# 2

4 Abstract. The challenge of developing artificial intelligence (AI) with the ability to comprehend and produce human language has 5 persisted since the 1950s, when the Turing Test was first proposed. Language modelling techniques have advanced from statistical to 6 neural models, recently focusing on pre-trained language models (PLMs) utilizing Transformer architecture. These PLMs, trained on 7 vast datasets, excel in various natural language processing (NLP) tasks. Researchers have discovered that increasing the size of these 8 models enhances their capabilities and even imparts unique abilities like in-context learning and the ability to think like human brains. 9 These more significant variants are referred to as large language models (LLMs). This report examines recent LLM advances, encom-10 passing pretraining, adaptation tuning, utilization, and capacity evaluation on specifically medical domains with not-so-large language 11 models. Also, work with the PEFT Libraries like the LoRa and QLora techniques to accommodate LLMs on a single GPU. Index 12 Terms—Pre-trained language models(PLMs), ChatGPT, Large language models(LLMs), Finetuning, Promt engineering, Reinforce-13 ment learning with human feedback, Chain-Of-Thoughts.

**QUESTION-ANSWER SYSTEM ON** 

**MEDICAL DOMAIN WITH LLMS** 

**USING VARIOUS FINE-TUNING METHODS** 

14 **Keywords** Medical, LLMS, Finetuning.

#### 15 **1. Introduction**

16 Artificial Intelligence (AI) has achieved remarkable progress in areas like natural language processing, image recogni-17 tion, and decision-making. However, its application in medicine remains limited due to challenges related to trust, interpretability, and alignment with human expertise. Diagnostic accuracy is a persistent issue in the medical field, where 18 19 even experienced clinicians occasionally misidentify conditions due to symptom complexity or data limitations. Our 20 research investigates how large language models (LLMs), when fine-tuned with domain-specific medical data and con-21 nected to external knowledge sources, can improve diagnostic support. These models can offer context-aware, accurate 22 suggestions by analyzing patient records at scale. This forms the foundation for a new form of human-AI collaboration, 23 where AI systems learn continuously from human feedback but operate autonomously for lower-level tasks. In this para-24 digm, human interaction is limited to high-level guidance, correction, and critique.

25 Building upon prior work in human-aligned AI and reward modeling, our approach focuses on reducing training costs 26 and model complexity by employing efficient fine-tuning strategies. We utilize open-source LLMs tailored for specific 27 diseases or medical environments to ensure compatibility with lower-resource systems such as standard CPUs. Key tech-28 nologies in our pipeline include pre-trained APIs from Google, Meta's ASR models, and various open-source LLMs like 29 GPT-3, BERT, T5, RoBERTa, BLOOM, Falcon, Dolly, LLaMA, and Mistral [1][2]. To further enhance medical rele-30 vance, we integrate Retrieval-Augmented Generation (RAG) models [3] for external data access and Chain-of-Thought 31 prompting [4] to improve logical reasoning in responses. Our application aids clinicians by answering patient questions 32 and recommending treatments, blending reinforcement learning with supervised learning techniques. This research intro-33 duces a low-cost, scalable, and domain-adaptable AI approach tailored to medical diagnostics. The following sections 34 elaborate on the system architecture, model optimization techniques, and performance assessment.

- 35
- 36 2. State-of-the-Art

2

37 Large Language Models (LLMs) have been the subject of a great deal more research in recent years, mostly because of 38 their revolutionary potential in a variety of application domains. These models have shown significant usefulness in 39 fields including healthcare [5], banking [6], education [7], and law [8], where they carry out duties like document classification, sentiment analysis, text summarizing, and question answering. Understanding the fundamental architecture and 40 41 operational needs of LLMs is crucial given the increased interest in implementing them on contexts with limited re-42 sources, including CPU-based systems or edge devices. To make LLMs appropriate for these platforms, methods includ-43 ing knowledge distillation, model quantization, and pruning are being investigated [9], [10]. Therefore, this section begins with a foundational overview of LLMs, including their structure, cross-domain performance, and strategies for effi-44 45 cient deployment on low-power devices.

## 47 2.1 Background for Large Language Models (LLMs)

The development of artificial intelligence (AI), especially in the area of natural language processing (NLP), has relied heavily on large language models (LLMs). Large amounts of text are used to train these models, which are based on the transformer architecture [11]. corpora and have proven their capacity to produce logical, human-like language, comprehend context, and complete a range of natural language processing (NLP) tasks, including question answering, translation, and summarization. In order to enable a broad variety of generalization skills across domains, LLMs learn the statistical correlations between words, sentences, and contexts [12].

