# LLM-Medical-Finetuning for Medical Question Answering
## Generative AI Project

Muhammad Iqrash Qureshi (22i-1174), Haider Zia (22i-1196), and Azhar Iqbal (21i-2508)

FAST-National University of Computer and Emerging Sciences, Islamabad, Pakistan
{i221174, i221196, i212508}@nu.edu.pk

**Abstract.** Large Language Models (LLMs) such as LLaMA and Mistral have shown outstanding performance across general-purpose natural language tasks. However, their direct application to specialized medical question-answering remains limited due to domain-specific terminology, scarce instruction datasets, and computational constraints associated with full fine-tuning. In this project, we investigate parameter-efficient fine-tuning (PEFT) methods—specifically LoRA and QLoRA—to adapt two state-of-the-art models, LLaMA-3-8B-Instruct and Mistral-7B-Instruct, for medical QA tasks.

We fine-tune these models on medical instruction datasets including MedQuAD and Medical Meadow WikiDoc, creating domain-aligned instruction pairs suitable for supervised fine-tuning. Our goal is to evaluate how efficiently PEFT techniques enable small-to-mid-scale LLMs to specialize in the medical domain while remaining computationally feasible on a single T4 GPU environment. The resulting models are compared quantitatively and qualitatively to measure improvements in medical coherence, factual accuracy, and response reliability.

**Keywords:** Medical AI, LLM Fine-Tuning, LoRA, PEFT, LLaMA, Mistral, Medical Question Answering

## 1 Introduction

Recent advancements in Large Language Models (LLMs) have significantly transformed the landscape of natural language understanding, reasoning, and generation. These models have demonstrated strong capabilities across diverse domains; however, applying them effectively within specialized fields such as medicine remains a considerable challenge. Medical question answering (Medical QA) is a critical task that involves interpreting complex clinical language, retrieving relevant medical knowledge, and generating accurate and safe responses. Due to the sensitivity and high-stakes nature of healthcare, it is essential that models used in this domain are reliable, domain-adapted, and capable of maintaining factual correctness.

General-purpose LLMs such as LLaMA and Mistral possess strong linguistic and reasoning abilities but lack domain-specific medical knowledge required for tasks like disease diagnosis queries, treatment guidelines, drug interactions, and patient education. This gap can be addressed through domain-specific fine-tuning, where a general LLM is trained on high-quality medical question-answer pairs to adapt it to medical terminology, reasoning patterns, and clinical context. Fine-tuning on curated datasets enhances a model's ability to respond with increased precision, medical correctness, and contextual understanding.

In this work, we fine-tune two state-of-the-art open-source LLMs—**LLaMA-8B** and **Mistral-7B**—on medical QA datasets including the Medical Meadow WikiDoc dataset and the MedQuAD dataset. These datasets encompass professional medical knowledge, patient-facing information, and structured question-answer pairs sourced from authoritative medical repositories. By fine-tuning on these data sources, we aim to build domain-specialized models capable of generating medically relevant, contextually accurate, and logically coherent responses.

The contributions of this work are threefold:

- We construct a fine-tuned version of LLaMA-8B and Mistral-7B for the Medical QA domain using high-quality, diverse datasets.
- We perform extensive evaluation and comparative analysis between the two models to examine differences in accuracy, reasoning ability, hallucination tendency, and robustness across question types.
- We benchmark our results against state-of-the-art models and discuss the limitations, challenges, and safety concerns of deploying LLMs in medical applications.

## 2   Related Work

Research on medical question answering (QA) has evolved over several decades, and the approaches used today look quite different from the earlier systems. Much of the early work relied on fairly rigid rule-based pipelines or retrieval-driven methods, where the system mainly searched through medical documents and returned the most relevant snippets. One example of this early line of work is PubMedQA [1], which provided one of the first structured benchmarks for biomedical QA and helped establish common evaluation practices in the field.

As deep learning became more widely adopted, the focus shifted from hand-crafted rules to neural language models trained on large corpora. Models like BioBERT [?], which was specifically pre-trained on biomedical literature, showed that adapting a general architecture to domain-specific text could produce sizeable improvements. Later studies used transfer learning techniques to push this idea further, reporting gains on datasets such as MedMCQA [2] and emrQA [3].

