

Fine-Tuning LLMs for Specialized Use Cases

D.M. Anisuzzaman, PhD, Jeffrey G. Malins, PhD, Paul A. Friedman, MD, Zachi I. Attia, PhD

MAYO CLINIC PROCEEDINGS:
DIGITAL HEALTH



PII: S2949-7612(24)00114-7

DOI: <https://doi.org/10.1016/j.mcpdig.2024.11.005>

Reference: MCPDIG 184

To appear in: *Mayo Clinic Proceeding: Digital Health*

Received Date: 2 August 2024

Revised Date: 6 November 2024

Accepted Date: 18 November 2024

Please cite this article as: Anisuzzaman DM, Malins JG, Friedman PA, Attia ZI, Fine-Tuning LLMs for Specialized Use Cases, *Mayo Clinic Proceeding: Digital Health* (2024), doi: <https://doi.org/10.1016/j.mcpdig.2024.11.005>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2024 Published by Elsevier Inc on behalf of Mayo Foundation for Medical Education and Research.

Fine-Tuning LLMs for Specialized Use Cases

D M Anisuzzaman, PhD; Jeffrey G. Malins, PhD; Paul A. Friedman, MD; Zachi I. Attia, PhD

Affiliation: Department of Cardiovascular Medicine, Mayo Clinic, Rochester, MN

Conflict of Interest Disclosures: D M Anisuzzaman, Jeffrey G. Malins, Paul A. Friedman, and Zachi I. Attia have invented algorithms licensed to UltraSight and may benefit from algorithm commercialization via Mayo Clinic. None of these relations with industry are related in any way to the content of the current submission.

Potential Competing Interest: Given their role as Editorial Board Members, Dr. Zachi I. Attia and Dr. Paul A. Friedman had no involvement in the peer-review of this article and have no access to information regarding its peer-review.

Abstract Presentation

None.

Correspondence:

Zachi I. Attia, Ph.D.

Assistant Professor of Medicine, Mayo Medical School

Co-Director, Artificial Intelligence in Cardiology

Mayo Clinic, 200 1st Street SW, Rochester, MN, 55905

E-mail: attia.itzhak@mayo.edu

Word count: 4394

Number of tables: 1

Number of figures: 2

Key Words: Large language models; artificial intelligence; deep learning; fine-tuning, ChatGPT

Abstract

Large language models (LLMs) are a type of artificial intelligence, and operate by predicting and assembling sequences of words that are statistically likely to follow from a given text input. With this basic ability, LLMs are able to answer complex questions and follow extremely complex instructions. Products created using LLMs such as ChatGPT by OpenAI and Claude by Anthropic have created a huge amount of traction and user engagements and revolutionized the way we interact with technology, bringing a new dimension to human-computer interaction. Fine-tuning is a process in which a pre-trained model, such as an LLM, is further trained on a custom dataset to adapt it for specialized tasks or domains. In this review, we outline some of the major methodological approaches and techniques that can be used to fine-tune LLMs for specialized use cases, and enumerate the general steps required for carrying out LLM fine-tuning. We then illustrate a few of these methodological approaches by describing several specific use cases of fine-tuning LLMs across medical subspecialties. Finally, we close with a consideration of some of the benefits and limitations associated with fine-tuning LLMs for specialized use cases, with an emphasis on specific concerns in the field of medicine.

Abbreviations and Acronyms:

AI = artificial intelligence

LLM = large language model

GPU = graphical processing unit

RLHF = reinforcement learning from human feedback

PPO = proximal policy optimization

AUC = area under the curve

PEFT = parameter-efficient fine-tuning (PEFT)

LoRA = low-rank adaptation

QLoRA = quantized low-rank adaptation

RAG = Retrieval augmented generation

CoT = chain-of-thought

Introduction

Large Language Models (LLMs), a specialized subset of artificial intelligence (AI), are designed to generate text through a process known as autoregression (often leading them to be termed autoregressive LLMs). These models operate by predicting and assembling sequences of words that are statistically likely to follow from a given text input, thereby enabling them to produce coherent and relevant sentences. The models can accept conversational input as text or via speech (using language recognition), and can generate outputs at various levels ranging from technical/professional to that of a high school education and more. They can summarize vast quantities of data, have access to unimaginably large volumes of information, and stand to make this available, easily, to the user. The public release of ChatGPT has opened the public's imagination, and given a glimpse into an information-rich future.

These capabilities allow LLMs to perform a variety of general-purpose tasks such as answering questions, completing sentences, and even generating entire articles. One of the breakthroughs that led to the creation of LLMs is the use of foundational models that process and comprehend natural language using deep learning methods. The two primary ideas of foundation models are self-supervised learning and scale. In self-supervision, instead of training a model to perform a task that requires explicit annotations, the model learns from the vast amounts of unlabeled data available, extracting patterns and understanding context without human intervention. In addition to being more scalable, self-supervised tasks can allow a model to anticipate a portion of the inputs, which makes the model richer and potentially more valuable than models trained on a more constrained label space. Once the model learns the foundational patterns of language, the same model can then be applied using transfer learning followed by fine-tuning, which enables the model to learn to perform more specific tasks using a smaller set

of labeled samples. For scale, the era of the internet provides a nearly limitless amount of data ¹ and, coupled with advances in computing power, enables the training of models on an unprecedented scale using Graphical Processing Units (GPUs). Together, these developments—enhanced by innovations such as the Transformer model architecture ²—have significantly propelled the capabilities and applications of LLMs. A general workflow of LLM fine-tuning for specialized use cases is shown in Figure 1.

Some existing LLMs to date are Alpaca ³, BERT ⁴, BLOOM ⁵, Claude ⁶, Cohere ⁷, Ernie ⁸, Falcon ⁹, Flan ¹⁰, Gemini ¹¹, Gemma ¹², GPT-3.5 ¹³, GPT-4 ¹⁴, LaMDA ¹⁵, LLaMA ¹⁶, Mistral ¹⁷, MPT ¹⁸, Orca ¹⁹, PaLM 2 ²⁰, Phi-1 ²¹, StableLM ²², T5 ²³, Vicuna ²⁴, and Zephyr ²⁵. All of these models were developed to handle language-related tasks by different for-profit and non-profit organizations such as Google, Meta, and Stanford. Though the majority of the models were created as general task models, some were developed for specialized tasks such as language translation, human-like chat, and code generation.

In addition to anticipating subsequent text, because models are trained with billions of tokens, many words map to multiple tokens (i.e., they are represented by word vectors), enabling mathematical connections between multiple meanings of a term. For example, Paris will have connections to France, city, capital, and so on, so that the relationships between Paris and France and London and England may be utilized by an LLM. A limitation of LLMs is that after training is completed, a model no longer "learns" or acquires new information, and the information it was trained on may be general (such as Wikipedia), but not well-suited to a specific task. These limitations can be mitigated with fine-tuning to better sculpt an LLM to address a specific field (such as medicine or law), and retrieval augmented generation, which provides additional

information that a model may use to address questions, and which is particularly useful if that additional information was not included in the model's training.