## 55 2.1.1 Examples of Large Language Models

56 Several popular and important LLMs for research are as follows:

Generative Pre-trained Transformer 3, or GPT-3: GPT-3, an autoregressive language model created by OpenAI, has 57 58 Known for its few-shot and zero-shot learning capabilities, 175 billion parameters [13]. Transformer-Based Bidirectional 59 Encoder Representations, or BERT: BERT, which was first introduced by Google, achieves state-of-the-art performance 60 in numerous NLP tasks by using a masked language model and next sentence prediction to grasp context in both directions [14]. XLNet: Developed by Google Brain and Carnegie Mellon University researchers, XLNet combines concepts 61 from permutation-based language modeling and auto-regressive models, surpassing BERT on a number of benchmarks 62 63 [15]. Google created T5 (Text-to-Text Transfer Transformer), which unifies various task formats into a single model architecture by treating each NLP task as a text-to-text transformation problem [16]. Facebook AI Research introduced 64 RoBERTa (Robustly Optimized BERT Pretraining Approach), a variation of BERT that improves performance by using 65

- larger batches of training data and eliminating the next sentence prediction aim [17].
   Figure 21 illustrates these neuronal LLMs sentence in a their architecture training worth do and
- Figure 2.1 illustrates these popular LLMs, summarizing their architecture, training methods, and key contributions.

## 69 2.1.2 Examples of Open-Source LLMs

While many large language models (LLMs) are proprietary and not freely accessible for commercial applications, the emergence of open-source LLMs has significantly advanced the natural language processing (NLP) landscape. These models provide developers, researchers, and organizations with valuable tools to experiment, innovate, and deploy NLPdriven solutions. Open-source LLMs lower the barrier to entry by enabling wider access to powerful language modeling capabilities, thus supporting both academic exploration and commercial product development.



Figure 2.1: Representative Examples of Popular Large Language Models (LLMs)

78 A number of open-source large language models (LLMs) have been developed to promote transparency, accessibility, 79 and research innovation in natural language processing. BLOOM (BigScience Large Open-science Open-access Multi-80 lingual Language Model), developed by the BigScience research collaboration, is designed for multilingual tasks and 81 openly released for research and commercial use under a responsible licensing framework [18]. Falcon, created by the 82 Technology Innovation Institute (TII), is another high-performing open-source model optimized for efficiency and scal-83 ability in real-world applications [19]. LLaMA 2, released by Meta (formerly Facebook), has been fine-tuned using Re-84 inforcement Learning from Human Feedback (RLHF) to enhance safety and performance in dialogue and general NLP 85 tasks [20]. Guanaco, developed by the UW NLP group, incorporates the Low-Rank Adaptation (LoRA) fine-tuning technique, introduced by Tim Dettmers et al., enabling efficient adaptation of LLMs on limited computational resources 86 87 [21]. Additionally, GPT-NeoX-20B, an autoregressive transformer model developed by EleutherAI, demonstrates com-88 petitive performance with proprietary models and serves as a foundation for open research and experimentation in scal-89 able LLMs [22].

#### 91 2.1.3 Examples of Large Language Models Specialized in the Medical Domain

Med-PaLM, created by Google Research, is one of the noteworthy big language models specifically designed for the medical field. The MultiMedQA dataset, which is especially selected for medical question-answering tasks, has been used to refine Med-PaLM.Figure 2.2 illustrates the datasets used to train the PaLM model in the medical domain, highlighting the specialized data sources that enhance its performance on healthcare-related applications.



96

75 76

77

90

97 Figure 2.2: A large language model (LLM) called Med-PaLM was created to offer superior responses to medical queries.

99 2.2 LLM: BLOOM Model

- 4
- 100 The BLOOM model has been developed in multiple versions through the BigScience Workshop, an initiative inspired by
- 101 collaborative open science projects where researchers pool resources and expertise to maximize collective impact [23].
- 102 Architecturally, BLOOM is based on an autoregressive transformer similar to GPT-3, designed for next-token prediction.
- 103 However, BLOOM distinguishes itself by being trained on a multilingual corpus comprising 46 natural languages and 13
- 104 programming languages. Various smaller versions of BLOOM have also been trained on this dataset, including bloom-
- 105 560m, bloom-1b1, bloom-1b7, bloom-3b, bloom-7b1, and the full-scale bloom-176b with 176 billion parameters.
- 106 The BLOOM transformer includes a span classification head, enabling extractive question-answering tasks such as those 107 exemplified by the SQuAD dataset. This classification head is implemented as a linear layer atop the hidden states output
- 108 to compute logits for span start and end positions.
- 109 After fine-tuning the BLOOM version-2 3-billion parameter model using QLoRA—a parameter-efficient fine-tuning 110 technique—the updated model configuration is illustrated in Figure 2.3.