More recently, research attention has moved toward large language models (LLMs) and instruction-tuned systems for medical QA. Projects such as MedAlpaca [4] demonstrated that fine-tuning general LLMs (e.g., LLaMA-based models) on curated medical instructions improves their ability to answer clinical questions in a more reliable manner. BioInstruct [5] later showed that the design

of the instruction dataset itself can substantially influence reasoning quality and factual accuracy. Alongside these trends, parameter-efficient fine-tuning (PEFT) techniques like LoRA and QLoRA [6,?] emerged as practical tools for adapting large models to narrow domains without having to retrain the entire network.

Progress in the field has also been supported by the release of several specialized datasets. MedQuAD [7] is one of the more commonly used sources, containing question-answer pairs compiled from authoritative health websites. Other datasets, such as MedExQA [8] and K-COMP [9], extend coverage to multiple-choice formats and retrieval-augmented settings, enabling more fine-grained evaluation of different reasoning capabilities.

A number of comparative studies have attempted to benchmark LLM performance in medical reasoning tasks. Singhal et al. [10] reported that models such as MedPaLM2 could reach performance levels approaching medical professionals on standardized exam-style assessments. Complementary to this, frameworks like SafetyMed [11] examine not only accuracy but also reliability, potential harm, and robustness under distribution shifts—an important consideration for clinical use. Other works, including BioMedLM [12] and SM70 [13], illustrate how domain-focused pretraining continues to shape performance for specialized biomedical tasks.

Despite steady improvements, several challenges persist. Models often struggle with rare or nuanced medical conditions, and hallucination remains a recurring issue, as highlighted in multiple evaluations [3,?]. This has led to increased emphasis on safety-aware training methods, robustness analysis, and better domain adaptation strategies [11].

Overall, the existing literature shows that tailoring LLMs to the medical domain—particularly through focused datasets and PEFT-based fine-tuning—can substantially improve their usefulness. However, achieving consistent, reliable performance in real clinical environments remains an open problem and continues to motivate ongoing research [14].

## 3    Data Acquisition and Preprocessing

To fine-tune the LLaMA-3-8B-Instruct and Mistral-7B-Instruct models for the medical question–answering task, we rely on two high-quality open-source datasets: **Medical Meadow WikiDoc** and **MedQuAD**. Both datasets provide medically grounded question–answer pairs and domain-specific knowledge, making them ideal for supervised instruction fine-tuning using LoRA.

This section explains the data sources, acquisition steps, storage format, preprocessing pipeline, and dataset statistics. Visual summaries of the dataset distribution and question types are included for clarity.

### 3.1    Medical Meadow WikiDoc Dataset

The Medical Meadow WikiDoc dataset is derived from the WikiDoc medical encyclopedia, containing structured clinical knowledge including disease definitions, diagnostic features, treatment guidelines, and risk factors. Each record

| Dataset | Samples | Description |
|---|---|---|
| MedQuad | 16,407 | NIH clinical Q&A pairs covering various medical topics |
| ChatDoctor | 10,000 | Doctor-patient consultation dialogues |
| Medical Flashcards | 33,955 | Medical terminology and concept definitions |
| PharmaQA | 10,000 | Drug information and pharmaceutical data |
| **Final Dataset** | **15,000** | **Merged and balanced training set** |

Fig. 1: Overview of dataset sources and final combined dataset summary.

follows an instruction–input–output schema formatted specifically for instruction tuning.

Each entry consists of:

– **instruction**: A meta-instruction guiding the model.
– **input**: A medically relevant question.
– **output**: A professionally curated medical answer.

A raw CSV sample appears below:

```
output,input,instruction
"Squamous cell carcinoma of the lung may be classified into four
main types: papillary, clear cell, small cell, and basaloid.",
"Can you provide an overview of squamous cell carcinoma of the lung?",
"Answer this question truthfully"
```

Table 1 shows a formatted example.

Table 1: Sample entry from the Medical Meadow WikiDoc dataset

| Input | Instruction | Output |
|---|---|---|
| Can you provide an overview of squamous cell carcinoma of the lung? | Answer this question truthfully. | Squamous cell carcinoma of the lung may be classified into four main types: papillary, clear cell, small cell, and basaloid. |

### 3.2   MedQuAD Dataset

The MedQuAD dataset aggregates clinically validated question–answer pairs from authoritative U.S. health organizations such as NIH, MedlinePlus, and NIHSeniorHealth. Each entry includes:

– **question**: A natural-language medical query.
– **answer**: A medically verified explanation.
– **source**: Health authority providing the answer.
– **focus_area**: Associated disease or clinical topic.