In the domain of healthcare, a number of LLMs have been fine-tuned to perform tasks associated with pre-consultation, diagnosis, management, and prediction of future medical outcomes, as well as medical education and medical writing ²⁶⁻²⁸. LLMs specific to the medical domain include BioBERT ²⁹, BioGPT ³⁰, BioMistral ³¹, ChatDoctor ³², Clinical Camel ³³, DoctorGLM ³⁴, Med-Alpaca ³⁵, Med-PaLM ³⁶, Med-PaLM 2 ³⁷, Med42-v2 ³⁸, Meditron-70b ³⁹, OpenBioLLM-70B ⁴⁰, and PMC-LLaMA ⁴¹. One particularly powerful use of models such as these is obtaining answers to questions rather than links to articles, with the caveat that using systems not designed to address medical questions may be inaccurate ⁴². LLMs can potentially be utilized for any task that requires "reading" text, and summarizing it, or extracting pertinent information. Examples of data extraction uses could include review of medical records to create a discharge summary, identifying and summarizing all the risks for stroke in a patient with atrial fibrillation, or determining preoperative surgical risk using standardized scoring criteria. A list of potential uses of LLMs in medicine along with specific examples is provided in Table 1.

Table 1. Uses of LLMs in medicine. The first three columns of this table were generated by ChatGPT 4.0 on May 6, 2024. The final column with example usages was added by the authors.

| Description of LLM Medical Uses | Strengths | Limitations | Example Usage |
|---------------------------------|---|--|--|
| Medical research assistance | Can quickly synthesize and summarize existing medical literature, helping researchers stay up-to-date with recent developments. | May not have access to the most recent studies due to training data cutoffs; could miss context or nuance in highly specialized areas. | Documentation for clinical trials ⁴³ |
| Clinical decision support | Provides support in diagnosing complex cases by suggesting possible diagnoses based on symptoms and medical history. | Relies on the data it was trained on, which may not include rare diseases or latest treatment modalities. | Differentiation between abdominal pathologies ⁴⁴ |
| Patient interaction automation | Handles routine inquiries from patients, such as explaining medical procedures and advising on medication schedules. | May lack the empathetic nuances that human interaction provides; risk of miscommunication in complex scenarios. | Answering cataract surgery-related questions ⁴⁵ |
| Medical education and training | Assists in the education of medical students and professionals by providing explanations, generating quizzes, and simulating patient cases. | Might not perfectly mimic the unpredictability of real-life medical cases; information may become outdated. | Interactive practice cases to evaluate medical reasoning ⁴⁶ |
| Documentation and reporting | Helps in generating and organizing medical reports, thereby reducing the administrative burden on healthcare providers. | Possible issues with accuracy and privacy concerns; needs constant verification. | Generation of radiology reports based on chest X-rays ⁴⁷ |
| Treatment plan management | Suggests treatment plans based on clinical guidelines and individual patient data. | May not incorporate experiential learning or adapt to unconventional cases as effectively as a human would. | Assistance with complex decision-making for breast cancer care ⁴⁸ |
| Support for remote areas | Provides medical information and support in remote areas where medical expertise is limited. | Dependence on internet connectivity; may not handle local medical practices or non-standard treatments well. | Providing community health workers with contextually appropriate medical knowledge ⁴⁹ |

In the following sections, we first outline some of the major approaches and techniques for fine-tuning LLMs in the medical domain, and touch on retrieval augmented generation. Then, we describe specialized use cases in which LLMs have been fine-tuned for medical applications across various medical subspecialties. Following this, we close with a consideration of some of the benefits and key limitations associated with fine-tuning LLMs in the medical domain.

Fine-tuning Methodology

Fine-tuning is a process in which a pre-trained model is adapted for particular tasks or domains by continuing to train the model using only a domain-specific dataset that is different than the original dataset used to train the base model. Various fine-tuning strategies and approaches are employed to adjust the model parameters to a specific need. Some fine-tuning approaches are briefly described here:

Supervised fine-tuning. With this approach, every input data point is linked to a label, and the model is trained on a task-specific labeled dataset. The model learns to modify its parameters to anticipate these labels as precisely as possible. Some supervised fine-tuning techniques are:

1. *Transfer learning*: In this approach, a model is first initialized with saved weights from a model pre-trained on a large, general dataset, and then is subsequently trained with limited task-specific data. Weights refer to the learned parameters of a model that has been trained on a large dataset for a specific task, which represent the knowledge the model has gained during its training process, encapsulating features and patterns relevant to the task it was originally trained on.

2. *Multi-task learning*: Here, models are fine-tuned on numerous related tasks, taking advantage of their similarities and differences, in order to maximize performance. For example, with a CNN model trained on a generic large dataset (e.g., KINETICS400), one can perform some specific tasks (e.g., estimating left ventricular ejection fraction, patient age, and patient sex from an echocardiogram) with a much smaller dataset by leveraging the generic features the model learned from the large dataset.
3. *Instruction-tuning*: Instruction-tuning involves fine-tuning a pre-trained LLM to follow specific task instructions, such as translation, summarization, or question answering. For example, in translation, the model is trained on examples in which each input includes an instruction like "Translate the following sentence from English to French," followed by an English sentence and its French translation. After fine-tuning, the model learns to follow translation instructions and can generalize to translate new sentences.

Reinforcement learning from human feedback (RLHF). This method uses the knowledge of human evaluators; in addition, it also allows the model to adjust and develop in response to real-world input, resulting in enhanced and more efficient applications. Some standard RLHF techniques are:

1. *Reward modeling*: In this method, the model generates multiple potential outputs or actions, which are subsequently assessed by human evaluators who assign a ranking or rating based on their quality. The model utilizes these human-provided assessments to generate predictions and adapt its behavior to optimize the anticipated rewards.
2. *Proximal policy optimization (PPO)*: PPO modifies the language model's policy to maximize the expected reward. A policy refers to the strategy or set of rules that a reinforcement learning agent uses to make decisions in an environment. For example, in

PPO, the policy determines how a robotic arm should move to pick up objects based on visual inputs from a camera. PPO's primary goal is to make policy improvements while ensuring the modifications don't deviate too much from the previous policy. To achieve this balance, the policy update process introduces a constraint that prohibits detrimental large updates while permitting advantageous minor updates. Compared to other reinforcement learning techniques, PPO is more reliable and effective.

3. *Comparative ranking*: In this method, the model produces several outputs or actions, which human investigators then rank according to compatibility or quality. The model then modifies its behavior to generate higher-ranked outputs. This method provides relative and better feedback to the model by ranking multiple outputs rather than individual outputs.
4. *Preference feedback*: This technique involves the model generating several outputs and human experts selecting among them, leading the model to modify its behavior accordingly. This method is useful when assigning a numeric value (reward) to an output is difficult. It is an effective method of fine-tuning the model in practical applications.