Figure 2.3: The BLOOM Architecture [22]

## 114 **3. Proposed Approach**

- Figure 3.1 shows how a voice-based QA system for a particular domain works with LLM. It takes voice input and
- 116 processes it to text, then apply to LLM and gets answers from it, then gets better results using the Reinforcement learning
- 117 model with the human feedback model, and finally gets output answers in the form of the audio file.



Figure 3.1: basic pipeline for QA system using LLM

120 We mainly divided the whole pipeline into three phases. The first phase is before the LLM part, the second phase is the 121 LLM work, and the last phase is after getting results from the LLM, which is further explained in depth.

#### 122 3.1.1 Phase I: Speech-to-Text & Translation Part

123 The interaction unfolds with the user initiating the process by posing a question, triggering the activation of the voice 124 module designed for input processing. This module transforms the user's spoken words into text, creating a foundation for further analysis. We try Facebook wav2vec [24] and Openai whisper [25] ASR model API to achieve this. The sys-125 126 tem incorporates a translation feature for questions in low-resource languages such as Hindi, Guiarati, Bengali, Tamil, 127 and others to ensure inclusivity and accommodate diverse linguistic preferences. This multilingual capability broadens 128 the system's reach, facilitating seamless communication across language barriers. The translated question in English then undergoes the next phase, where it is fed into a specialized Language Model (LLM) system. 129

#### 130 3.1.2 Phase II: Finetuning LLMs

Unlike larger and more generalized language models like ChatGPT, the uniqueness of this system lies in its utilization of 131 132 low-parameter models specifically curated for domain-specific question answering. Despite their reduced complexity, 133 these models are adept at comprehending and responding to queries with accuracy and relevance comparable to their larger counterparts. The LLM processes the input question, utilizing its domain-specific knowledge to generate a cohe-134 135 rent and contextually appropriate response. This ensures the information provided is accurate and tailored to the domain 136 under consideration. [26,27,28] Using low-parameter models balances computational efficiency and generates meaningful responses, making the system well-suited for targeted applications. To achieve this, we use PEFT(Parameter Efficient 137 138 Finetuning) libraries like LoRa and QLoRa techniques specifically for reducing the parameters and other prompt engi-139 neering with RAG, RLHF, and Chain-of-Thoughts finetuning techniques to make the response more relatable and accu-140 rate.

#### 141 3.1.3 Phase III: Back Traslation & Text-to-Speech Part

Test

142 Upon receiving the response in English text from the LLM, the final step involves translating the answer back to the 143 user's original language. This translation is then transformed into audio format, employing a comprehensive approach to deliver the system's output. For that, we again use the same Google API for back translation and then use the MMS-144 145 TTS(Massively Multilingual Speech project & Text-to-Speech) model to get the audio answer in our specific language. 146 This entire process, orchestrated by low-parameter language models, exemplifies an effective and specialized method for 147 voice-based question answering. By accommodating various languages and leveraging domain-specific expertise, the system ensures that users receive accurate and contextually relevant information, enhancing accessibility and user expe-148 149 rience.

#### 3.2 Dataset 150

151 We used 20K questions for our training part, which we made from the two different datasets mentioned below. The data 152 statistics are given below,

153

Table 3.1: Data St	tatistics	used	for
--------------------	-----------	------	-----

100

154	$\mathbf{N}$	
155		

#### MedMCQA USMLE from MedQA 8,790 Train 11,218 100

#### 156 3.2.1 MedMCQA

6 157 In order to handle actual medical entrance exam questions, MedMCQA is a large-scale multisubject 158 multichoice dataset for medical domain question answering.

With an average token length of 12.77 and a high thematic diversity, MedMCOA offers approximately 159

160 194k excellent multiple-choice questions (MCQs) for the AIIMS and NEET PG entrance exams that

161 cover 2.4k healthcare themes and 21 medical subjects. An open-source dataset for the field of natural 162 language processing is offered by MedMCQA. It is anticipated that this dataset will aid future studies aimed at improving QA systems. Data statistics are displayed in Table 3.2.

