis and treatment, including with robots¹¹⁸. Instead of pre-training or fine-tuning on general medical data, it is often better to use medical QA data, because the LLM's main task will be talking to the patient, which involves back-and-forth conversation in the format of QA¹¹⁹. PsyChat¹²⁰ is a client-centric LLM dialogue system that provides psychological support comprising five modules: client behaviour recognition, counsellor strategy selection, input packer, response generator and response selection. Specifically, the response generator is fine-tuned with ChatGLM-6B on a vast dialogue data set, Xingling and SmileChat¹²⁰. The system demonstrated improved performance in metrics including empathy, relevance and therapeutic alignment compared with base LLM ChatGLM. Similarly, ChatCounselor initializes from Vicuna and fine-tunes from an 8k-sized instruct-tuning data set collected from real-world counselling dialogue examples¹¹⁷. Psy-LLM is meant to be an assistive mental health tool to support the workflow of professional counsellors, tailored for depression or anxiety cases¹¹⁹. A clinical trial (NCT06346496)¹²¹ in China is assessing the effectiveness of using LLMs for depression and anxiety symptoms in young adults over a 28-day period.

Fine-tuning on a variety of data sets can improve an LLM's capability on multiple mental-health-specific tasks across different data sets simultaneously¹²². For example, Mental-Alpaca and Mental-FLAN-T5 are instruction-fine-tuned on mental health data sets for tasks such as depression detection, stress prediction and suicide risk prediction¹²². Automated evaluations of mental health measure the relevance, coherence and empathy of the generated responses using metrics BLEU and accuracy. Mental health professionals conduct human evaluations through simulated counselling sessions, assessing the clinical appropriateness and therapeutic potential of the models' responses. Various evaluation frameworks have also been introduced that integrate text generation (conversational response)¹¹⁹, QA¹¹⁷ and mental health prediction¹²². For example, GPT-4 has been used as an evaluator to assess the proposed LLM ChatCounselor against traditional mental health chatbots¹¹⁷. The evaluation focuses on criteria including empathy, safety and therapeutic alignment, with ChatCounselor demonstrating superior performance in empathetic understanding and adherence to therapeutic principles.

The biggest risks in using LLMs for mental health support are the lack of emotional understanding and the inappropriate or harmful responses¹²³. LLMs might struggle to fully grasp and respond to the complex emotional states and needs of individuals seeking mental health support and might not be able to provide the same level of empathy and human connection required in therapeutic interactions. Moreover, if not properly trained or controlled, LLMs might generate responses that are inappropriate, insensitive or even harmful to individuals in vulnerable emotional states¹²⁴. They might provide advice that is not grounded in evidence-based psychological practices or that goes against established mental health guidelines. Addressing these challenges requires rigorous training of LLMs in evidence-based practices, ethical considerations, and risk assessment protocols, as well as collaboration between mental health professionals and AI (artificial intelligence) researchers.

Medical inquiry and response

LLMs are also suitable for tasks such as answering real-time patient inquiries and assisting physicians in documentation¹²⁵. Instead of

relying solely on rule-based algorithms or limited data sets, these systems leverage the vast knowledge and reasoning capabilities of LLMs to engage in diagnostic conversations and provide personalized recommendations. For example, Healthcare Copilot¹²⁵ combines dialogue management modules, patient history tracking mechanisms and information processing units to enable safe patient–LLM interactions, enhance conversations with historical data and summarize consultations. Similarly, Google’s Articulate Medical Intelligence Explorer (AMIE)¹²⁶ uses a self-play-based simulated environment with automated feedback mechanisms to enable the model to learn and adapt across different medical scenarios. Current evaluation often involves the calculation of metrics such as accuracy, precision, recall and F1-score¹²⁶. Multi-dimensional assessments, including inquiry capability, conversational fluency, response accuracy and safety using benchmarks, and comparisons with human experts or well-established models like ChatGPT, have also been conducted¹²⁵. However, these metrics alone are not sufficient, and evaluation should also focus on diagnostic accuracy, patient satisfaction and adherence to medical guidelines¹²⁷. The clinical trial ChiCTR2400081938 is currently assessing ChatGPT as an online consultant to assist physicians in remote diagnosis and treatment of hypertension in young adults¹²⁸; they focus on patient satisfaction, patient management efficiency, doctor work efficiency and quality of response information for evaluation.

Still, integrating medical LLMs into existing health-care workflows and infrastructure will require substantial technical and organizational efforts. Privacy and security concerns surrounding patient data must also be carefully considered and addressed. Ensuring transparency, explainability and accountability in the decision-making processes is crucial to maintaining trust and facilitating informed consent from patients¹²⁹.

Challenges

In this section, we address the challenges and discuss solutions to the adoption of LLMs in an array of medical applications.

Hallucination

Hallucination refers to situation in which the generated output contains inaccurate or nonfactual information. It can be categorized into intrinsic and extrinsic hallucinations¹¹⁴; intrinsic hallucinations generate outputs logically contradicting factual information, such as wrong mathematical calculations¹¹⁴. Extrinsic hallucinations occur when the generated output cannot be verified, such as ‘faking’ citations that do not exist or ‘dodging’ the question. When integrating LLMs into the medical domain, fluent but nonfactual LLM hallucinations can lead to the dissemination of incorrect medical information, causing misdiagnoses and inappropriate treatments.

Current solutions to mitigate LLM hallucination can be categorized into training-time correction, generation-time correction and retrieval-augmented correction. Training-time correction adjusts model parameter weights by including factually consistent reinforcement learning¹³⁰ and contrastive learning¹³¹. Generation-time correction adds a ‘reasoning’ process to the LLM inference to ensure reliability, using techniques such as sampling multiple outputs¹³² or a confidence score to identify hallucination before the final generation. Retrieval-augmented correction instead uses external resources to mitigate hallucination, such as using factual documents as prompts¹³³ or chain-of-retrieval prompting technique¹³⁴. For example, training-time correction is particularly suitable for specialized medical tasks like radiology reporting, in which consistent patterns exist in the training

data. Generation-time correction works well for general medical consultations in which multiple perspectives need to be considered. Retrieval-augmented correction is essential for tasks requiring up-to-date medical knowledge, such as treatment recommendations, in which external verification against current medical guidelines is crucial.

Lack of evaluation benchmarks and metrics

Current benchmarks and metrics often fail to evaluate LLMs’ overall capabilities, especially in the medical domain. For example, MedQA (USMLE)¹³ and MedMCQA⁶⁹ offer extensive coverage on QA tasks but fail to evaluate trustworthiness, helpfulness, explainability and faithfulness⁹⁵. Although HealthSearchQA provides some improvement by evaluating LLMs on common health queries that reflect real-world information needs, it still lacks comprehensive assessment of the aforementioned crucial aspects⁹. Benchmarks such as TruthfulQA¹³⁵ and HaluEval¹³⁶ evaluate metrics such as truthfulness but do not cover the medical domain. Future research is necessary in this space.