Fine-Tuning Pipeline

To carry out a fine-tuning process for a specialized use case, there are several generic steps that include the following ⁵⁰:

1. *Data preparation*: Dataset preparation for LLM model fine-tuning is entirely task-specific. In general, the model must be presented with some blocks of text. Many datasets are available to fine-tune an LLM ⁵¹⁻⁵³. One must follow the instructions given for each dataset to prepare the data for fine-tuning. For custom datasets, depending on the task,

the dataset preparation may include data cleaning, normalizing for missing values, and formatting the text to align with the model's input requirements.

2. *Selecting the appropriate pre-trained model:* There are several LLMs to date (BERT, Cohere, GPT-4, LLaMA, Mistral, etc.; access to non-open source models will require working with the owners of the models), and choosing the appropriate one that complies with the demands of the target task is essential. For fine-tuning an LLM model for a specific dataset or task, one must have a good grasp of the model architecture and the input and output requirements. Depending on the available resources, the model weights and number of parameters should be considered when selecting a model. Finally, the performance of the model on the relevant target task should be considered during model selection.
3. *Fine-tuning the model:* LLM fine-tuning also includes basic hyperparameter tuning, adjustments of the learning rate, batch size, regularization, optimizer, number of epochs, etc. As LLMs are trained on vast amounts of data, overfitting a small dataset for a specific task could be a likely event. Careful tuning of the hyperparameters guarantees that the model learns efficiently and doesn't overfit when applied to new data.
4. *Validation:* Validation of an LLM is complex. In predictive AI, a specific input and output are expected. For example, a neural network might assess a medical image such as an ECG, or a medical video such as an echocardiogram, and be tasked with determining the ejection fraction (heart pump strength). The ejection fraction can be manually measured by a human to assess the performance of the AI tool. In contrast, generative AI, such as LLMs, generate new text or images, the performance of which may be harder to grade. If asked to create a poem, how does one assess the quality of it?

In general practice, there are two types of validation to perform: (1) internal validation, which is used to select the best model, monitor the model's learning process, and callback and stop model training with certain criteria; (2) validation on a hold-out test set to evaluate model performance for real world applications. Although in general, some commonly used metrics for model evaluation include accuracy, area under the curve (AUC), precision, recall, etc.; LLM model evaluation may require some careful task-specific metric selection⁵⁴. Some key performance metrics used in LLM evaluation are:

- Accuracy: measures the model's ability to produce correct responses to prompts.
- Perplexity: measures uncertainty in predicting the next token.
- ROUGE scores: compares an LLM's output with a set of reference summaries.
- Diversity: evaluates the variety of responses generated.
- Disparity analysis: identifies and mitigates biases within model responses.
- Coh-Metrix: analyzes logical consistency and clarity over longer stretches of text.
- Human evaluation: subjective assessment by human judges.

In medical LLM development, in cases for which models are fine-tuned for clinical prediction tasks in which the ground truth labels are well defined (e.g., predicting discharge events), evaluation typically involves statistical performance metrics like accuracy, precision, and so on. For more generative tasks in which the ground truth labels are not well defined (e.g., medical report summarization), human or domain expert evaluation is crucial to ensure that model outputs are clinically accurate and safe for real-world applications. For example, Singhal et al. developed Med-PaLM for answering medical questions, and had clinicians review outputs to ensure that the responses were

medically sound and factually accurate³⁶. Similarly, Serapio et al. fine-tuned LLMs for generating radiological impressions from chest CT scans, and had model outputs assessed by board-certified radiologists⁵⁵.

Beyond these general approaches, a number of specific techniques can be applied to fine-tune LLMs for specialized use cases. Some example techniques include the following:

1. *In-context learning*: In this approach, a pre-trained LLM is induced to perform a task using prompted examples. An example of this is few-shot learning, which involves giving the model a few "shots" or instances to learn a new task during inference. Few-shot learning aims to direct the model's predictions by providing examples and context specifically in the prompt, but importantly does not involve gradient-based training⁵⁶.
2. *Hyperparameter tuning*: This is a straightforward method that consists of manually modifying basic hyperparameters (i.e., learning rate, batch size, optimizer, number of epochs, etc.) of the model until the desired performance is obtained. This changes how the model "learns"; that is, how fast it learns, how to decide when training is completed, etc.
3. *Parameter-efficient fine-tuning (PEFT)*: PEFT is an efficient technique in which only a small portion of the parameters of an LLM are selectively modified during fine-tuning, typically by adding new layers or modifying existing ones in a task-specific manner. This method drastically lowers computational and storage needs while keeping performance comparable to complete fine-tuning. Some PEFT techniques are low-rank adaptation (LoRA)⁵⁷, quantized low-rank adaptation (QLoRA)⁵⁸, Prefix tuning, Prompt tuning, etc.

QLoRA is a very popular technique used for LLM fine-tuning due to its power of using much smaller amounts of memory than a full fine-tuning approach with the price of sacrificing some performance. For example, a full fine-tuning of the LLaMA 65B parameter model requires more than 780 GB of GPU memory, whereas using the QLoRA technique requires only 48GB of GPU memory ⁵⁸. This powerful technique is based on these highly technical ingredients:

- i. 4-bit NormalFloat representation of model parameters, whereas typically, parameters of trained models are stored in a 32-bit format. This technique divides model parameters into equally-sized buckets instead of equally-spaced buckets.
- ii. Double Quantization, a method that quantizes the quantization constants. In general, quantization converts datatypes with a larger number of bits to fewer bits (e.g., FP32 to 8-bit Integers). QLoRA uses the block-wise quantization technique which requires more memory than standard quantization but reduces bias significantly, thus retaining good performance.
- iii. LoRA ⁵⁷, which freezes the pre-trained model weights and injects trainable rank decomposition matrices into each layer of the Transformer architecture, greatly reducing the number of trainable parameters for downstream tasks.

In simpler terms, low rank adaptation finds a more compressed version of the LLM weights, and updates those weights. Although the compression may lose some data, under the assumption a lot of the model weights are redundant, leading only to a small decrease in performance relative to savings in memory and required compute power.

4. *Retrieval augmented generation (RAG)*: RAG is a technique that combines the capabilities of neural language models with information retrieval systems to enhance the generation of contextually rich and accurate responses. In RAG, when a query is received, the model first uses a retrieval system to fetch relevant documents or snippets from a large corpus, such as a database of scientific literature. These retrieved texts are then fed into a generative model, typically a transformer-based neural network, which integrates the retrieved information with its pre-trained knowledge to produce a coherent and informed response. This approach is particularly useful in domains where accuracy and specificity are critical, such as scientific research or technical support, as it allows the model to base its answers on up-to-date and source-specific data, providing citations and grounding its responses in existing literature. A comparison of the outputs generated by presenting the same medical question to a search engine, LLM, and RAG-based system is shown in Figure 2.