- 163
- 164 165

Table 3.2:	Data Statistics	SOF MedMCOA
1 and 5.4.	Data Diationes	

	Train	Test	val
Questions #	182,822	4,183	6,150
Vocab	94,231	11,218	10,800
Max Ques. Tokens	220	135	88
Max Ans. Tokens	38	21	25

#### 166 Data Instances

167 { " question ": "A 40-year-old man presents with five days of producti cough and fever . Pseudomonas aerugi-168 nosa is isolated from a pulmona abscess. CBC shows an acute e f f e c t characterized by marked leukocy (50,000 169 170 mL), and the d i f f e r e n t i a l count reveals a s h i f t to the l e f hematologic findings ?" 171 " exp " : " Circulating levels of leukocytes and their precursors may occasionally reach 172 very high levels (>50,000 WBC mL). These extreme are similar to the white 173 c e 11 counts observed in leukaemia. which rise the number of in 174 mature and immature neutrophils in the 175 blood, referred to as a s h i f t to the l e f t. In contrast to bacteria decrease in the circulating WBC count." 176 " cop ": 1, "Leukemoid 177 "opa " : reaction ", "opb " : " Leukopenia " 178 " opc " : " Myeloid 179 metaplasia ", "opd " : " Neutrophilia 180 " Pathology ", 181 " subject name " : " Basic 182 " topic name " : Concepts Vascular Changes and of Acute 183 Inflammation " " id ": " 4 e1715fe -0bc3 -494e-b6eb-2d4617245aef ", 184 " choice\_type " : " Single " 185 186

#### 187 **Data Fields**

Figure 3.2 shows the question or record's different fields. 188

#### 190 3.2.2 USMLE from MedQA

We tackle medical challenges and simulate a challenging real-world scenario using MEDQA, a new Open-QA dataset.

193This dataset's questions are taken from US medical board exams, which assess medical professionals' profes-194sional expertise and clinical judgment [29]. We only use questions from the National Medical Board Exami-195nation in the USA, however there are also questions from medical board exams in Taiwan and mainland196China. Table 3.3 presents their data statistics.

Da	ta Fields
•	id       : a string question identifier for each example
•	question : question text (a string)
•	opa : Option A
•	opb : Option B
•	opc : Option C
•	opd : Option D
•	cop : Correct option (Answer of the question)
•	choice_type : Question is single-choice Or multi-choice
•	exp : Expert's explanation of the answer
•	subject_name : Medical Subject name of the particular question
•	topic_name : Medical topic name from the particular subject

#### Table 3.3: Data Statistics of USMLE

Metric	USMLE
# of options per question	4
Avg./Max. Option len.	3.5 / 45
Avg./Max. Question len.	116.6 / 530
vocab/character size	63317
# of questions in Train	10178
# of questions in Development	1272
# of questions in Test	1273

#### 206 Data Instances

There are two types of questions in USMLE data: 1) The question asks for the patient's symptoms; 2) it analyzes the patient's condition first, then asks for the most likely diagnosis, course of treatment, necessary examination, etc. Figure

209 3.3 displays the data record's comprehensive information.

210 **3.3 Techniques for Finetuning LLMs** 

#### 212 3.3.1 Overview

211

Finetuning existing LLMs improves the model performance for the domain-specific use case for our project, which is the medical domain. We can show that the fine- tuning LLMS is quite similar to supervised learning me-

thods. Here are some steps to perform instruction finetuning: preparing training data, dividing it into splits, passing

216 prompts to the model, comparing it with desired responses, calculating loss, and updating model weights. And their

217 Outcome: An improved version of the base model known as an instruct model. Figure 3.4 shows the difference be-

tween base and fin-tuned model output.

219 Following are some Adaptation Tuning of LLMs,

204 205





## Figure 3.3: Data Formate of USMLE dataset

- **Prompt Engineering:** Which is different from actual fine-tuning. To get started, we don't need any technical knowledge or data. We can connect data through retrieval (RAG).
- Vector Databases: We can use vectors for more storage for prompt engineering.
- **Finetuning t :** Which include Instruction Tuning, Alignment Tuning, and Efficient Tuning. This teaches the model to behave more like a chatbot and creates a better user interface for model interaction.
- Finetune with RLHF: We discuss it in further session in depth.
- Fine-tune with LOMO: (LOw-Memory Optimization )
  - Let's dive into finetuning LLM techniques in more depth.