Domain data limitations

Current training data sets in the medical domain (Table 1) remain relatively small compared with those for general-purpose LLMs (Box 1). This results in medical-specified LLMs exhibiting extraordinary performance on open benchmarks with extensive data coverage yet falling short on real-life tasks such as differential diagnosis and personalized treatment planning¹⁰. Although the volume of medical and health data is large, most such data require extensive ethical, legal and privacy procedures to be accessed. Moreover, these data are often unlabelled, and solutions such as human labelling and unsupervised learning¹³⁷ are hindered by a lack of human expertise and small margins of error.

Current state-of-the-art approaches^{10,14} typically fine-tune the general LLMs on smaller open-sourced data sets to improve their domain-specific performance. Another solution is to generate high-quality synthetic data sets using LLMs to broaden the knowledge coverage; however, training on generated data sets causes models to experience catastrophic forgetting, in which they lose their original pretrained knowledge and capabilities due to the limited diversity and context in synthetic data¹³⁸.

New knowledge adaptation

Once trained, it is expensive and inefficient to inject new knowledge into an LLM through re-training. However, it is sometimes necessary to update on a new adverse effect of a medication or a new disease. Two problems arise during such knowledge updates; the first is how to make LLMs ‘forget’ the ‘old knowledge’, as it is almost impossible to remove it all from the training data, and the discrepancy between new and old knowledge can cause unintended association and bias¹³⁹. The second problem is the timeliness of the additional knowledge to ensure the model is updated in real time¹⁴⁰. Current solutions to knowledge adaptation can be categorized into model editing and RAG. Model editing¹⁴¹ alters the knowledge of the model by modifying its parameters; however, this method does not generalize well, meaning that its effectiveness is often limited to specific scenarios or model architectures, and it may not perform consistently across different tasks or domains. By contrast, RAG provides external knowledge sources as prompts during model inference; for example, by updating the model’s external knowledge memory¹⁴². Although RAG does not directly solve the ‘forget’ issue, it addresses the ‘timeliness’ problem by enabling quick updates of external knowledge without altering the model’s core parameters.

Behaviour alignment

Behaviour alignment refers to the process of ensuring that the LLMs' behaviours align with the objectives of its task, which is often to mimic general human behaviour. For example, ChatGPT demonstrates general conversational capabilities in answering human inquiries¹⁴³, but its answers to medical consultations are not as concise and professional as those of human experts^{9,10}.

Current solutions include instruction fine-tuning, reinforcement learning from human feedback¹⁴³ and prompt tuning^{43,50}. Instruction fine-tuning³² refers to improving the performance of LLMs on specific tasks based on explicit instructions¹⁴³ to generate better outputs. Reinforcement learning from human feedback uses human feedback to evaluate and align the LLMs' outputs, which could then be used as chatbots¹⁴⁴ and decision-making agents¹⁴⁵. Prompt tuning can also align LLMs to the expected output format; for example, chain of hindsight prompting enables the LLMs to review their initial responses, identify potential errors and generate corrected outputs¹⁴⁶.

Ethical and safety concerns

Concerns have been raised regarding the use of LLMs in the medical domain¹⁴⁷, with a focus on ethics, accountability and safety. For example, the scientific community has disapproved of using ChatGPT in writing biomedical research papers¹²⁹. Tracking the accountability of using LLMs as clinical assistants is also challenging^{33,148}; for example, prompt injection can cause the LLM to leak personally identifiable information (such as email addresses) from its training data¹⁴⁹. Such leakage has been attributed to the mismatched generalization between safety and capability objectives, that is, the pre-training of LLMs uses a larger and more varied data set than the data set used for safety training, resulting in many of the model's capabilities not being covered by safety training¹⁵⁰. A potential solution is to increase the safety training data set and develop comprehensive safety training to cover the model's behaviours and capabilities.

Regulatory challenges

The regulatory landscape of LLMs presents distinct challenges owing to their large scale, broad applicability and varying reliability across applications. As LLMs progressively permeate the fields of medicine and health care, their versatility allows a single LLM family to facilitate a multitude of tasks across a broad spectrum of interest groups. This is different from previous AI-based medical technologies, which were typically tailored to meet specific medical needs and cater to particular interest groups^{71,151}. This divergence requires regulators to develop fast and adaptable frameworks to ensure the safety, ethical standards and privacy of LLM-powered medical technologies without compromising innovation. For example, creating a dedicated regulatory category and incorporating patient-centred design principles in LLM development can help to ensure that decisions align with patient welfare and clinical best practices¹⁵¹. Other suggestions include assessing LLM-enabled applications in real-world settings, introducing obligations of transparency of data and algorithms, performing adaptive risk assessment and mitigation processes, and conducting continuous testing and refinement of audited technologies^{151–153}. Such proactive regulatory adaptations are crucial to maintaining high standards of safety, ethics and trustworthiness of medical technology.

Outlook

Although LLMs have already made an impact on people's lives through chatbots and search engines, their integration into medical practices is

still in its infancy. In particular, existing benchmarks fall short in evaluating LLMs for clinical applications¹⁵⁴; traditional benchmarks mainly gauge accuracy in medical QA and do not capture the full spectrum of clinical skills required⁹. Criticisms have been levelled against the use of human-centric standardized medical exams for LLM evaluation, arguing that passing these tests does not necessarily reflect an LLM's proficiency in the nuanced expertise required in real-world clinical settings⁹. Therefore, more comprehensive benchmarks should be developed to assess capabilities, such as sourcing from authoritative medical references, adapting to the evolving landscape of medical knowledge and communicating uncertainties clearly^{9,18}. New benchmarks should also incorporate scenarios that test an LLM's ability through simulation of real-world applications and adjust to feedback from clinicians while maintaining robustness. Moreover, these benchmarks should also assess parameters such as fairness, ethics and equity, which are currently evaluated through basic metrics like demographic parity but require more sophisticated measures incorporating contextual considerations⁹. For example, AMIE uses real physician evaluations and comprehensive criteria, including clinical reasoning, patient communication and professional behaviour as reflected in the objective structured clinical examination. However, these benchmarks are still not adaptive, scalable and robust enough for different and personalized applications. Future research could focus on using synthetic alongside real-world data, incorporating clinical guidelines such as patient safety protocols and cost-effectiveness, and developing interactive evaluation systems in which clinicians provide real-time feedback and assess model–physician collaboration.

Although LLMs mainly address NLP tasks, multimodal LLMs (MLLMs) or large multimodal models (text and visual data)¹⁵⁵ support a broader range of tasks, such as comprehending the underlying meaning of a meme and generating website codes from images. Several MLLM-based frameworks integrating vision and language, such as Med-Flamingo³⁵, LLaVA-Med³⁶ and Med-Gemini³⁷, adopt medical image–text pairs for fine-tuning, thus enabling the medical LLMs to understand medical images (for example, radiology). For example, integrating vision, audio and language inputs for automated dental diagnosis has shown promising results for clinical assessment¹⁵⁶. However, only very few medical LLMs^{157,158} can process time series data such as electrocardiograms¹⁵⁷ and sphygmomanometres¹⁵⁸, despite their importance in diagnosis and monitoring. MLLMs trained at scale could potentially generalize across various domains and modalities outside of NLP tasks. However, the training of MLLMs at scale faces challenges in aligning and processing multiple modalities simultaneously, leading to computational constraints that result in smaller model sizes than single-modality LLMs. Future research could focus on improving processing, representation and learning of multi-modal data and knowledge, and cost-effective training of MLLMs, especially more resource-demanding modalities such as videos and images, and the collection and access to multi-modal clinical data.