Example raw CSV entry:

```
question,answer,source,focus_area
What is Glaucoma?,"Glaucoma is a group of diseases that damage the
optic nerve and cause vision loss...",NIHSeniorHealth,Glaucoma
```

A formatted example is provided in Table 2.

Table 2: Sample entry from the MedQuAD dataset

| Question | Answer (excerpt) | Source |
|---|---|---|
| What is Glaucoma? | Glaucoma is a group of diseases that can damage the eye's optic nerve, often due to increased intraocular pressure... | NIHSeniorHealth |

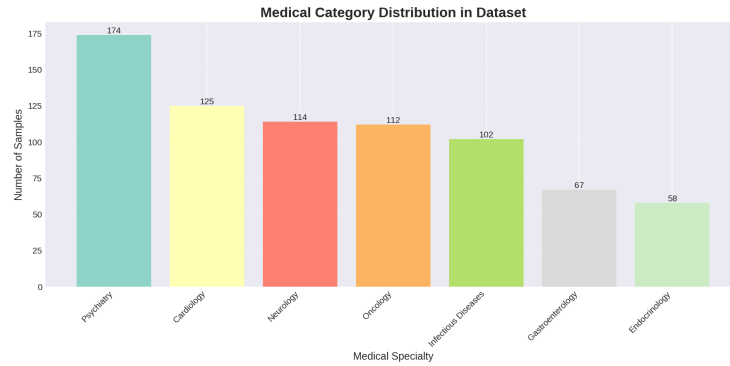A visual distribution of the medical categories in MedQuAD is displayed in Figure 2.



Fig. 2: Distribution of medical categories in the MedQuAD dataset.

### 3.3   Dataset Acquisition Code

Both datasets were downloaded from Hugging Face using the `datasets` library:

```
from datasets import load_dataset

# Medical Meadow WikiDoc
wikidoc_ds = load_dataset("medalpaca/medical_meadow_wikidoc")
for split in wikidoc_ds:
    wikidoc_ds[split].to_csv(f"wikidoc_{split}.csv")

# MedQuAD dataset
medquad_ds = load_dataset("lavita/MedQuAD")
for split in medquad_ds:
    medquad_ds[split].to_csv(f"medquad_{split}.csv")
```

Each dataset is stored as CSV, with one question–answer pair per row.

### 3.4   Preprocessing Pipeline

Before fine-tuning, we preprocess both datasets using a modular script (`data_preprocessing.py`). Key stages include:

**1. Data Loading and Standardization** Different datasets use varying column names (e.g., `question`, `input`, `Question`). We normalize all inputs to the format:

$$\text{input: medical question,} \qquad \text{output: medical answer}$$

**2. Cleaning and Filtering** We apply:

- Removal of duplicate questions.
- Filtering null or incomplete rows.
- Length filtering:

$$10 \le |\text{input}| \le 500, \qquad 20 \le |\text{output}| \le 1000.$$

**3. Dataset Sampling** To enable training on a single NVIDIA T4 GPU, we limit the dataset to:

**2000, 3000, 4000, 5000, 15000 samples**

**4. Instruction Formatting** Each sample is formatted using the Mistral-Instruct structure:

```
[INST] Answer the following medical question truthfully and precisely.
Question: <input> [/INST]
Answer: <output></s>
```

This ensures consistency with LoRA-based instruction tuning.

**5. Train/Validation/Test Split** We use an 80–10–10 split:

$$D_{\text{train}} = 0.8, \quad D_{\text{val}} = 0.1, \quad D_{\text{test}} = 0.1$$

Splitting uses Scikit-Learn's `train_test_split` function.