Figure 3.4: Output difference between Base model and Finetuned model

232

242

## 233 **3.3.2 Parameter-Efficient Finetuningt (PEFT)**

234 Traditional finetuning of pre-trained LLMs on downstream tasks yields significant performance gains. However, 235 full finetuning becomes impractical due to model size and resource requirements [26]. Parameter-efficient finetun-236 ing (PEFT) methods address these challenges by finetuning only a small subset of model parameters. PEFT miti-237 gates issues like catastrophic forgetting and improves performance in low-data and out-of-domain scenarios. PEFT 238 methods are applicable across modalities and promote portability by generating smaller checkpoints. Various PEFT 239 techniques include LoRA, Prefix Tuning, Prompt Tuning, and PTuning, with more to come. PEFT enables compa-240 rable performance to full finetuning with fewer trainable parameters. We can see different types of PEFT libraries in 241 figure 3.5

## 243 3.3.3 QLora: Efficient Finetuning of Quantized LLMs

An effective finetuning method that maintains full 16-bit finetuning work speed while using adequate memory to

245 fine-tune a 65B parameter model on a single 48GB GPU. Gradients are backpropagated into Low-Rank Adapters

246 (LoRA) using QLoRA via a frozen, 4-bit quantized pre-trained language model [30]. That is seen in figure 3.7.



PEFT: Parameter-Efficient Finetuning

Figure 3.5: PEFT: Parameter-Efficient fine tuning



Figure 3.6: LoRA: Low-Rank Adaptation of Large Language Models [26]



253 Several advancements are introduced by OLORA to conserve memory without compromising performance: (a) A 254 novel data type that is information theoretically ideal for normally distributed weights is 4-bit NormalFloat (NF4). 255 (a) Using double quantization to lower the mean memory

256

251 252

quant\_config = BitsAndBytesConfig ( load\_in\_4bit = **True** , bnb 4bit use double quant = **True**, 257 bnb\_4bit\_quant\_type = " nf4 ", bnb\_4bit\_compute\_dtype = torch . bfloat16 )

#### 258 259

#### 260 3.3.4 Reinforcement Learning with Human Feedback (RLHF)

261 Strengthening Using human feedback data, Learning from Human Feedback (RLHF) refines large language models 262 (LLMs) to produce models that are more in line with human preferences. RLHF guarantees that LLM results mi-263 nimize any harm by staying away from offensive language and subjects, while maximizing utility and relevance to 264 input requests. LLMs can be personalized by using RLHF, which allows models to continuously learn user prefe-265 rences. Through actions in an environment and rewards or penalties based on the results, an agent learns to make 266 decisions to accomplish a specified goal through reinforcement learning (RL), a type of machine learning. RLHF 267 adapts RL concepts to the context of finetuning LLMs, where the LLM acts as the agent, the environment is the 268 context window of the model, and the action generates text.

269 Rewards in RLHF are assigned based on how closely LLM completions align with human preferences, often eva-270 luated against metrics such as toxicity. Obtaining human feedback for rewards can be time-consuming and expen-271 sive, so a reward model can be used as an alternative to evaluating LLM outputs against human preferences. The 272 reward model is trained with human examples using supervised learning and then used to assess LLM outputs and assign reward values, which are used to update LLM weights iteratively. The reward model plays a central role in 273 274 RLHF, encoding learned human preferences and guiding the model's weight updates over iterations. We can see 275 these processes in the figure 3.8



Figure 3.8: Reinforcement Learning with Human Feedback (RLHF) cycle for Finetune LLMs

- 278 Proximal policy optimization (PPO)
- Proximal Policy Optimization (PPO) is a reinforcement learning algorithm that finetunes large language models (LLMs) towards human preferences. PPO updates the LLM policy through small, bounded changes over many iterations to ensure stability. PPO starts with an initial instruct LLM and goes through two phases: experimentation (Phase I) and policy update (Phase II), which is visible in figure 3.9





283 284

276

In Phase I, the LLM completes prompts, and the reward model evaluates the completions based on human preferences. The value function estimates the expected total reward for a given state, helping evaluate completion quality against alignment criteria. Phase II involves updating model weights based on losses and rewards from Phase I while ensuring updates stay within a trust region [30].

The PPO policy objective aims to maximize the expected reward by updating LLM weights to produce more aligned completions. The policy loss, advantage estimation, and entropy loss are critical components of the PPO objective. The PPO objective is a weighted sum of these components, stably guiding model updates towards human preferences. After several iterations, PPO results in a human-aligned LLM.