Specialized Use Cases in Medicine

Across many of the major subspecialties of medicine, LLMs are being fine-tuned to address specific issues, and practitioners are posing questions about how fine-tuned LLMs could revolutionize their fields. These subspecialties include but are not limited to: cardiology⁵⁹⁻⁶¹, dermatology⁶², digital pathology⁶³, gastroenterology and hepatology⁶⁴⁻⁶⁶, hematology⁶⁷, neurology⁶⁸⁻⁷⁰, obstetrics and gynecology^{71,72}, oncology⁷³⁻⁷⁵, ophthalmology⁷⁶, orthopedics⁷⁷, pediatrics^{78,79}, psychiatry⁸⁰⁻⁸², radiology^{83,84}, surgery^{85,86}, and urology^{87,88}. Although there are nuances according to specific subspecialties, many practitioners highlight the potential of fine-tuned LLMs to aid clinicians in areas such as clinical decision support, treatment planning, and patient consultation, as well as alleviate administrative burden associated with tasks such as

generating clinical notes, discharge reports, and medical billing. At the same time, many are concerned about the ethical, legal, and social implications of using such models.

In the subsequent section, we review some of these concerns. Prior to that, below we highlight a few specific methodologies and use cases that illustrate the general framework outlined in the previous sections.

Using RLHF to fine-tune LLMs in medicine. Mukherjee et al.⁸⁹ developed a constellation system called Polaris, which was composed of several agents. Their primary agent (focused on patient-friendly conversation) was developed in three stages: General Instruction Tuning, Conversation and Agent Tuning, and RLHF. The RLHF step was performed by registered nurses, who gave preference feedback on multiple responses. Zhao et al. developed Aquilia-Med⁹⁰, a bilingual medical LLM, using supervised fine-tuning and RLHF to tackle medical challenges. It was trained on large-scale Chinese and English medical datasets, with RLHF further aligning the model to improve performance in medical dialogues and multiple-choice questions.

Integration of LLMs into Electronic Health Record (EHR) systems. Several LLMs have been applied to Electronic Health Record (EHR) systems, providing benefits such as generating patient summaries from EHRs, assisting healthcare providers with more efficient decision-making, named entity recognition, medical note summarization, and predictive diagnosis⁹¹. Zhang et al. investigated the application of LLM fine-tuning to EHR audit log data for clinical prediction tasks, with a focus on discharge predictions⁹². Cui et al. evaluated the zero-shot and few-shot performance of LLMs on EHR-based disease prediction tasks, and proposed a novel approach that leverages collaborative LLM agents to enhance predictive performance⁹³. Li et al. fine-tuned an LLM named LlamaCare, and evaluated it on various

clinical tasks, such as generating discharge summaries, predicting mortality and length of stay, and more ⁹⁴.

Generation of echocardiography reports to streamline workflows. Echocardiography (echo) is one of the most widely employed imaging techniques for gaining insights into the structure and function of the heart. A typical echo report includes numerous measurements as well as text-based statements, or findings. These findings are summarized by a clinician to give an overall set of final impressions for the study. This is a time-consuming and error-prone process. To address this issue, Chao et al. ⁹⁵ leveraged several open-source LLMs to generate echo reports using either zero-shot learning (for Flan-T5, MedAlpaca, Llama-2, and Zephyr) or QLoRA fine-tuning (Llama-2 and Zephyr). Using a training dataset of 95,506 echo reports, the authors observed that EchoGPT, which is a Llama-2 model trained using instruction fine-tuning with QLoRA, outperformed other LLMs on critical performance metrics. In addition, when four echocardiography-board-certified cardiologists were asked to rate reports generated by EchoGPT for 30 randomly selected cases, the generated reports were rated similarly to reports generated by cardiologists for these same cases (in completeness, conciseness, correctness, and clinical utility). Based on these results, the authors argue that EchoGPT could be used as a 'co-pilot' for report generation, which would allow for considerable streamlining of the echo report workflow. With that said, the authors stress that draft reports generated by EchoGPT should still be reviewed and approved by clinicians, noting that some hallucinations were observed in reports generated by EchoGPT (albeit not as many as were observed for zero-shot learning).

Identifying eligible patients for clinical trials. Randomized clinical trials are a cornerstone of medical research, yet it can be a challenge to identify patients who meet all the inclusion and exclusion criteria for a clinical trial. To leverage the power of LLMs to assist with

participant recruitment for clinical trials, Guan et al.⁹⁶ developed CohortGPT, which is built upon ChatGPT and GPT-4. CohortGPT can take input text from unstructured or semi-structured data, such as clinical notes and radiology reports, in order to designate disease labels associated with the input text. To develop this model, the authors made use of a technique called chain-of-thought (CoT) prompting, which is a type of in-context learning that guides LLMs to learn task-specific logical chains, which detail how correct answers are deduced from given information. Using the CoT technique in conjunction with reinforcement learning, Guan et al.⁹⁶ trained a policy model to dynamically select chain-of-thought (CoT) samples. They then presented these CoT samples to a prompt model alongside knowledge graphs, which can be thought of as rules detailing the relationships between different concepts, such as that cardiomegaly is a type of heart disease, or that scoliosis is a type of spine disease. Using thousands of publicly available radiology reports in the Indiana chest X-ray collection⁹⁷ and MIMIC-CXR⁹⁸ datasets, Guan et al.⁹⁶ demonstrated that CohortGPT can reliably classify report text as being associated with specific disease labels. Based on these results, the authors argue that CohortGPT can be useful not only for patient recruitment for clinical trials, but also for other medical applications such as diagnosis and treatment optimization. Furthermore, even though CohortGPT was built upon ChatGPT and GPT-4, the model can be implemented in any open-source LLM.

Benefits and Limitations of Fine-Tuning LLMs for Specialized Use Cases

To build a reliable real world LLM-based application, fine-tuning is a necessary and crucial step, as it fills in the gap between general knowledge and domain-specific expertise for that application. Some benefits of LLM fine-tuning are the following:

1. *Domain-specific knowledge*: general purpose LLMs may not have enough domain-specific knowledge.
2. *Specific task optimization*: general purpose LLMs can be optimized for specific tasks (e.g., health report summarization, disease detection from a report, etc.).
3. *Data efficiency*: fine-tuning works well with smaller quantities of labeled data as it involves using pre-trained LLM(s) trained on huge datasets.
4. *Better performance*: fine-tuning often leads to improved performance as the model learns domain-specific knowledge to perform relevant tasks while preserving out-of-domain knowledge.
5. *Resource efficiency*: fine-tuning requires less resources in terms of time and memory than training a general purpose LLM from scratch.