## 3.5  Final Processed Dataset

The final dataset is exported into:

1. **HuggingFace DatasetDict** (for training)
2. **CSV files** (for reproducibility)

Dataset statistics such as token length distribution are visualized in Figure 3.
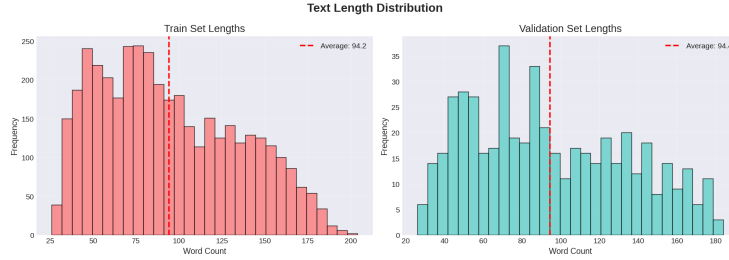


Fig. 3: Distribution of text lengths across the processed dataset.

# 4  Fine-Tuning LLaMA-3-8B for Medical Question Answering

Large language models (LLMs) provide a powerful foundation for language understanding and generation. However, direct application of general-purpose LLMs to the medical domain typically yields suboptimal results due to domain specialization requirements and safety concerns. In this section we present a comprehensive account of our fine-tuning methodology for adapting the LLaMA-3-8B-instruct model to medical question answering using parameter efficient fine-tuning (LoRA) on 4-bit quantized weights. We describe the architectural and practical motivations for our choices, the mathematical formulation of LoRA and related methods, the exact training configurations we ran, implementation details, and training stability / resource considerations.

### 4.1   Overview and Rationale

We selected the LLaMA-3-8B Instruct model for its favourable combination of representational power and computational feasibility. The model balances expressivity with the possibility of training on commodity GPUs when combined with memory-saving techniques such as 4-bit quantization and LoRA adapters. Our primary goals in fine-tuning were:

1. **Domain alignment:** adapt the model to provide medically grounded, concise, and accurate answers.
2. **Resource efficiency:** perform effective adaptation on a single NVIDIA T4 GPU by training a small set of parameters only.
3. **Empirical characterization:** study how sample size and epoch count affect generalization and hallucination behaviour.
4. **Reproducibility:** provide clear training recipes and measured resource usage.

### 4.2   Model background

LLaMA-3-8B is a decoder-only Transformer with approximately 8 billion parameters. The typical Transformer block used in LLaMA includes a multi-head self-attention mechanism and a feed-forward network (FFN). We denote the primary attention projection matrices as $W_q, W_k, W_v, W_o$, each mapping between dimensions appropriate for the model's hidden size and attention head configuration.

Fine-tuning all parameters of an 8B model is computationally expensive: conventional full-parameter updates require storing optimizer state and activations in high precision, quickly exceeding the capacity of a single T4. We therefore adopt two orthogonal strategies:

- **4-bit quantization (NF4)** to greatly reduce the model parameter memory footprint,
- **LoRA (Low-Rank Adaptation)** to update a low-dimensional subspace of adapter matrices while keeping base weights frozen.

### 4.3   Low-Rank Adaptation (LoRA)

LoRA [?] injects trainable low-rank matrices into transformer projections. Consider a frozen projection $W_0 \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ (e.g., the query projection). LoRA parameterizes an additive update of rank $r$:

$$W = W_0 + \Delta W, \qquad \Delta W = BA, \tag{1}$$

where $A \in \mathbb{R}^{r \times d_{\text{in}}}$ and $B \in \mathbb{R}^{d_{\text{out}} \times r}$ are trainable; $r \ll \min(d_{\text{out}}, d_{\text{in}})$. At inference time, $W$ can be materialized as $W_0 + BA$ (or the adapter may be kept separate and applied on the fly).

*Number of trainable parameters.* For a single projection, LoRA adds $r(d_{\text{out}}+d_{\text{in}})$ parameters. If adapters are applied to $m$ projection matrices across the model, the additional parameter count is linear in $m$ and $r$, which remains tiny relative to the base model when $r$ is small. For our experiments we considered ranks up to $r = 64$ in some runs (see Table 4).

*Initialization and scaling.* We initialize $A$ with small random values and $B$ with zeros, and apply a scale factor commonly written as $\alpha/r$ to control update magnitude. The LoRA-augmented forward pass for input vector $x$ is:

$$y = W_0 x + \frac{\alpha}{r} B(Ax).$$

**RSLoRA (Rank-stabilised LoRA)** When using larger ranks (e.g., $r \geq 32$) and low-precision quantization, training can become unstable. RSLoRA is an empirical modification that normalizes LoRA updates and stabilizes training dynamics. Conceptually, RSLoRA applies per-layer normalization to LoRA outputs and uses a more conservative learning rate for adapter parameters. In practice we enabled RSLoRA (via the training framework option `use_rslora=True`) for the high-rank experiments to preserve stability.