Another promising line of research are LLM-based agents¹⁴⁵, which can be seen as autonomous systems that combine LLMs' reasoning capabilities with the ability to interact with external tools and environments to achieve specific goals. These agents use LLMs as controllers to leverage their reasoning capabilities. By integrating LLMs with external tools and multimodal perceptions, these agents can interact with environments, learn from feedback and acquire new skills to solve complex tasks (for example, software design or molecular dynamics simulation) through human-like behaviours, such as role-playing and communication^{123,145}. For example, Chat-Orthopedist⁶² interacts with

external knowledge bases, such as UpToDate.com, acquiring up-to-date adolescent idiopathic scoliosis domain knowledge to provide accurate and comprehensible responses to patient inquiries. However, integrating these agents in the medical domain is challenging, as the medical field involves numerous roles^{123,145} and decision-making processes. For example, disease diagnosis often requires series of tests such as computed tomography scans, ultrasounds, electrocardiograms and blood tests. LLMs could be used to model each of these roles/expertise and create collaborative medical agents to provide a more holistic and accurate diagnosis. Not only do these agents interpret individual medical reports, but they can also integrate these interpretations to form a cohesive medical opinion. Future research in this space could explore seamless data pipelines that collect data from various devices and transform them into data formats that are compatible with LLMs; improve communication and collaboration between agents, especially in areas such as ensuring truthfulness during communication, dispute resolution between agents and role-based data security measures; carry out real-time decision-making using data collected from remote monitoring devices; and utilize adaptive learning to prepare for unforeseen medical and health-care situations. Finally, current medical LLM research has largely focused on general medicine, likely due to the greater availability of data in this area^{10,148}. This has resulted in the under-representation of specialized fields such as ‘rehabilitation therapy’ or ‘sports medicine’.

So far, the medical community has primarily adopted LLMs provided by companies without questioning their data training, ethical protocols or privacy protection. Medical professionals are therefore encouraged to actively participate in creating and deploying medical LLMs by providing relevant training data, defining the desired benefits of LLMs and conducting tests in real-world scenarios to evaluate these benefits^{18,80}. Such assessments would help to determine the legal and medical risks associated with LLM use in medicine and inform strategies to mitigate LLM hallucination¹⁵⁹. Moreover, training ‘bilingual’ professionals – those versed in both medicine and LLM technology – will be increasingly important. Future research may explore interdisciplinary frameworks to facilitate the sharing of localized data from rural clinics; the creation of ‘bilingual education programmes’ that offer training from both worlds, as demonstrated in emerging medical AI curricula¹⁶⁰; and the development of institutional data management protocols and privacy protection mechanisms to help hospitals and physicians ‘guard’ patient data from corporations, without stifling innovation.