With that said, there are several critical limitations to LLMs to consider when fine-tuning for specialized tasks^{99,100}. A few of these are:

1. *Hallucinations*: these refer to situations in which model output contains inaccurate or non-factual information^{42,99}. In the medical domain for example, these could consist of findings that are not actually present in a study report. Addressing hallucinations could involve processes such as inducing a model to provide a reasoning process or confidence score associated with model output. For example, the Medical Domain Hallucination Test (Med-HALT) has been designed to evaluate and reduce hallucinations in the

medical domain, and includes metrics for hallucinations associated with reasoning and memory ¹⁰¹.

2. *Legal and safety concerns*: for example, in the medical domain, the data to be used for fine-tuning may contain sensitive patient information that needs to be safeguarded. In addition, if model output is used to guide treatment decisions for patients, incorrect output (such as hallucinations) could be harmful. This is why authors such as Chao et al. ⁹⁵ emphasize the critical need for human review of model output. In addition, cybersecurity measures such as the use of pseudonyms can enhance the privacy and security of patient data ¹⁰².
3. *Biases in training datasets*: fine-tuned LLMs can inherit biases from the pre-trained models upon which they are built, and there is a critical need to employ techniques that mitigate this bias ^{103,104}. In medicine, this bias has the potential to exacerbate health inequities if not addressed ¹⁰⁵. Some techniques for mitigating bias include prompt engineering, debiasing algorithms, and continuous monitoring of model performance ¹⁰⁶.
4. *Lack of domain-specific data*: depending on the extent to which a specific use case is specialized, there may not be sufficient quantities of domain-specific data to fine-tune an LLM using certain approaches. Here, techniques such as in-context learning or PEFT may be more appropriate than full fine-tuning.
5. *Data leakage*: many of the pre-trained LLMs do not report which data was used for training, so if open datasets are utilized for fine-tuning, this data may have already been used for training the base model. This can lead to data leakage from the validation set to the model, resulting in overly optimistic performance. Addressing this concern will involve greater transparency on the part of developers when describing training datasets,

and careful selection of pre-trained LLMs that provide information about the source, quality, and quantity of training data ¹⁰².

Conclusion

Large language models are poised to transform medicine. In written form or verbally, they can summarize vast amounts of information, may prevent important pieces of information from being missed, and can meaningfully tap into vast stores of literature to inform clinicians at the point of care when meeting with a patient. However, much remains unproven including how to ensure the information is reliable, privacy is preserved, and answers are tuned to usefully guide medical professionals.

References

1. Bommasani R, Hudson DA, Adeli E, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:210807258*. 2021;
2. Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need. *Advances in Neural Information Processing Systems*. 2017;30
3. Taori R, Gulrajani I, Zhang T, et al. Stanford alpaca: An instruction-following llama model. 2023.
4. Devlin J, Chang M-W, Lee K, Toutanova K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*. 2018;
5. Le Scao T, Fan A, Akiki C, et al. Bloom: A 176b-parameter open-access multilingual language model. 2023;
6. Anthropic. Introducing Claude. Accessed April 24, 2024. <https://www.anthropic.com/news/introducing-claude>
7. Cohere. Cohere: The Leading enterprise AI platform. Accessed April 24, 2024. <https://cohere.com/>
8. BaiduResearch. ERNIE Bot: Baidu's knowledge-enhanced large language model built on full AI stack technology. Accessed April 25, 2024. <http://research.baidu.com/Blog/index-view?id=183>
9. ZXhang YX, Haxo YM, Mat YX. Falcon llm: A new frontier in natural language processing. *AC Investment Research Journal*. 2023;220(44)
10. GoogleResearch. Introducing FLAN: More generalizable language models with instruction fine-tuning. Accessed April 25, 2024. <https://research.google/blog/introducing-flan-more-generalizable-language-models-with-instruction-fine-tuning/>
11. Team G, Anil R, Borgeaud S, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*. 2023;
12. Team G, Mesnard T, Hardin C, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*. 2024;
13. OpenAI. Models - GPT 3.5 Turbo. Accessed April 25, 2024. <https://platform.openai.com/docs/models/gpt-3-5-turbo>
14. Achiam J, Adler S, Agarwal S, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*. 2023;
15. Thoppilan R, De Freitas D, Hall J, et al. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*. 2022;
16. Touvron H, Lavril T, Izacard G, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*. 2023;
17. Jiang AQ, Sablayrolles A, Mensch A, et al. Mistral 7B. *arXiv preprint arXiv:2310.06825*. 2023;
18. HuggingFace. MPT. Accessed April 25, 2024. https://huggingface.co/docs/transformers/main/model_doc/mpt
19. KDnuggets. Orca LLM: Simulating the reasoning processes of ChatGPT. Accessed April 25, 2024. <https://www.kdnuggets.com/2023/06/orca-llm-reasoning-processes-chatgpt.html>
20. Anil R, Dai AM, Firat O, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*. 2023;
21. Gunasekar S, Zhang Y, Aneja J, et al. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*. 2023;
22. Bellagente M, Tow J, Mahan D, et al. Stable LM 2 1.6 B technical report. *arXiv preprint arXiv:2402.17834*. 2024;
23. Raffel C, Shazeer N, Roberts A, et al. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*. 2020;21(140):1-67.
24. Chiang W-L, Li Z, Lin Z, et al. Vicuna: An open-source chatbot impressing GPT-4 with 90%* ChatGPT quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023). 2023;2(3):6.
25. Tunstall L, Beeching E, Lambert N, et al. Zephyr: Direct distillation of lm alignment. *arXiv preprint arXiv:2310.16944*. 2023;