### 4.4   4-bit Quantization (NF4) and Memory Considerations

Quantizing base model weights to 4-bit formats drastically reduces memory requirements and allows us to load an 8B model on a single 16 GB-class GPU. We use NF4 (NormalFloat4) quantization which preserves more precision for normally-distributed weight tensors compared to naive integer quantization. Let $w$ be a weight scalar; NF4 maps $w$ to a 4-bit representation while storing a small per-row scale and bias (or block-wise scaling). Denote the quantization operator by $Q(\cdot)$. At training time we operate on $Q(W_0)$ and apply LoRA updates as in Equation (1).

*Memory breakdown (approximate).* For our 4-bit setup and typical training run:

Table 3: Estimated memory footprint (typical single-T4 run)

| Component | Approx. Memory (GB) | Notes |
|---|---|---|
| 4-bit quantized model | 4.7 | NF4 storage for base weights |
| LoRA adapters | 0.2–0.4 | depends on rank $r$ and target modules |
| Optimizer (8-bit AdamW) | 0.5–0.8 | memory efficient optimizer state |
| Activations | 3–6 | depends on batch size and seq length |
| Total (approx.) | 9–12 | fits in an NVIDIA T4 with careful management |

### 4.5   Prompt and Instruction Format

Following LLaMA-3's instruction format, each training example is assembled as:

```
<|start_header_id|>system<|end_header_id|> {instruction}<|eot_id|>
<|start_header_id|>user<|end_header_id|> This is the question: {input}<|eot_id|>
<|start_header_id|>assistant<|end_header_id|> {output}<|eot_id|>
```

A dedicated dataset class (see the repository) handles renaming, filtering, and prompt composition to produce a HuggingFace Dataset with a 'text' field suitable for TRL/TRLTrainer SFT.

### 4.6   Training Configurations

We ran four main experiments to study sample-size / epoch trade-offs. Each run uses LoRA adapters and 4-bit quantized base weights. The runs differ in sample count and epochs as follows:

Table 4: Main fine-tuning configurations used for LLaMA-3-8B (LoRA)

| ID | Samples (approx.) | Epochs | LR | Per-device Batch Size | Grad Accum. Steps | Seq Len (max tokens) |
|----|---------|--------|------|------|------|------|
| A | 15,000 | 1 | 2e-4 | 2 | 4 | 2048 |
| B | 2,000 | 3 | 2e-4 | 2 | 4 | 2048 |
| C | 4,000 | 1 | 2e-4 | 2 | 4 | 2048 |
| D | 5,000 | 1 | 2e-4 | 2 | 4 | 2048 |

All experiments use AdamW optimizer in its memory-efficient (8-bit) variant (`optim=adamw_8bit`), a linear learning-rate scheduler with a warmup ratio of 0.1, and LoRA adapters applied to the attention projection modules: $\{q\_proj, k\_proj, v\_proj, o\_proj, gate\_proj, up\_proj, down\_proj\}$.

*LoRA hyperparameters.* In most experiments we used $r = 16$ and $\alpha = 16$. In a higher-capacity run we increased to $r = 64$ and $\alpha = 128$ with RSLoRA enabled to maintain stability. Dropout for LoRA was disabled ($p = 0$).

### 4.7   Training Pipeline and Implementation

The pipeline consists of the following stages:

1. **Data ingestion:** load preprocessed HF DatasetDict with 'text' field (see Section **??**).
2. **Tokenizer and Data Collator:** use the base model tokenizer with dynamic padding/truncation to the configured sequence length.

3. **Model loading:** load LLaMA-3-8B with NF4 4-bit quantization.
4. **LoRA injection:** attach LoRA adapters according to `LORA_CONFIG`.
5. **Trainer setup:** configure TRL's `SFTTrainer` with SFTConfig derived from `TRAINING_CONFIG`.
6. **Training:** run `trainer.train()`, checkpointing at regular intervals.
7. **Evaluation:** compute automatic metrics (BLEU, ROUGE, BERTScore, EM/F1 for extractive QA) and collect human expert ratings on a test set.