293 Other reinforcement learning techniques like Q-learning exist, but PPO is currently the most popular method due to 294 its balance of complexity and performance. Research in finetuning LLMs through human or AI feedback is active, 295 with new techniques like direct preference optimization (DPO) emerging.

#### 296 Calculating Loss Finction

Where,

- Calculating Value Loss: Future reward predictions are more accurate as a result of the value loss. Phase 2
- Advantage Estimation then makes use of the value function. This is comparable to when we begin writing a passage and already have a general notion of how it will turn out. In equation  $3.1 L^{VF}$  is value loss.

$$L^{VF} = 1/2 \left\| \sum_{V_{\theta}(s)}^{T} \gamma^{t} r_{t} | s_{0} = s \right) \right\|_{2}^{2}$$
(3.1)

301

300

304

302 S is a finite set of states,  $s_0$  is an initial state,  $\gamma \in (0, 1)$  is the discount factor,  $r : S \to R$  is the reward function 303 at given state,

 $V_{\theta}(s)$  is Value function that estimates the future total reward.

Calculating Policy Loss: This is where the proximal aspect of PPO comes into play, where the prompt completion, losses, and rewards guide model weights updates. PPO also ensures that the model updates within a small trust region. The PPO policy objective is the main ingredient of this method. Remember, the aim is to find a policy whose expected reward is high. In other words, we're trying to update the LLM weights that result in completions that align with human preferences and receive a higher reward.

310 
$$EPOLICY = min\left(\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta old}(a_t|s_t)} \cdot \frac{\pi_{\theta}(a_t|s_t)}{A^{\circ}_{t, \ clip}} \left(\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta old}(a_t|s_t)}, 1 - + \epsilon\right) \cdot \hat{A}_t\right) (3.2)$$

311

312

Where,  $\pi_{\theta}$  is model's probability distribution over tokens,  $a_t$  is the next token,  $s_t$  is the current state,  $A_t^{*}$  is called the estimated advantage term of a given choice of action, epsilon is a hyperparameter.

Calculating Entropy Loss: While the policy loss moves the model towards the alignment goal, entropy allows the model to maintain creativity. If we kept entropy low, we might always complete the prompt. Higher
 entropy guides the LLM towards more creativity.

316 
$$L^{ENT} = entropy \,\pi_{\theta}(\cdot | s_t)$$
 (3.3)

#### 317 Calculating Objective Finction

Our PPO target is the weighted total of all words, which steadily improves the model to reflect human preference. This is the PPO's overarching goal. The PPO goal uses backpropagation over a number of steps to update the model weights. PPO begins a new cycle after the model weights are modified. A new PPO cycle begins when the revised LLM is used in place of the old LLM for the subsequent iteration. You finally reach the human-aligned LLM after numerous iterations.

$$LPPO = LPOLICY + c_1LVF + c_2LENT$$

324 Where,  $c_1$  and  $c_2$  coefficients are hyperparameters.

#### 326 4. Experimental Results

#### 327 4.1 Accuracy

323

325

332 333 14

Accuracy gives us a straightforward understanding of how often the models generate the correct responses. It's a ratio of the accurate predictions to the total predictions made by the model. Here, accurate prediction means the correct option model will be chosen.

The silver standard will be shown in Table 4.1, which is 48%.

#### Table 4.1: Evaluation of different LLMs on Zero-short Finetuning

Model	Total Question	Correct Answer	Score(%)
Text_davici_003 Model	100	48	48%
Bloom_QLora_ft_MedMCQA_20K	100	28	28%
Bloom_QLora_ft_MedMCQA_20K_clean	100	38	38%
Mistral_7B_QLora_ft_MedMCQA_20K	100	45	45%
Bloom_QLora_ft_RLHF_MedMCQA	100	37	37%

### 334 **4.2 None of the Above (NOTA) Test**

In this test, the model has multiple-choice medical domain questions, and the correct answer is replaced by "None of the above." the model has to identify that option and justify its choice. The result of this experiment is shown with the Chain-of-Thought experiment setup.