Published online: 7 April 2025

References

- Zhao, W. X. et al. A survey of large language models. Preprint at arXiv <https://doi.org/10.48550/arXiv.2303.18223> (2023).
- Yang, J. et al. Harnessing the power of LLMs in practice: a survey on ChatGPT and beyond. *ACM Trans. Knowl. Discov. Data* **18**, 160 (2024).
- Chowdhery, A. et al. PaLM: scaling language modeling with pathways. Preprint at arXiv <https://doi.org/10.48550/arXiv.2204.02311> (2022).
- Touvron, H. et al. LLaMA: open and efficient foundation language models. Preprint at arXiv <https://doi.org/10.48550/arXiv.2302.13971> (2023).
- Touvron, H. et al. Llama 2: open foundation and fine-tuned chat models. Preprint at arXiv <https://doi.org/10.48550/arXiv.2307.09288> (2023).
- Brown, T. et al. Language models are few-shot learners. In *Proc. 34th Int. Conf. Neural Inform. Process. Syst.* (eds Larochelle, H. et al.) 1877–1901 (2020).
- OpenAI et al. GPT-4 technical report. Preprint at arXiv <https://doi.org/10.48550/arXiv.2303.08774> (2023).
- Du, Z. et al. GLM: general language model pretraining with autoregressive blank infilling. In *Proc. 60th Annu. Meet. Assoc. Comput. Linguist.* (eds Muresan, S. et al.) 320–335 (ACL, 2022).
- Singhal, K. et al. Large language models encode clinical knowledge. *Nature* **620**, 172–180 (2023).
- Singhal, K. et al. Toward expert-level medical question answering with large language models. *Nat. Med.* <https://doi.org/10.1038/s41591-024-03423-7> (2025).
- Nori, H. et al. Can generalist foundation models outcompete special-purpose tuning? Case study in medicine. Preprint at arXiv <https://doi.org/10.48550/arXiv.2311.16452> (2023).
- Wu, C. et al. PMC-LLaMA: toward building open-source language models for medicine. *J. Am. Med. Inform. Assoc.* **31**, 1833–1843 (2024).
- Jin, D. et al. What disease does this patient have? A large-scale open domain question answering dataset from medical exams. *Appl. Sci.* **11**, 6421 (2021).
- Li, Y. et al. ChatDoctor: a medical chat model fine-tuned on a large language model meta-AI (LLaMA) using medical domain knowledge. *Cureus* **15**, 6 (2023).
- Han, T. et al. MedAlpaca — an open-source collection of medical conversational AI models and training data. Preprint at arXiv <https://doi.org/10.48550/arXiv.2304.08247> (2023).
- Wang, H. et al. HuaTuo: tuning LLaMA model with Chinese medical knowledge. Preprint at arXiv <https://doi.org/10.48550/arXiv.2304.06975> (2023).
- Toma, A. et al. Clinical Camel: an open-source expert-level medical language model with dialogue-based knowledge encoding. Preprint at arXiv <https://doi.org/10.48550/arXiv.2305.12031> (2023).
- Thirunavukarasu, A. J. et al. Large language models in medicine. *Nat. Med.* **29**, 1930–1940 (2023).
- Patel, S. B. & Lam, K. ChatGPT: the future of discharge summaries? *Lancet Digit. Health* **5**, e107–e108 (2023).
- Yang, X. et al. A large language model for electronic health records. *npj Digit. Med.* **5**, 194 (2022).
- Abd-Alrazaq, A. et al. Large language models in medical education: opportunities, challenges, and future directions. *JMIR Med. Educ.* **9**, e48291 (2023).
- Peng, C. et al. A study of generative large language model for medical research and healthcare. *npj Digit. Med.* **6**, 210 (2023).
- Alsentzer, E. et al. Publicly available clinical BERT embeddings. Preprint at arXiv <https://doi.org/10.48550/arXiv.1904.03323> (2019).
- Johnson, A. E. W. et al. MIMIC-III, a freely accessible critical care database. *Sci. Data* **3**, 160035 (2016).
- Wu, J. et al. Clinical text datasets for medical artificial intelligence and large language models — a systematic review. *NEJM AI* **1**, A2400012 (2024).
- Lee, J. et al. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics* **36**, 1234–1240 (2020).
- Gu, Y. et al. Domain-specific language model pretraining for biomedical natural language processing. *ACM Trans. Comput. Healthc.* **3**, 1–23 (2021).
- Luo, R. et al. BioGPT: generative pre-trained transformer for biomedical text generation and mining. *Brief. Bioinform.* **23**, bbac409 (2022).
- Ye, Q. et al. Qilin-Med: multi-stage knowledge injection advanced medical large language model. Preprint at arXiv <https://doi.org/10.48550/arXiv.2310.09089> (2023).
- Xiong, H. et al. DoctorGLM: fine-tuning your Chinese doctor is not a herculean task. Preprint at arXiv <https://doi.org/10.48550/arXiv.2304.01097> (2023).
- Yang, S. et al. Zhongjing: enhancing the Chinese medical capabilities of large language model through expert feedback and real-world multi-turn dialogue. In *Proc. AAAI Conf. Artif. Intell.* (eds Wooldridge, M. J. et al.) 19368–19376 (AAAI, 2023).
- Zhang, S. et al. Instruction tuning for large language models: a survey. Preprint at arXiv <https://doi.org/10.48550/arXiv.2308.10792> (2023).
- He, K. et al. A survey of large language models for healthcare: from data, technology, and applications to accountability and ethics. *Inf. Fusion* **118**, 102963 (2025).
- Byambasuren, O. et al. Preliminary study on the construction of Chinese medical knowledge graph. *J. Chin. Inf. Process.* **33**, 1–9 (2019).
- Moor, M. et al. Med-Flamingo: a multimodal medical few-shot learner. In *Proc. 3rd Mach. Learn. Health Symp.* (eds Hegselmann, S. et al.) 353–367 (PMLR, 2023).
- Li, C. et al. LLaVA-Med: training a large language-and-vision assistant for biomedicine in one day. In *Ann. Conf. Neural Inform. Process. Syst.* (eds Oh, A. et al.) 28541–28564 (Curran Associates, 2023).
- Saab, K. et al. Capabilities of Gemini models in medicine. Preprint at arXiv <https://doi.org/10.48550/arXiv.2404.18416> (2024).
- Hyland, S. L. et al. MAIRA-1: a specialised large multimodal model for radiology report generation. Preprint at arXiv <https://doi.org/10.48550/arXiv.2311.13668> (2023).
- Wu, C., Zhang, X., Zhang, Y., Wang, Y. & Xie, W. Towards generalist foundation model for radiology. Preprint at arXiv <https://doi.org/10.48550/arXiv.2308.02463> (2023).
- Zhang, X. et al. AlpaCare: instruction-tuned large language models for medical application. Preprint at arXiv <https://doi.org/10.48550/arXiv.2310.14558> (2023).
- Hu, E. J. et al. LoRA: low-rank adaptation of large language models. Preprint at arXiv <https://doi.org/10.48550/arXiv.2106.09685> (2021).
- Li, X. L. & Liang, P. Prefix-tuning: optimizing continuous prompts for generation. In *Proc. 59th Annu. Meet. Assoc. Comput. Linguist.* (eds Zong, C. et al.) 4582–4597 (ACL, 2021).
- Liu, X. et al. P-Tuning: prompt tuning can be comparable to fine-tuning across scales and tasks. In *Proc. 60th Annu. Meet. Assoc. Comput. Linguist.* (eds Muresan, S. et al.) 61–68 (ACL, 2022).
- Houlsby, N. et al. Parameter-efficient transfer learning for NLP. In *Proc. 36th Int. Conf. Mach. Learn.* (eds Chaudhuri, K. & Salakhutdinov, R.) 2790–2799 (PMLR, 2019).
- Xu, C., Guo, D., Duan, N. & McAuley, J. Baize: an open-source chat model with parameter-efficient tuning on self-chat data. In *Proc. 2023 Conf. Empir. Methods Nat. Lang. Process.* (eds Bouamor, H. et al.) 6268–6278 (ACL, 2023).