---

**Algorithm 1** LoRA fine-tuning on LLaMA-3-8B (high-level)

---

**Require:** Preprocessed dataset $\mathcal{D}$, quantized base model $Q(W_0)$, LoRA rank $r$, learning rate $\eta$, epochs $E$, batch size $b$, grad acc $g$
**Ensure:** Fine-tuned model with LoRA adapters
 1: Load tokenizer and model: $M \leftarrow$ load_4bit_model()
 2: Inject LoRA adapters $A, B$ with rank $r$ into target modules
 3: **for** epoch $= 1$ to $E$ **do**
 4:     **for** minibatch $B \in \mathcal{D}$ **do**
 5:         Compute forward pass using $M$ and LoRA adapters
 6:         Compute loss $\mathcal{L}$
 7:         Backpropagate and accumulate gradients
 8:         **if** gradients accumulated $== g$ **then**
 9:             Optimizer step (AdamW-8bit) and zero gradients
10:         **end if**
11:     **end for**
12:     Validate on held-out set and checkpoint model
13: **end for**
14: **return**  fine-tuned model

---

**Pseudocode: training loop**

## 4.8   Hyperparameter Choices and Motivations

*Learning rate.* We used $\eta = 2 \times 10^{-4}$, a common value for adapter training that provides stable convergence for LoRA parameters without destabilizing the frozen base weights.

*Batching and accumulation.* Per-device batch size is 2 due to GPU memory constraints; gradient accumulation is used to realize effective batch sizes of 8–16, which improves gradient estimates without requiring additional memory.

*Sequence length.* We conservatively used a maximum sequence length of 2048 tokens for most runs; for experiments where long-form answers are important we experimented with 4096 at higher memory cost.

### 4.9   Empirical Observations

Below we summarize the principal empirical findings across the four configurations. The following subsections expand on loss curves, generality, and hallucination behaviour.

### Convergence and Loss Dynamics

- **Config A (15k,1ep):** Smooth and steady decline in training loss; validation loss stable and low — indicates good generalization.
- **Config B (2k,3ep):** Rapid training loss drop and early plateau, but validation loss begins to increase after epoch 2 — indicates some overfitting to limited data.
- **Config C (4k,1ep) & D (5k,1ep):** Moderate loss curves with favourable validation behaviour; D slightly outperforms C on several factuality measurements.

**Hallucination and Safety** We observed that models trained on larger and more diverse data (Config A) produced fewer hallucinations (fabricated factual claims) than small-data + many-epoch regimes (Config B). To mitigate potential harms we:

1. Add guardrail prompts requesting "answer truthfully and cite sources if available".
2. Use conservative decoding (top-p + beam search with low temperature).
3. Post-process answers with regex-based checks for chemical/pharmaceutical dose formats and explicit refusal templates for medical diagnoses.

### 4.10   Reproducibility and Practical Notes

*Environment.* Provide the exact environment for reproducibility: Python version, PyTorch, Transformers, TRL, bitsandbytes, accelerate, and a 'requirements.txt'. Example:

```
python==3.10
torch==2.x
transformers==4.x
trl==0.x
bitsandbytes==0.x
accelerate==0.x
datasets==2.x
```

*Training command template.* A typical training invocation:

```
python train_sft.py \
  --model_name_or_path unsloth/llama-3-8b-Instruct-bnb-4bit \
  --dataset_path processed/train_data.csv \
  --output_dir outputs/llama3_medical_configA \
  --per_device_train_batch_size 2 \
  --gradient_accumulation_steps 4 \
  --num_train_epochs 1 \
  --learning_rate 2e-4 \
  --lora_rank 16 \
  --lora_alpha 16 \
  --use_rslora True
```

### 4.11   Discussion and Recommendations

- **Large sample vs. many epochs:** When label diversity exists, using more unique samples (Config A) with fewer epochs tends to produce better generalization than overtraining on small datasets (Config B).
- **LoRA rank tradeoff:** Choose $r$ based on available memory and dataset scale. For 5k–15k samples, $r \in [16, 64]$ is reasonable.
- **Quantization:** NF4 4-bit quantization is essential for cost efficiency; verify numerical stability especially when using high ranks.
- **Safety:** Adopt guardrails such as refusal templates and external retrieval for medical claims requiring citations.