26. Yang R, Tan TF, Lu W, et al. Large language models in health care: Development, applications, and challenges. *Health Care Science*. 2023;2(4):255-263.
27. Thirunavukarasu AJ, Ting DSJ, et al. Large language models in medicine. *Nature medicine*. 2023;29(8):1930-1940.
28. Kraljevic Z, Bean D, Shek A, et al. Foresight—a generative pretrained transformer for modelling of patient timelines using electronic health records: a retrospective modelling study. *The Lancet Digital Health*. 2024;6(4):e281-e290.
29. Lee J, Yoon W, Kim S, et al. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*. 2020;36(4):1234--1240.
30. Luo R, Sun L, Xia Y, et al. BioGPT: generative pre-trained transformer for biomedical text generation and mining. *Briefings in bioinformatics*. 2022;23(6):bbac409.
31. Labrak Y, Bazoge A, Morin E, et al. Biomistral: A collection of open-source pretrained large language models for medical domains. *arXiv preprint arXiv:240210373*. 2024;
32. Li Y, Li Z, Zhang K, et al. Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge. *Cureus*. 2023;15(6)
33. Toma A, Lawler PR, Ba J, et al. Clinical camel: An open expert-level medical language model with dialogue-based knowledge encoding. *arXiv preprint arXiv:230512031*. 2023;
34. Xiong H, Wang S, Zhu Y, et al. Doctorglm: Fine-tuning your chinese doctor is not a herculean task. *arXiv preprint arXiv:230401097*. 2023;
35. Han T, Adams LC, Papaioannou J-M, et al. MedAlpaca—an open-source collection of medical conversational AI models and training data. *arXiv preprint arXiv:230408247*. 2023;
36. Singhal K, Azizi S, Tu T, et al. Large language models encode clinical knowledge. *Nature*. 2023;620(7972):172--180.
37. Singhal K, Tu T, Gottweis J, et al. Towards expert-level medical question answering with large language models. *arXiv preprint arXiv:230509617*. 2023;
38. Christophe C, Kanithi PK, Raha T, Khan S, Pimentel MAF. Med42-v2: A suite of clinical llms. *arXiv preprint arXiv:240806142*. 2024;
39. Chen Z, Cano AH, Romanou A, et al. Meditron-70b: Scaling medical pretraining for large language models. *arXiv preprint arXiv:231116079*. 2023;
40. Pal A, Sankarasubbu M. OpenBioLLMs: Advancing Open-Source Large Language Models for Healthcare and Life Sciences. Hugging Face. Accessed September 30, 2024, <https://huggingface.co/aaditya/OpenBioLLM-Llama3-70B>
41. Wu C, Lin W, Zhang X, et al. PMC-LLaMA: toward building open-source language models for medicine. *Journal of the American Medical Informatics Association*. 2024:ocae045.
42. Siontis KC, Attia ZI, Asirvatham SJ, Friedman PA. ChatGPT hallucinating: can it get any more humanlike? *Eur Heart J*. Feb 1 2024;45(5):321-323. doi:10.1093/eurheartj/ehad766
43. Markey N, El-Mansouri I, Renzonnet G, van Langen C, Meier C. From RAGs to riches: Using large language models to write documents for clinical trials. *arXiv preprint arXiv:240216406*. 2024;
44. Hager P, Jungmann F, Holland R, et al. Evaluation and mitigation of the limitations of large language models in clinical decision-making. *Nature medicine*. 2024;30(9):2613--2622.
45. Ramjee P, Sachdeva B, Golechha S, et al. CataractBot: An LLM-Powered Expert-in-the-Loop Chatbot for Cataract Patients. *arXiv preprint arXiv:240204620*. 2024;
46. Safranek CW, Sidamon-Eristoff AE, Gilson A, Chartash D. The role of large language models in medical education: applications and implications. *JMIR medical education*2023. p. e50945.
47. Wang Z, Liu L, Wang L, Zhou L. R2gengpt: Radiology report generation with frozen llms. *Meta-Radiology*. 2023;1(3):100033.
48. Griewing S, Knitza J, Boekhoff J, et al. Evolution of publicly available large language models for complex decision-making in breast cancer care. *Archives of Gynecology and Obstetrics*. 2024:1--14.

49. Gangavarapu A. Introducing L2M3, A Multilingual Medical Large Language Model to Advance Health Equity in Low-Resource Regions. *arXiv preprint arXiv:240408705*. 2024;
50. Turing. Fine-tuning LLMs: Overview, methods, and best practices. Accessed April 26, 2024. <https://www.turing.com/resources/finetuning-large-language-models>
51. Zhao J. LLMDataHub: Awesome datasets for LLM training. Accessed April 26, 2024. <https://github.com/Zjh-819/LLMDataHub>
52. HuggingFace. Datasets (filter Other by name "llm"). Accessed April 26, 2024. <https://huggingface.co/datasets?other=llm>
53. Liu Y, Cao J, Liu C, Ding K, Jin L. Datasets for Large Language Models: A Comprehensive Survey. *arXiv preprint arXiv:240218041*. 2024;
54. Aisera. LLM Evaluation metrics: Performance benchmark. Accessed April 26, 2024. <https://aisera.com/blog/llm-evaluation/>
55. Serapio A, Chaudhari G, Savage C, et al. An Open-source Fine-tuned Large Language Model for Radiological Impression Generation: A Multi-reader Performance Study. 2024;
56. Liu H, Tam D, Muqeeth M, et al. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *Advances in Neural Information Processing Systems*. 2022;35:1950-1965.
57. Hu EJ, Shen Y, Wallis P, et al. LoRA: Low-rank adaptation of large language models. *arXiv preprint arXiv:210609685*. 2021;
58. Dettmers T, Pagnoni A, Holtzman A, Zettlemoyer L. QLoRA: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*. 2024;36
59. Gendler M, Nadkarni G, Sudri K, et al. Large Language Models in Cardiology: A Systematic Review. *medRxiv*. 2024:2024--09.
60. Novak A, Rode F, Lisii. The Pulse of Artificial Intelligence in Cardiology: A Comprehensive Evaluation of State-of-the-art Large Language Models for Potential Use in Clinical Cardiology. *medRxiv*. 2023:2023--08.
61. Boonstra MJ, Weissenbacher D, Moore JH, Gonzalez-Hernandez G, Asselbergs FW. Artificial intelligence: revolutionizing cardiology with large language models. *European Heart Journal*. 2024;45(5):332--345.
62. Gui H, Omiye JA, Chang CT, Daneshjou R. The Promises and Perils of Foundation Models in Dermatology. *Journal of Investigative Dermatology*. 2024;
63. Ullah E, Parwani A, Baig MM, Singh R. Challenges and barriers of using large language models (LLM) such as ChatGPT for diagnostic medicine with a focus on digital pathology--a recent scoping review. *Diagnostic pathology*. 2024;19(1):43.
64. Shahab O, El Kurdi B, Shaukat A, Nadkarni G, Soroush A. Large language models: a primer and gastroenterology applications. *Therapeutic Advances in Gastroenterology*. 2024;17:17562848241227031.
65. Omar Sr M, Sharif Sr K, Glicksberg Sr BS, Nadkarni G, Klang Sr E. Emerging Applications of NLP and Large Language Models in Gastroenterology and Hepatology: A Systematic Review. *medRxiv*. 2024:2024--06.
66. Giuffre M, Kresevic S, Pugliese N, You K, Shung DL. Optimizing large language models in digestive disease: strategies and challenges to improve clinical outcomes. *Liver International*. 2024;
67. Mudrik A, Nadkarni GN, Efros O, et al. Exploring the role of Large Language Models (LLMs) in hematology: a systematic review of applications, benefits, and limitations. *medRxiv*. 2024:2024--04.
68. Barrit S, El Hadwe S, Carron R, Madsen JR. Rise of large language models in neurosurgery. *Journal of Neurosurgery*. 2024;1(aop):1--2.
69. Chiang C-C, Fries JA. Exploring the Potential of Large Language Models in Neurology, Using Neurologic Localization as an Example. *Neurology: Clinical Practice* 2024. p. e200311.