46. Shoham, O. B. & Rappoport, N. CLLM: clinical prediction with large language models. *PLoS Digit. Health* **3**, e0000680 (2024).
47. Dong, Q. et al. A survey on in-context learning. In *Proc. Conf. Empir. Methods Nat. Lang. Process.* (eds Al-Onaizan, Y. et al.) 1107–1128 (ACL, 2024).
48. Wei, J. et al. Chain-of-thought prompting elicits reasoning in large language models. In *Proc. 36th Int. Conf. Neural Inform. Process. Syst.* (eds Koyejo, S. et al.) 24824–24837 (Curran Associates, 2022).
49. Liu, Z. et al. DelD-GPT: zero-shot medical text de-identification by GPT-4. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2303.11032> (2023).
50. Lester, B., Al-Rfou, R. & Constant, N. The power of scale for parameter-efficient prompt tuning. In *Proc. Conf. Empir. Methods Nat. Lang. Process.* (eds Moens, M.-F. et al.) 3045–3059 (ACL, 2021).
51. Gao, Y. et al. Retrieval-augmented generation for large language models: a survey. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2312.10997> (2023).
52. Luo, Y. et al. An empirical study of catastrophic forgetting in large language models during continual fine-tuning. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2308.08747> (2023).
53. Xiong, G., Jin, Q., Lu, Z. & Zhang, A. Benchmarking retrieval-augmented generation for medicine. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2402.13178> (2024).
54. Li, X. & Li, J. AngLE-optimized text embeddings. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2309.12871> (2023).
55. Wang, G. et al. Voyager: an open-ended embodied agent with large language models. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2305.16291> (2023).
56. Chen, J. et al. M3-Embedding: multi-linguality, multi-functionality, multi-granularity text embeddings through self-knowledge distillation. In *Findings Assoc. Comput. Linguist.* (eds Ku, L. et al.) 2318–2335 (ACL, 2024).
57. Shao, Z. et al. Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy. In *Findings Assoc. Comput. Linguist.* (eds Bouamor, H. et al.) 9248–9274 (ACL, 2023).
58. Trivedi, H., Balasubramanian, N., Khot, T. & Sabharwal, A. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. In *Proc. 61st Annu. Meet. Assoc. Comput. Linguist.* (eds Rogers, A. et al.) 10014–10037 (ACL, 2023).
59. Asai, A., Wu, Z., Wang, Y., Sil, A. & Hajishirzi, H. Self-rag: learning to retrieve, generate, and critique through self-reflection. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2310.11511> (2023).
60. Zakka, C. et al. Almanac—retrieval-augmented language models for clinical medicine. *NEJM AI* **1**, A0a2300068 (2024).
61. Kim, J. & Min, M. From RAG to QA-RAG: integrating generative AI for pharmaceutical regulatory compliance process. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2402.01717> (2024).
62. Shi, W. et al. Retrieval-augmented large language models for adolescent idiopathic scoliosis patients in shared decision-making. In *Proc. 14th ACM Int. Conf. Bioinform. Comput. Biol. Health Inform.* (ACM, 2023).
63. Tang, L. et al. Evaluating large language models on medical evidence summarization. *npj Digit. Med.* **6**, 158 (2023).
64. Van Veen, D. et al. Adapted large language models can outperform medical experts in clinical text summarization. *Nat. Med.* **30**, 1134–1142 (2024).
65. Ondov, B., Attal, K. & Demner-Fushman, D. A survey of automated methods for biomedical text simplification. *J. Am. Med. Inform. Assoc.* **29**, 1976–1988 (2022).
66. Liu, F. et al. Retrieve, reason, and refine: generating accurate and faithful patient instructions. In *Proc. 36th Int. Conf. Neural Inform. Process. Syst.* (eds Koyejo, S. et al.) 18864–18877 (Curran Associates, 2022).
67. Joseph, S. et al. Multilingual simplification of medical texts. In *Proc. Conf. Empir. Methods Nat. Lang. Process.* (eds Bouamor, H. et al.) 16662–16692 (ACL, 2023).
68. Jin, Q., Dhingra, B., Liu, Z., Cohen, W. W. & Lu, X. PubMedQA: a dataset for biomedical research question answering. In *Proc. Conf. Empir. Methods Nat. Lang. Process. & 9th Int. Joint Conf. Nat. Lang. Process.* (eds Inui, K. et al.) 2567–2577 (ACL, 2019).
69. Pal, A., Umapathi, L. K. & Sankarasubbu, M. MedMCQA: a large-scale multi-subject multi-choice dataset for medical domain question answering. In *Proc. Conf. Health Inference Learn.* (eds Flores, G. et al.) 248–260 (PMLR, 2022).
70. Omar, M., Nadkarni, G. N., Klang, E. & Glucksberg, B. S. Large language models in medicine: a review of current clinical trials across healthcare applications. *PLoS Digit. Health* **3**, e0000662 (2024).
71. Rajpurkar, P., Chen, E., Banerjee, O. & Topol, E. J. AI in health and medicine. *Nat. Med.* **28**, 31–38 (2022).
72. Jiang, L. Y. et al. Health system-scale language models are all-purpose prediction engines. *Nature* **619**, 357–362 (2023).
73. Gao, Y. et al. Leveraging a medical knowledge graph into large language models for diagnosis prediction. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2308.14321> (2023).
74. Chung, H. W. et al. Scaling instruction-finetuned language models. *J. Mach. Learn. Res.* **25**, 1–53 (2024).
75. McDuff, D. et al. Towards accurate differential diagnosis with large language models. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2312.00164> (2023).
76. Kraljevic, Z. et al. Foresight—a generative pretrained transformer for modelling of patient timelines using electronic health records: a retrospective modelling study. *Lancet Digit. Health* **6**, e281–e290 (2024).
77. Jin, Q. et al. Matching patients to clinical trials with large language models. *Na. Commun.* **15**, 9074 (2024).
78. US National Library of Medicine. *ClinicalTrials.gov* <https://clinicaltrials.gov/study/NCT06002425> (2024).
79. German Clinical Trials Register. *drks.de* <https://drks.de/search/en/trial/DRKS00033775> (2024).
80. Wang, S., Zhao, Z., Ouyang, X., Wang, Q. & Shen, D. Interactive computer-aided diagnosis on medical image using large language models. *Commun. Eng.* **3**, 133 (2024).
81. Huang, C.-W., Tsai, S.-C. & Chen, Y.-N. PLM-ICD: automatic ICD coding with pretrained language models. In *Proc. 4th Clin. Nat. Lang. Process. Workshop* (eds Naumann, T. et al.) 10–20 (ACL, 2022).
82. Wang, H., Gao, C., Dantona, C., Hull, B. & Sun, J. DRG-LLaMA: tuning LLaMA model to predict diagnosis-related group for hospitalized patients. *npj Digit. Med.* **7**, 16 (2024).
83. Liu, J., Yang, S., Peng, T., Hu, X. & Zhu, Q. ChatICD: prompt learning for few-shot ICD coding through ChatGPT. In *2023 IEEE Int. Conf. Bioinform. Biomed.* (eds Jiang, X. et al.) 4360–4367 (IEEE, 2023).
84. Yang, Z., Batra, S. S., Stremmel, J. & Halperin, E. Surpassing GPT-4 medical coding with a two-stage approach. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2311.13735> (2023).
85. Elkin, P. L. & Brown, S. H. in *Terminology, Ontology and their Implementations* 2nd edn (ed. Elkin, P. L.) 367–370 (Springer, 2023).
86. Liu, Y. et al. RoBERTa: a robustly optimized BERT pretraining approach. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.1907.11692> (2019).
87. Liu, F., Wu, X., Ge, S., Fan, W. & Zou, Y. Exploring and distilling posterior and prior knowledge for radiology report generation. In *2021 IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)* 13748–13757 (IEEE, 2021).
88. Liu, X. et al. Reporting guidelines for clinical trials evaluating artificial intelligence interventions are needed. *Nat. Med.* **25**, 1467–1469 (2019).
89. Ma, C. et al. An iterative optimizing framework for radiology report summarization with ChatGPT. In *IEEE Trans. Artif. Intell.* 4163–4175 (IEEE, 2024).
90. Van Veen, D. et al. RadAdapt: radiology report summarization via lightweight domain adaptation of large language models. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2305.01146> (2023).
91. Lin, C.-Y. ROUGE: a package for automatic evaluation of summaries. In *Proc. Workshop Text Summarization Branches Out* 74–81 (ACL, 2004).
92. Papineni, K., Roukos, S., Ward, T. & Zhu, W. BLEU: a method for automatic evaluation of machine translation. In *Proc. 40th Annu. Meet. Assoc. Comput. Linguist.* (eds Isabelle, P. et al.) 311–318 (ACL, 2002).
93. Yu, F. et al. Evaluating progress in automatic chest X-ray radiology report generation. *Patterns* **4**, 100802 (2023).
94. US National Library of Medicine. *ClinicalTrials.gov* <https://clinicaltrials.gov/study/NCT06263855> (2024).
95. Xie, Q. et al. Faithful AI in medicine: a systematic review with large language models and beyond. Preprint at *medRxiv* <https://doi.org/10.1101/2023.04.18.23288752> (2023).
96. Dupont, P. E. A decade retrospective of medical robotics research from 2010 to 2020. *Sci. Robot.* **6**, eabi8017 (2021).
97. Xu, H. et al. Enhancing surgical robots with embodied intelligence for autonomous ultrasound scanning. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2405.00461> (2024).
98. Wang, J. et al. Large language models for robotics: opportunities, challenges, and perspectives. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2401.04334> (2024).
99. Moghani, M. et al. SuFIA: language-guided augmented dexterity for robotic surgical assistants. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2405.05226> (2024).
100. Killeen, B. D., Chaudhary, S., Osgood, G. & Unberath, M. Take a shot! natural language control of intelligent robotic X-ray systems in surgery. *Int. J. Comput. Assist. Radiol. Surg.* **19**, 1165–1173 (2024).
101. Weeraratna, I. N., Raymond, D. & Luharia, A. Human-robot collaboration for healthcare: a narrative review. *Cureus* **15**, e49210 (2023).
102. Garcia-Ferrero, I. et al. Medical mT5: an open-source multilingual text-to-text LLM for the medical domain. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2404.07613> (2024).
103. Wang, X. et al. Apollo: lightweight multilingual medical LLM towards democratizing medical AI to 6b people. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2403.03640> (2024).
104. Pieri, S. et al. BiMedix: bilingual medical mixture of experts LLM. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2402.13253> (2024).
105. Tang, C., Wang, S., Goldsack, T. & Lin, C. Improving biomedical abstractive summarisation with knowledge aggregation from citation papers. In *Proc. Conf. Empir. Methods Nat. Lang. Process.* (eds Bouamor, H. et al.) 606–618 (ACL, 2023).
106. Guo, Y., Qiu, W., Leroy, G., Wang, S. & Cohen, T. Retrieval augmentation of large language models for lay language generation. *J. Biomed. Inform.* **149**, 104580 (2024).
107. Chen, Y., Arunasalam, A. & Celik, Z. B. Can large language models provide security & privacy advice? Measuring the ability of LLMs to refute misconceptions. In *Proc. 39th Annu. Comput. Secur. Appl. Conf.* 366–378 (ACL, 2023).
108. Karabacak, M. et al. The advent of generative language models in medical education. *JMIR Med. Educ.* **9**, e48163 (2023).
109. Biri, S. K. et al. Assessing the utilization of large language models in medical education: insights from undergraduate medical students. *Cureus* **15**, e47468 (2023).
110. Ahn, S. The impending impacts of large language models on medical education. *Korean J. Med. Educ.* **35**, 103–107 (2023).