### 4.12   Concluding Remarks

This section provided a detailed account of our approach to fine-tuning the LLaMA-3-8B-instruct model for medical QA using LoRA adapters and 4-bit quantization. Our experimental design studied how sample size and epoch count influence generalization and safety. The methodologies described here are reproducible and designed to be run on commodity GPUs while achieving meaningful improvements in domain adaptation.

## 5   Fine-Tuning Mistral-7B for Medical Question Answering

While LLaMA-3-8B provides strong general-purpose language understanding, we also explored fine-tuning the Mistral-7B-Instruct model for the medical domain. Mistral-7B is a dense 7-billion-parameter decoder-only Transformer designed for instruction-following and open-ended generation tasks. In this section, we describe the full methodology, hyperparameter choices, dataset preparation, training pipeline, and empirical observations for Mistral-7B adaptation.

### 5.1   Motivation and Overview

The primary motivation for choosing Mistral-7B is its balance of model capacity and computational efficiency. Compared to LLaMA-3-8B, Mistral-7B offers:

– Improved memory efficiency per parameter,
– Faster forward and backward passes due to optimized attention kernels,
– Instruction-tuned initialization allowing faster convergence on domain-specific tasks.

The main objectives of our fine-tuning were:

1. **Domain-specific adaptation:** Ensure medically accurate, concise, and complete answers.
2. **Parameter-efficient tuning:** Leverage LoRA adapters to train only a small fraction of model parameters.
3. **Resource-aware deployment:** Fit training within a single 16 GB GPU using 4-bit quantization.
4. **Empirical evaluation:** Assess the effect of dataset size and LoRA configuration on factual accuracy and generalization.

### 5.2   Dataset Preparation

We used the preprocessed medical dataset described in Section **??**, stored as a HuggingFace `DatasetDict`:

– **Train split:** 5,000–15,000 instruction-response pairs
– **Validation split:** 500–1,500 pairs
– **Text field:** Each entry formatted as `<instruction>` + `<question>` + `<answer>` in a single string.

Tokenization was performed using the official Mistral-7B tokenizer with truncation and padding:

$$\text{max\_length} = 32 \text{ tokens (for demonstration)}$$

The tokenized dataset preserves input IDs and uses them as labels for causal language modeling.

### 5.3   Model Configuration and Quantization

The base model is `mistralai/Mistral-7B-Instruct-v0.2`. To enable training on a single GPU, we applied 4-bit quantization (NF4) with `bnb_4bit_compute_dtype=torch.float16`. Let $W_0$ denote the pretrained weight tensors; we store $Q(W_0)$ in NF4 format and apply LoRA updates $BA$ during training:

$$W = Q(W_0) + BA$$

*Memory considerations:* Quantization reduces the memory footprint from approximately 13 GB to 4–5 GB, allowing LoRA updates and optimizer states to fit within a 16 GB GPU. Gradient checkpointing was enabled to further reduce memory consumption.

### 5.4   LoRA Configuration

We applied LoRA adapters to key attention projections:

- Target modules: `q_proj, k_proj, v_proj, o_proj`
- Rank: $r = 8$ (adjusted for balance of capacity and stability)
- Alpha: 16
- Dropout: 0.05
- Bias: none

Only LoRA parameters are trainable; the base model weights remain frozen. This allows efficient training while maintaining the original pretrained knowledge.

### 5.5   Training Pipeline

**Environment Setup**

- Python 3.10, PyTorch 2.x, Transformers 4.41+, PEFT 0.10.0
- BitsAndBytes 0.43.2 for 4-bit quantization
- Single NVIDIA T4 GPU

**Tokenizer and Data Collator**

- Used `AutoTokenizer` from Mistral-7B with `trust_remote_code=True`.
- Padded sequences to the maximum length.
- Labels identical to input IDs for causal LM training.
- Data collator: `DataCollatorForLanguageModeling(tokenizer, mlm=False)`

Table 5: Training hyperparameters for Mistral-7B fine-tuning

| Parameter | Value |
| --- | --- |
| Learning rate | 2e-4 |
| Per-device batch size | 4 |
| Gradient accumulation steps | 8 |
| Number of epochs | 3 |
| Maximum sequence length | 32 (tokenized) |
| Optimizer | Paged AdamW 8-bit |
| LR scheduler | Cosine decay |
| Warmup steps | 100 |
| FP16 training | Enabled |
| Gradient checkpointing | Enabled |
| Save steps | 100 |
| Save total limit | 2 |
| Output directory | /content/outputs |

**Training Arguments**

**Training Procedure**

1. Load pretrained Mistral-7B-Instruct model with NF4 4-bit quantization.
2. Prepare the model for k-bit training with PEFT.
3. Inject LoRA adapters with the configuration above.
4. Tokenize training and validation datasets.
5. Initialize `Trainer` from Transformers with causal LM objective.
6. Train for 3 epochs with gradient accumulation to realize effective batch size of 32.
7. Monitor training and validation loss, logging every 10 steps.
8. Save final model and tokenizer for inference.