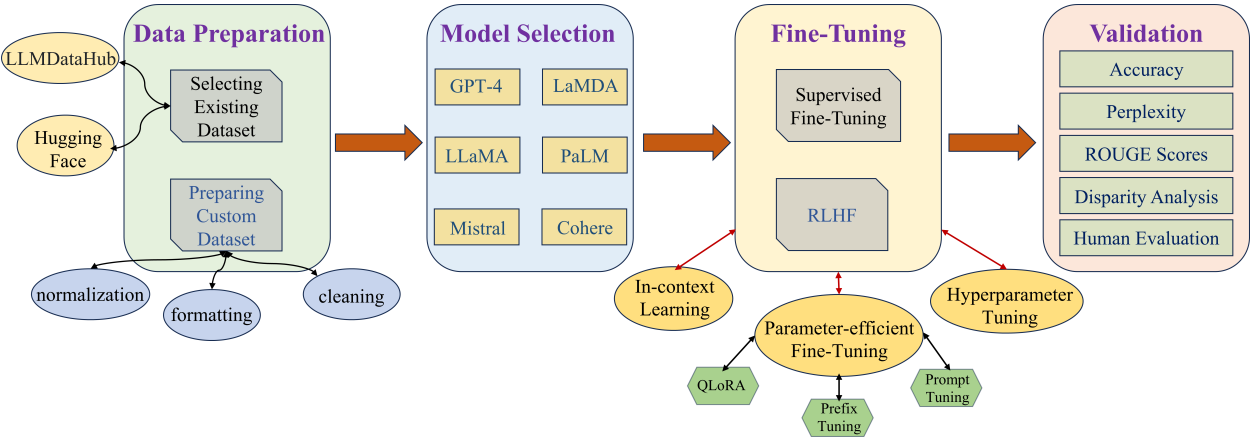
70. Romano MF, Shih LC, Paschalidis IC, Au R, Kolachalama VB. Large language models in neurology research and future practice. *Neurology*. 2023;101(23):1058--1067.
71. Bachmann M, Duta I, Mazey E, et al. Exploring the capabilities of ChatGPT in women's health: obstetrics and gynaecology. *npj Women's Health*. 2024;2(1):26.
72. Mudrik A, Tsur A, Nadkarni G, et al. Leveraging Large Language Models in Gynecologic Oncology: A Systematic Review of Current Applications and Challenges. *medRxiv*. 2024:2024--08.
73. Rydzewski NR, Dinakaran D, Zhao SG, et al. Comparative evaluation of LLMs in clinical oncology. *Nejm Ai*. 2024;1(5):Aloa2300151.
74. Lawson McLean A, Wu Y, Lawson McLean AC, Hristidis V. Large language models as decision aids in neuro-oncology: a review of shared decision-making applications. *Journal of Cancer Research and Clinical Oncology*. 2024;150(3):139.
75. Benary M, Wang XD, Schmidt M, et al. Leveraging large language models for decision support in personalized oncology. *JAMA Network Open*. 2023;6(11):e2343689--e2343689.
76. Luo M-J, Pang J, Bi S, et al. Development and evaluation of a retrieval-augmented large language model framework for ophthalmology. *JAMA ophthalmology*. 2024;142(9):798--805.
77. Chatterjee S, Bhattacharya M, Pal S, Lee S-S, Chakraborty C. ChatGPT and large language models in orthopedics: from education and surgery to research. *Journal of Experimental Orthopaedics*. 2023;10(1):128.
78. Sisk BA, Antes AL, DuBois JM. An Overarching Framework for the Ethics of Artificial Intelligence in Pediatrics. *JAMA pediatrics*. 2024;178(3):213--214.
79. Wyatt KD, Alexander N, Hills GD, et al. Making sense of artificial intelligence and large language models—including ChatGPT—in pediatric hematology/oncology. *Pediatric Blood & Cancer*. 2024;71(9):e31143.
80. Obradovich N, Khalsa SS, Khan WU, et al. Opportunities and risks of large language models in psychiatry. *NPP—Digital Psychiatry and Neuroscience*. 2024;2(1):8.
81. Volkmer S, Meyer-Lindenberg A, Schwarz E. Large Language Models in Psychiatry: Opportunities and Challenges. *Psychiatry Research*. 2024:116026.
82. Omar Sr M, Soffer Sr S, Charney Sr A, et al. Applications of Large Language Models in Psychiatry: A Systematic Review. *medRxiv*. 2024:2024--03.
83. Liu Z, Zhong A, Li Y, et al. Tailoring large language models to radiology: A preliminary approach to llm adaptation for a highly specialized domain. 2023;
84. D'Antonoli TA, Stanzione A, Bluethgen C, et al. Large language models in radiology: fundamentals, applications, ethical considerations, risks, and future directions. *Diagnostic and Interventional Radiology*. 2024;30(2):80.
85. Lee J, Sharma I, Arcaro N, et al. Automating surgical procedure extraction for society of surgeons adult cardiac surgery registry using pretrained language models. *JAMIA open*. 2024;7(3)
86. Oh N, Choi G-S, Lee WY. ChatGPT goes to the operating room: evaluating GPT-4 performance and its potential in surgical education and training in the era of large language models. *Annals of Surgical Treatment and Research*. 2023;104(5):269--273.
87. Adhikari K, Naik N, Hameed BMZ, Raghunath SK, Somani BK. Exploring the ethical, legal, and social implications of ChatGPT in urology. *Current Urology Reports*. 2024;25(1):1--8.
88. Gupta R, Pedraza AM, Gorin MA, Tewari AK. Defining the role of large language models in urologic care and research. *European Urology Oncology*. 2024;7(1):1--13.
89. Mukherjee S, Gamble P, Ausin MS, et al. Polaris: A Safety-focused LLM Constellation Architecture for Healthcare. *arXiv preprint arXiv:240313313*. 2024;
90. Zhao L, Zeng W, Shi X, et al. Aquilia-Med LLM: Pioneering Full-Process Open-Source Medical Language Models. *arXiv preprint arXiv:240612182*. 2024;