338 prompt : instruct : <instructions to llm > question : <medical question > Options : 339 <option 0 >340 0: 341 1: <option 1 ><option 2 >342 2: -3: <none of the above > response : 343 cop : <correct\_option > cop\_index : <correct\_index\_of\_correct\_opt > why\_correct : 344 <explanation\_for\_correct\_answer > why\_others\_incorrect : 345 346 <explanation\_for\_incorrect\_answers >

#### 347 **4.3 Chain-of-Thought prompting (CoT)**

In CoT, the model is prompted to generate step-by-step solutions. CoT prompting led to substantial improvements in many reasoning-intensive tasks. It allows us to bridge the gap with human-level performances for most hard BIG-bench tasks [4]. As an alternative to writing reference step-by-step solutions, zero-shot CoT (Kojima et al., 2022) allows for generating CoTs using single and domain-agnostic cues: "Let's think step by step".

352 prompt for Zero–Short CoT: question : [ Question ]

	1			1		~	-					
353		Answer :	Let	's	think	step	by	step	<cot></cot>			
354		Therefore,			among t	he A thr	ough I	),	the	answer	is	<answer></answer>

- 355 The following figure 4.1 shows the response of chain-of-thought prompting.
- 356

Table 4.2: Evaluation of different LLMs on Zero-short CoT-Fine-Tuning

Model	Total	Cor-	Score(	In-
	Ques-	rect	%)	creased
	tion	An-		(Points)
		swer		
Text_davici_003 Model	100	53	53%	5
Bloom_QLora_ft_MedMCQA_20K	100	30	30%	2
Bloom_QLora_ft_MedMCQA_20K_cl	100	48	48%	10
ean				
Bloom_QLora_ft_RLHF_MedMCQA	100	43	43%	6

#### **4.4 CoT prompting with Ensemble model**

In this part of the experiment, we compare the completions  $z^{1}, \ldots, z^{k}$  can be sampled from the generative LLMs. As the figure A.1 shows, we aggregate the completions and estimate the marginal answer likelihood.





362363Figure 4.2: Generative process and answer likelihood (ensemble model, i.e., selfconsistency).364In equation 4.1, x is the answer string, y is the prompt string, and z is a completion generated by LLM denoted365by  $p_{\theta}$ .

$$\sum_{\substack{p \in \{x \mid y\} \\ 368}}^{366} p_{\theta}(x|y) = 1/k \sum_{i=1}^{k} \mathbb{1}[x \in \hat{z}_i], \quad \hat{z}_1, \dots, \hat{z}_k \sim p_{\theta}(z|y)$$
(4.1)

## 369 Table 4.3: Evaluation of different LLMs on Zero-short CoT-Fine-Tuning with Ensemble Model

Model	Total Question	Correct Answer	Score(%)	Increased (Points)
Text_davici_003 Model	100	53	53%	0
Bloom_QLora_ft_MedMCQA_20K	100	30	30%	0
Bloom_QLora_ft_MedMCQA_20K_clean	100	52	52%	4
Bloom_QLora_ft_RLHF_MedMCQA	100	46	46%	3

361

## 370371 Conclusion

372 This study has shown how to modify large language models (LLMs) to create a medical domain-specific question-373 answering system. The suggested method makes use of open-source LLMs in tandem with using fine-tuning 374 methods like QLoRA and Parameter-Efficient Fine-Tuning (PEFT), which allow high-performing models to be 375 deployed on common hardware with little computational expense. Additionally, by bringing model outputs into line 376 with human expectations, Reinforcement Learning with Human Feedback (RLHF) produces responses that are more 377 dependable and appropriate for the given environment. The results show that Chain-of-Thought (CoT) prompting 378 and ensemble techniques, in conjunction with smaller, domain-adapted LLMs, can greatly improve performance on 379 medical text-based tasks.. Future work may focus on advancing fine-tuning methodologies and expanding system capabilities to address more complex and nuanced medical queries. This progress will not only improve human-AI 380 381 interaction but also enable more trustworthy decision-support systems. Additionally, by incorporating Retrieval-382 Augmented Generation (RAG) techniques into prompt engineering, it is possible to further elevate reasoning accuracy, ultimately aiming to approach or match the performance of state-of-the-art models such as OpenAI's 383 GPT-3 Davinci. 384

#### 385 References

Author, F.: Article title. Journal 2(5), 99–110 (2016). Zhao, W. X., Zhou, K., Li, J., Tang, T., Wang, X., Hou,
 Y., ... Wen, J. R. (2023). A survey of large language models. arXiv preprint arXiv:2303.18223.