111. Peacock, J., Austin, A., Shapiro, M., Battista, A. & Samuel, A. Accelerating medical education with ChatGPT: an implementation guide. *MedEdPublish* **13**, 64 (2023).
112. Tian, Q. et al. Iteratively refined ChatGPT outperforms clinical mentors in generating high-quality interprofessional education clinical scenarios: a comparative study. Preprint at *Res. Sq.* <https://doi.org/10.21203/rs.3.rs-4637356/v1> (2024).
113. Veras, M. et al. Usability and efficacy of artificial intelligence chatbots (ChatGPT) for health sciences students: protocol for a crossover randomized controlled trial. *JMIR Res. Protoc.* **12**, e51873 (2023).
114. Rawte, V., Sheth, A. & Das, A. A survey of hallucination in large foundation models. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2309.05922> (2023).
115. Vaidyam, A. N., Wisniewski, H., Halamka, J. D., Kashavan, M. S. & Torous, J. B. Chatbots and conversational agents in mental health: a review of the psychiatric landscape. *Can. J. Psychiatry* **64**, 456–464 (2019).
116. Stock, A., Schlögl, S. & Groth, A. Tell me, what are you most afraid of? Exploring the effects of agent representation on information disclosure in human–chatbot interaction. In *Proc. Int. Conf. Hum. Comput. Interact.* (eds Degen, H. et al.) 179–191 (Springer, 2023). https://doi.org/10.1007/978-3-031-35894-4_13.
117. Liu, J. M. et al. ChatCounselor: a large language models for mental health support. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2309.15461> (2023).
118. Robinson, N., Connolly, J., Suddrey, G. & Kavanagh, D. J. A brief wellbeing training session delivered by a humanoid social robot: a pilot randomized controlled trial. *Int. J. Soc. Robot.* **16**, 937–951 (2024).
119. Lai, T. et al. Psy-LLM: scaling up global mental health psychological services with AI-based large language models. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2307.11991> (2023).
120. Qiu, H., Li, A., Ma, L. & Lan, Z. PsyChat: a client-centric dialogue system for mental health support. In *Proc. 27th Int. Conf. Comput. Support. Coop. Work Des. (CSCWD)* 2979–2984 (IEEE, 2024).
121. US National Library of Medicine. *ClinicalTrials.gov* <https://clinicaltrials.gov/study/NCT06346496> (2024).
122. Xu, X. et al. Mental-LLM: leveraging large language models for mental health prediction via online text data. In *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* Vol. 8 1–32 (ACM, 2024).
123. Ma, Z., Mei, Y. & Su, Z. Understanding the benefits and challenges of using large language model-based conversational agents for mental well-being support. *AMIA Annu. Symp. Proc.* **2023**, 1105 (2023).
124. Chung, N. C., Dyer, G. & Brocki, L. Challenges of large language models for mental health counseling. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2311.13857> (2023).
125. Ren, Z., Zhan, Y., Yu, B., Ding, L. & Tao, D. Healthcare Copilot: eliciting the power of general LLMs for medical consultation. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2402.13408> (2024).
126. Tu, T. et al. Towards conversational diagnostic AI. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2401.05654> (2024).
127. Hager, P. et al. Evaluation and mitigation of the limitations of large language models in clinical decision-making. *Nat. Med.* **30**, 2613–2622 (2024).
128. Chinese Clinical Trial Register. *ChiCTR.org.cn* <https://www.chictr.org.cn/showproj.html?proj=220887> (2024).
129. Stokel-Walker, C. ChatGPT listed as author on research papers: many scientists disapprove. *Nature* **613**, 620–621 (2023).
130. Roit, P. et al. Factually consistent summarization via reinforcement learning with textual entailment feedback. In *Proc. 61st Annu. Meet. Assoc. Comput. Linguist.* (eds Rogers, A. et al.) 6252–6272 (ACL, 2023).
131. Chern, I.-C. et al. Improving factuality of abstractive summarization via contrastive reward learning. In *Proc. 3rd Workshop Trustworthy Nat. Lang. Process.* (eds Ovalle, A. et al.) 55–60 (ACL, 2023).
132. Manakul, P., Liusie, A. & Gales, M. J. SelfCheckGPT: zero-resource black-box hallucination detection for generative large language models. In *Proc. 2023 Conf. Empir. Methods Nat. Lang. Process.* (eds Bouamor, H. et al.) 9004–9017 (ACL, 2023).
133. Shuster, K., Poff, S., Chen, M., Kiela, D. & Weston, J. Retrieval augmentation reduces hallucination in conversation. In *Find. Assoc. Comput. Linguist.* (eds Moens, M. et al.) 3784–3803 (ACL, 2021).
134. Dhuliawala, S. et al. Chain-of-verification reduces hallucination in large language models. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2309.11495> (2023).
135. Lin, S., Hilton, J. & Evans, O. TruthfulQA: measuring how models mimic human falsehoods. In *Proc. 60th Annu. Meet. Assoc. Comput. Linguist.* (eds Muresan, S. et al.) 3214–3252 (ACL, 2022).
136. Li, J., Cheng, X., Zhao, W. X., Nie, J.-Y. & Wen, J.-R. HaluEval: a large-scale hallucination evaluation benchmark for large language models. In *Proc. Conf. Empir. Methods Nat. Lang. Process.* (eds Bouamor, H. et al.) 6449–6464 (ACL, 2023).
137. Liu, F. et al. Auto-encoding knowledge graph for unsupervised medical report generation. In *Proc. 35th Int. Conf. Neural Inform. Process. Syst.* (eds Ranzato, M. et al.) 16266–16279 (Curran Associates, 2021).
138. Shumailov, I. et al. Model dementia: generated data makes models forget. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2305.17493> (2023).
139. Hoelscher-Obermaier, J., Persson, J., Kran, E., Konstas, I. & Barez, F. Detecting edit failures in large language models: an improved specificity benchmark. In *Find. Assoc. Comput. Linguist.* (eds Rogers, A. et al.) 11548–11559 (ACL, 2023).
140. Liu, F. et al. A medical multimodal large language model for future pandemics. *npj Digit. Med.* **6**, 226 (2023).
141. Yao, Y. et al. Editing large language models: problems, methods, and opportunities. In *Proc. Conf. Empir. Methods Nat. Lang. Process.* (eds Bouamor, H. et al.) 10222–10240 (ACL, 2023).
142. Lewis, P. et al. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Adv. Neural Inform. Process. Syst.* (eds Larochelle, H. et al.) 9459–9474 (Curran Associates, 2020).
143. Ouyang, L. et al. Training language models to follow instructions with human feedback. In *Proc. 36th Int. Conf. Neural Inf. Process. Syst.* (eds Koyejo, S. et al.) 27730–27744 (Curran Associates, 2022).
144. Glaese, A. et al. Improving alignment of dialogue agents via targeted human judgements. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2209.14375> (2022).
145. Xi, Z. et al. The rise and potential of large language model based agents: a survey. *Sci. China Inf. Sci.* **68**, 121101 (2025).
146. Liu, H., Sferazza, C. & Abbeel, P. Chain of hindsight aligns language models with feedback. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2302.02676> (2023).
147. Sallam, M. ChatGPT utility in healthcare education, research, and practice: systematic review on the promising perspectives and valid concerns. *Healthcare* **11**, 887 (2023).
148. Tian, S. et al. Opportunities and challenges for ChatGPT and large language models in biomedicine and health. *Brief. Bioinform.* **25**, bbad493 (2024).
149. Li, H. et al. Multi-step jailbreaking privacy attacks on ChatGPT. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2304.05197> (2023).
150. Wei, A., Haghtalab, N. & Steinhardt, J. Jailbroken: how does LLM safety training fail? In *Adv. Neural Inform. Process. Syst.* (eds Oh, A. et al.) 80079–80110 (Curran Associates, 2023).
151. Meskó, B. & Topol, E. J. The imperative for regulatory oversight of large language models (or generative AI) in healthcare. *npj Digit. Med.* **6**, 120 (2023).
152. Derraz, B. et al. New regulatory thinking is needed for AI-based personalised drug and cell therapies in precision oncology. *npj Precis. Oncol.* **8**, 23 (2024).
153. Mökander, J., Schuett, J., Kirk, H. R. & Floridi, L. Auditing large language models: a three-layered approach. *AI Ethics* **4**, 1085–1115 (2024).
154. Liu, F. et al. Large language models are poor clinical decision-makers: a comprehensive benchmark. In *Proc. Conf. Empir. Methods Nat. Lang. Process.* (eds Al-Onaizan, Y. et al.) 13696–13710 (ACL, 2024).
155. Yin, S. et al. A survey on multimodal large language models. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2306.13549> (2023).
156. Huang, H. et al. ChatGPT for shaping the future of dentistry: the potential of multi-modal large language model. *Int. J. Oral Sci.* **15**, 29 (2023).
157. Li, J., Liu, C., Cheng, S., Arcucci, R. & Hong, S. Frozen language model helps ECG zero-shot learning. *Proc. Mach. Learn. Res.* **227**, 402–415 (2023).
158. Enghardt, Z. et al. Exploring and characterizing large language models for embedded system development and debugging. In *Proc. Extend. Abstr. CHI Conf. Hum. Factor. Comput. Syst.* (eds Mueller, F. et al.) 150:1–150:9 (ACM, 2024).
159. Mello, M. M. & Guha, N. ChatGPT and physicians’ malpractice risk. *JAMA Health Forum* **4**, e231938 (2023).
160. Mekki, Y. M. & Zughaier, S. M. Teaching artificial intelligence in medicine. *Nat. Rev. Bioeng.* **2**, 450–451 (2024).
161. Beltagy, I., Lo, K. & Cohan, A. SciBERT: a pretrained language model for scientific text. In *Proc. Conf. Empir. Methods Nat. Lang. Process. & 9th Int. Joint Conf. Nat. Lang. Proces.* (eds Inui, K. et al.) 3615–3620 (ACL, 2019).
162. Alrowili, S. & Shanker, V. Large biomedical question answering models with ALBERT and ELECTRA. In *Conf. Labs Eval. Forum* 213–220 (2021).
163. Gururangan, S. et al. Don’t stop pretraining: adapt language models to domains and tasks. In *Proc. 58th Annu. Meet. Assoc. Comput. Linguist.* (eds Jurafsky, D. et al.) 8342–8360 (ACL, 2020).
164. Yasunaga, M., Leskovec, J. & Liang, P. Linkbert: pretraining language models with document links. In *Proc. 60th Annu. Meet. Assoc. Comput. Linguist.* Vol. 1 (eds Muresan, S. et al.) 8003–8016 (ACL, 2022).
165. Peng, Y., Yan, S. & Lu, Z. Transfer learning in biomedical natural language processing: an evaluation of BERT and ELMo on ten benchmarking datasets. In *Proc. 18th BioNLP Workshop Shared Task* (eds Demner-Fushman, D. et al.) 58–65 (ACL, 2019).
166. Phan, L. N. et al. SciFive: a text-to-text transformer model for biomedical literature. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2106.03598> (2021).
167. Lu, Q., Dou, D. & Nguyen, T. ClinicalT5: a generative language model for clinical text. In *Find. Assoc. Comput. Linguist.* (eds Goldberg, Y. et al.) 5436–5443 (ACL, 2022).
168. Jin, Q. et al. MedCPT: contrastive pre-trained transformers with large-scale PubMed search logs for zero-shot biomedical information retrieval. *Bioinformatics* **39**, btad651 (2023).
169. Yasunaga, M. et al. Deep bidirectional language-knowledge graph pretraining. In *Proc. 36th Int. Conf. Neural Inform. Process. Syst.* (eds Koyejo, S. et al.) 37309–37323 (Curran Associates, 2022).
170. Venigalla, A., Frankle, J. & Carbin, M. BioMedLM: a domain-specific large language model for biomedical text. *MosaicML* <https://medium.com/@MosaicML/pubmed-gpt-a-domain-specific-large-language-model-for-biomedical-text-567b18e2b11> (2022).
171. Gao, W. et al. OphGLM: training an ophthalmology large language-and-vision assistant based on instructions and dialogue. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2306.12174> (2023).
172. Chen, Y. et al. BianQue: balancing the questioning and suggestion ability of health LLMs with multi-turn health conversations polished by ChatGPT. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2310.15896> (2023).