91. Li L, Zhou J, Gao Z, et al. A scoping review of using Large Language Models (LLMs) to investigate Electronic Health Records (EHRs). *arXiv preprint arXiv:240503066*. 2024;
92. Zhang X, Yan C, Yang Y, et al. Optimizing Large Language Models for Discharge Prediction: Best Practices in Leveraging Electronic Health Record Audit Logs. *medRxiv*. 2024:2024--09.
93. Cui H, Shen Z, Zhang J, et al. LLMs-based Few-Shot Disease Predictions using EHR: A Novel Approach Combining Predictive Agent Reasoning and Critical Agent Instruction. *arXiv preprint arXiv:240315464*. 2024;
94. Li R, Wang X, Yu H. LlamaCare: An Instruction Fine-Tuned Large Language Model for Clinical NLP. 2024;
95. Chao C-J, Banerjee I, Arsanjani R, et al. EchoGPT: A Large language model for echocardiography report summarization. *medRxiv*. 2024:2024.01. 18.24301503.
96. Guan Z, Wu Z, Liu Z, et al. CohortGPT: An enhanced gpt for participant recruitment in clinical study. *arXiv preprint arXiv:230711346*. 2023;
97. Demner-Fushman D, Kohli MD, Rosenman MB, et al. Preparing a collection of radiology examinations for distribution and retrieval. *Journal of the American Medical Informatics Association*. 2016;23(2):304-310.
98. Johnson AE, Pollard TJ, Berkowitz SJ, et al. MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*. 2019;6(1):317.
99. Zhou H, Gu B, Zou X, et al. A survey of large language models in medicine: Progress, application, and challenge. *arXiv preprint arXiv:231105112*. 2023;
100. Haltaufderheide J, Ranisch R. The ethics of ChatGPT in medicine and healthcare: a systematic review on Large Language Models (LLMs). *NPJ Digital Medicine*. 2024;7(1):183.
101. Pal A, Umapathi LK, Sankarasubbu M. Med-halt: Medical domain hallucination test for large language models. *arXiv preprint arXiv:230715343*. 2023;
102. Ong JCL, Chang SY-H, William W, et al. Ethical and regulatory challenges of large language models in medicine. *The Lancet Digital Health*. 2024;6(6):e428--e432.
103. Goh E, Bunning B, Khoong E, et al. ChatGPT Influence on Medical Decision-Making, Bias, and Equity: A Randomized Study of Clinicians Evaluating Clinical Vignettes. *Medrxiv*. 2023:2023.11.24.23298844.
104. Schmidgall S, Harris C, Essien I, et al. Addressing cognitive bias in medical language models. *arXiv preprint arXiv:240208113*. 2024;
105. Perez-Downes JC, Tseng AS, McConn KA, et al. Mitigating Bias in Clinical Machine Learning Models. *Current Treatment Options in Cardiovascular Medicine*. 2024:1-17.
106. Omar Sr M, Sorin Sr V, Apakama DU, et al. Evaluating and Addressing Demographic Disparities in Medical Large Language Models: A Systematic Review. *medRxiv*. 2024:2024--09.

Figure 1: A general workflow of LLM fine-tuning for specialized use cases. RLHF: Reinforcement Learning from Human Feedback; QLoRA: Quantized Low-Rank Adaptation.

Figure 2. Comparison of responses using different systems to the same question: "What are the interactions between propafenone and colchicine?".

Journal Pre-proof



Google Search Engine

Drugs.com
https://www.drugs.com › Interactions

Drug Interactions between colchicine and propafenone

Propafenone may increase the blood levels of colchicine to dangerous levels, increasing the risk of serious side effects that can affect your muscles, ...

People also ask

- Does colchicine interact with heart medication?
- What medications should not be taken with colchicine?
- What medications should not be taken with propafenone?
- Does colchicine interact with beta blockers?

National Institutes of Health (NIH) (.gov)
https://www.ncbi.nlm.nih.gov › articles › PMC9754312

Colchicine Drug Interaction Errors and Misunderstandings

by PD Hansten · 2023 · Cited by 13 — Propafenone did not appear to affect colchicine pharmacokinetics [30]. However, one cannot rule out that a larger daily dose of propafenone ...

Colchicine Toxicity · Mechanisms of Colchicine... · Risk Factors for Colchicine DDI

Drugs.com
https://www.drugs.com › Propafenone

Propafenone Interactions Checker

409 medications are known to interact with propafenone. Includes metoprolol, sertraline, albuterol.

DrugBank
https://go.drugbank.com › drugs

Propafenone: Uses, Interactions, Mechanism of Action

Propafenone may increase the arrhythmogenic activities of Ajmaline, Albendazole, The metabolism of Propafenone can be increased when combined with Albendazole.

Pharmacology · Interactions · Products · Categories

ChatGPT 4.0

Propafenone and colchicine can interact in a way that may increase the risk of side effects. Here's a breakdown of how these medications might affect each other:

- Increased Levels of Colchicine:** Propafenone can inhibit the enzyme CYP3A4 in the liver, which is responsible for metabolizing colchicine. When taken together, propafenone can increase the levels of colchicine in the blood, potentially leading to increased toxicity.
- Risk of Toxicity:** Elevated levels of colchicine can lead to serious side effects, including muscle damage, nerve damage, and gastrointestinal symptoms. In severe cases, colchicine toxicity can be life-threatening, particularly due to its effects on the muscles and nerves.

Because of these interactions, it's crucial for individuals who are prescribed both medications to be closely monitored by their healthcare provider. Adjustments in dosage or alternative medications may be necessary to minimize the risk of adverse effects. If you or someone you know is taking these medications together, it's important to consult with a healthcare professional to ensure safe and effective treatment.

OpenEvidence (LLM-Based RAG System):

Propafenone is metabolized by cytochrome P450 enzymes, including CYP2D6 and CYP3A4.^[1-2] Colchicine is a substrate for both CYP3A4 and P-glycoprotein (P-gp).^[3-4] The concomitant use of propafenone, which is a CYP3A4 substrate, with colchicine, which is a CYP3A4 and P-gp substrate, may potentially increase plasma levels of both drugs due to competitive inhibition of the metabolic pathway. This could increase the risk of adverse reactions associated with each drug, including propafenone-related proarrhythmia and colchicine toxicity, which can manifest as myopathy, rhabdomyolysis, and multiorgan failure.^[5-6]

The medical literature indicates that colchicine toxicity can be increased when co-prescribed with drugs that inhibit CYP3A4 and P-gp, and it is recommended to avoid such interactions or to monitor for toxicity.^[5] Given that propafenone is a CYP3A4 substrate, there is a potential for interaction, although the specific interaction between propafenone and colchicine is not explicitly detailed in the provided medical literature. Therefore, when considering co-prescription of these medications, it is prudent to closely monitor patients for signs of toxicity and consider alternative therapies if appropriate.

- Propafenone HCl. [Label via DailyMed](#). Food and Drug Administration (DailyMed)
- PROPAPENONE HYDROCHLORIDE. [Label via DailyMed](#). Food and Drug Administration (DailyMed)
- Colcrys. [Label via DailyMed](#). Food and Drug Administration (DailyMed)
- Gloperba. [Label via DailyMed](#). Food and Drug Administration (DailyMed)
- [Drugs That Interact With Colchicine via Inhibition of Cytochrome P450 3A4 and P-Glycoprotein: A Signal Detection Analysis Using a Database of Spontaneously Reported Adverse Events \(FAERS\)](#). Gómez-Lumbreras A, Boyce RD, Villa-Zapata L, et al. The Annals of Pharmacotherapy. 2023;57(10):1137-1146. doi:10.1177/10600280221148031. New Research
- [Colchicine Drug Interaction Errors and Misunderstandings: Recommendations for Improved Evidence-Based Management](#). Hansten PD, Tan MS, Horn JR, et al. Drug Safety. 2023;46(3):223-242. doi:10.1007/s40264-022-01265-1.