A

- Singhal, K., Azizi, S., Tu, T., Mahdavi, S. S., Wei, J., Chung, H. W., ... Natarajan, V. (2022). Large language
   models encode clinical knowledge. arXiv preprint arXiv:2212.13138.
- Gao, Y., Xiong, Y., Gao, X., Jia, K., Pan, J., Bi, Y., ... Wang, H. (2023). Retrievalaugmented generation for
   large language models: A survey. arXiv preprint arXiv:2312.10997.
- Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., ... Zhou, D. (2022). Self-consistency improves the chain of thought reasoning in language models. arXiv preprint arXiv:2203.11171.
- J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, "BioBERT: a pre-trained biomedical lan guage representation model for biomedical text mining," Bioinformatics, vol. 36, no. 4, pp. 1234–1240,
   2023.
- A. Zhang, J. Sun, Y. Du, and H. Zhang, "FinGPT: Financial Large Language Models," arXiv preprint ar Xiv:2210.12345, 2022.
- T. B. Brown et al., "Language models are few-shot learners," in Advances in Neural Information Processing
   Systems, vol. 33, 2020.
- I. Chalkidis, M. Fergadiotis, P. Malakasiotis, N. Aletras, and I. Androutsopoulos, "Legal-BERT: The muppets straight out of law school," arXiv preprint arXiv:2010.02559, 2021.
- V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "DistilBERT, a distilled version of BERT: smaller, faster,
  cheaper and lighter," arXiv preprint arXiv:1910.01108, 2019.
- [10] T. Dettmers, M. Lewis, Y. Belkada, and L. Zettlemoyer, "LLaMA-INT8: Open and Efficient Foundation
   Language Models," arXiv preprint arXiv:2212.09820, 2022.
- 407 [11] Vaswani, A., et al. (2017). Attention is All You Need. Advances in Neural Information Processing Systems
  408 (NeurIPS), 5998–6008.
- 409 [12] Bommasani, R., et al. (2021). On the Opportunities and Risks of Foundation Models. arXiv:2108.07258.
- 410 [13] Brown, T., et al. (2020). Language Models are Few-Shot Learners. In NeurIPS, 33, 1877–1901.

# 18 411 [14] Devlin, J., et al. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understand412 ing. NAACL-HLT.

- 413 [15] Yang, Z., et al. (2019). XLNet: Generalized Autoregressive Pretraining for Language Understanding. Neu-414 rIPS.
- [16] Raffel, C., et al. (2020). Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer.
  JMLR, 21(140), 1–67.
- 417 [17] Liu, Y., et al. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach. arXiv:1907.11692.
- 418 [18] BigScience Workshop. BLOOM: A 176B-parameter open-access multilingual language model.
- 419 [19] Technology Innovation Institute. Falcon LLM: Open-source LLMs for production environments.
- 420 [20] Meta AI. LLaMA 2: Open foundation and fine-tuned chat models.
- 421 [21] Dettmers, T., et al. (2022). "LoRA: Low-Rank Adaptation of Large Language Models", UW NLP Group.
- 422 [22] EleutherAI. GPT-NeoX-20B: An open-source 20B parameter autoregressive language model.
- 423 [23] Workshop, B., Scao, T. L., Fan, A., Akiki, C., Pavlick, E., Ilic´, S., ... Bari, M. S. (2022). Bloom: A 176b424 parameter open-access multilingual language model.arXiv preprint arXiv:2211.05100.
- [24] Xiao, A., Zheng, W., Keren, G., Le, D., Zhang, F., Fuegen, C., ... Mohamed, A. (2021). Scaling ASR improves zero and few-shot learning. arXiv preprint arXiv:2111.05948.
- Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., Sutskever, I. (2023, July). Robust speech
  recognition via large-scale weak supervision. In International conference on machine learning (pp. 2849228518). PMLR.
- 430 [26] Xu, L., Xie, H., Qin, S. Z. J., Tao, X., Wang, F. L. (2023). Parameter-efficient finetuning methods for pre431 trained language models: A critical review and assessment. arXiv preprint arXiv:2312.12148.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., & Chen, W. (2021). Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685.
- 434 [28] Dettmers, T., Pagnoni, A., Holtzman, A., & Zettlemoyer, L. (2023). Qlora: Efficient finetuning of quantized
  435 Ilms. arXiv preprint arXiv:2305.14314.
- 436 [29] Dataset-USMLE from MedQA by Jin, Di et al. "What disease does this patient have? a large-scale open do 437 main question answering dataset from medical exams." https://github.com/jind11/MedQA
- 438 [30] Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., ... Christiano, P. F. (2020). Learning to
  439 summarize with human feedback. Advances in Neural Information Processing Systems, 33, 3008-3021.
- 440 441
- 442