173. Wang, G., Yang, G., Du, Z., Fan, L. & Li, X. ClinicalGPT: large language models finetuned with diverse medical data and comprehensive evaluation. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2306.09968> (2023).
174. Zhang, H. et al. HuatuoGPT, towards taming language model to be a doctor. In *Find. Assoc. Computat. Linguist.* (eds Bouamor, H. et al.) 10859–10885 (ACL, 2023).
175. Luo, Y. et al. BioMedGPT: open multimodal generative pre-trained transformer for biomedicine. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2308.09442> (2023).
176. Ferber, D. et al. Gpt-4 for information retrieval and comparison of medical oncology guidelines. *NEJM AI* **1**, Alcs2300235 (2024).
177. Chen, Z. et al. MEDITRON-70B: scaling medical pretraining for large language models. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2311.16079> (2023).
178. He, X., Zhang, Y., Mou, L., Xing, E. & Xie, P. PathVQA: 30000+ questions for medical visual question answering. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2003.10286> (2020).
179. Johnson, A. E. et al. MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports. *Sci. Data* **6**, 317 (2019).
180. Yang, L. et al. Advancing multimodal medical capabilities of Gemini. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2405.03162> (2024).
181. Liévin, V., Hother, C. E., Motzfeldt, A. G. & Winther, O. Can large language models reason about medical questions? *Patterns* **5**, 100943 (2024).
182. Sun, Z., Luo, C., Liu, Z. & Huang, Z. Conversational disease diagnosis via external planner-controlled large language models. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2404.04292> (2024).
183. Dong, H. et al. Automated clinical coding: what, why, and where we are? *npj Digit. Med.* **5**, 159 (2022).
184. D'Onofrio, G. et al. Emotion recognizing by a robotic solution initiative (EMOTIVE project). *Sensors* **22**, 2861 (2022).
185. Bengio, Y., Ducharme, R. & Vincent, P. A neural probabilistic language model. In *Proc. 14th Int. Conf. Neural Inform. Process. Syst.* (eds Leen, T. K. et al.) 893–899 (MIT press, 2000).
186. Mikolov, T., Karafiát, M., Burget, L., Černocký, J. & Khudanpur, S. Recurrent neural network based language model. In *Proc. 11th Annu. Conf. Int. Speech Commun. Assoc.* (eds Kobayashi, T. et al.) 1045–1048 (ISCA, 2010).
187. Sundermeyer, M., Ney, H. & Schlüter, R. From feedforward to recurrent LSTM neural networks for language modeling. In *IEEE/ACM Transact. Audio Speech Lang. Process.* (ed. Li, H.) 517–529 (IEEE, 2015).
188. Hochreiter, S. & Schmidhuber, J. Long short-term memory. *Neural Comput.* **9**, 1735–1780 (1997).
189. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proc. Conf. North Am. Chapt. Assoc. Comput. Linguist.* (eds Burstein, J. et al.) 4171–4186 (ACL, 2019).
190. Vaswani, A. et al. Attention is all you need. In *Proc. 31st Int. Conf. Neural Inform. Process. Syst.* (eds von Luxburg, U. et al.) 6000–6010 (Curran Associates, 2017).
191. Kaplan, J. et al. Scaling laws for neural language models. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2001.08361> (2020).
192. Hoffmann, J. et al. An empirical analysis of compute-optimal large language model training. In *Proc. 36th Int. Conf. Neural Inform. Process. Syst.* (eds Koyejo, S. et al.) 30016–30030 (Curran Associates, 2022).
193. He, P., Liu, X., Gao, J. & Chen, W. DeBERTa: decoding-enhanced BERT with disentangled attention. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2006.03654> (2021).
194. Radford, A. et al. Language models are unsupervised multitask learners. *OpenAI blog* **1**, 9 (2019).
195. The Vicuna team. Vicuna: an open-source chatbot impressing GPT-4 with 90% ChatGPT quality. *LMSYS ORG* <https://lmsys.org/blog/2023-03-30-vicuna/> (2023).
196. Jiang, A. Q. et al. Mistral 7B. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2310.06825> (2023).
197. Bai, J. et al. Qwen technical report. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2309.16609> (2023).
198. Lewis, M. et al. BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proc. 58th Annu. Meet. Assoc. Comput. Linguist.* (eds Jurafsky, D. et al.) 7871–7880 (ACL, 2020).
199. Tay, Y. et al. UL2: unifying language learning paradigms. Preprint at *arXiv* <https://doi.org/10.48550/arXiv.2205.05131> (2022).

Acknowledgements

This work was supported in part by the Pandemic Sciences Institute at the University of Oxford, the National Institute for Health Research (NIHR) Oxford Biomedical Research Centre, an NIHR Research Professorship, a Royal Academy of Engineering Research Chair, the Wellcome Trust-funded VITAL project, the UK Research and Innovation, the Engineering and Physical Sciences Research Council, and the InnoHK Hong Kong Centre for Cerebro-cardiovascular Engineering (COCHE), the Clarendon Fund, and the Magdalen Graduate Scholarship.

Author contributions

H.Z. and F.L. conceived and designed the study. H.Z., F.L., B.G., X.Z., J.H. and J.W. conducted the literature review, performed data analysis and drafted the manuscript. All authors contributed to the interpretation and final manuscript preparation. All authors read and approved the final manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s44222-025-00279-5>.

Peer review information *Nature Reviews Bioengineering* thanks Jakob Kather and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.

Related links

Alpaca: https://github.com/tatsu-lab/stanford_alpaca
Bard: <https://gemini.google.com/>
ChatGPT: <https://chat.openai.com/>
Claude-3: <https://www.anthropic.com/news/claude-3-family>
HealthcareMagic: <https://www.healthcaremagic.com/>
iCliniq: <https://www.icliniq.com/>
LLaMA-3: <https://github.com/meta-llama/llama3>
MedLlama3-v20: <https://huggingface.co/ProbeMedicalYonseiMILab/medllama3-v20>
OpenBioLLM: <https://huggingface.co/aaditya/OpenBioLLM-Llama3-70B>
PubMed: <https://pubmed.ncbi.nlm.nih.gov/>
PubMed Central (PMC): <https://www.ncbi.nlm.nih.gov/pmc/>
ShareGPT: <https://sharegpt.com/>

© Springer Nature Limited 2025