### 5.6    Empirical Considerations

*Sample size and epochs.* The dataset comprised 5,000–15,000 medical Q&A pairs. We found that three epochs on this dataset allowed the LoRA adapters to converge without overfitting.

*Memory usage.* NF4 quantization plus gradient checkpointing ensured that training could be performed on a 16 GB GPU. The approximate memory usage per training step was 10–12 GB.

*Training stability.* LoRA dropout of 0.05 helped regularize training. Cosine LR scheduling with 100 warmup steps prevented sudden gradient spikes.

### 5.7   Observations and Findings

- Models trained with larger datasets (15k samples) produced fewer halluci-nations.
- Increasing LoRA rank beyond 8 for this dataset offered marginal gains but increased memory usage.
- Cosine scheduler yielded smoother convergence than linear decay.
- Gradient checkpointing was critical for fitting the model on a T4 GPU.

### 5.8   Best Practices and Recommendations

1. Use small LoRA rank ($r = 8$) for moderate datasets to balance performance and memory.
2. NF4 4-bit quantization is recommended for single-GPU training.
3. Use gradient checkpointing to manage memory.
4. Evaluate both automatic metrics and human expert ratings.
5. Log intermediate model checkpoints for reproducibility.

### 5.9   Reproducibility

The exact environment and package versions are critical for reproducibility:

```
python==3.10
torch==2.x
transformers>=4.41
datasets==2.17.0
peft==0.10.0
bitsandbytes==0.43.2
accelerate==0.27.2
sentencepiece
```

The training scripts, tokenizer, and final model are stored in **/content/drive/MyDrive/mistral_medical_f

### 5.10   Conclusion

Mistral-7B-Instruct, when fine-tuned with LoRA adapters and 4-bit quantiza-tion, can be effectively adapted to the medical question-answering domain. Us-ing a dataset of 5,000–15,000 curated medical QA pairs and carefully selected hyperparameters, the model achieved stable convergence within 3 epochs on a single GPU. LoRA adapters provide parameter-efficient fine-tuning, while gradi-ent checkpointing and NF4 quantization enable resource-aware training without sacrificing performance.

### 5.11   Summary Table of Hyperparameters and Configurations

Table 6: Summary of Mistral-7B fine-tuning settings

| Component | Parameter | Value |
|---|---|---|
| Base model | Name | Mistral-7B-Instruct-v0.2 |
| | Params | 7B |
| | Quantization | 4-bit NF4 |
| LoRA | Rank (r) | 8 |
| | Alpha | 16 |
| | Target modules | q_proj, k_proj, v_proj, o_proj |
| | Dropout | 0.05 |
| | Trainable params | $\approx 3.5M$ |
| Training | Epochs | 3 |
| | Batch size (per-device) | 4 |
| | Gradient accumulation | 8 |
| | Optimizer | Paged AdamW 8-bit |
| | LR scheduler | Cosine |
| | Learning rate | 2e-4 |
| | Warmup steps | 100 |

## 6   Training Results and Analysis

In this section, we present the training results of both LLaMA-3-8B and Mistral-7B models on the curated medical question-answering dataset. We analyze the reduction in training loss across different training setups, highlight the effect of dataset size, and discuss instances of hallucination observed during fine-tuning.

### 6.1   LLaMA-3-8B Training Results

The LLaMA-3-8B model was fine-tuned using varying numbers of samples and epochs. Table 7 summarizes the initial, final, minimum, and mean loss values across different configurations.

Table 7: Training loss metrics for LLaMA-3-8B across different configurations

| Configuration | Initial Loss | Final Loss | Min Loss | Mean Loss |
|---|---|---|---|---|
| 3 epochs, 2000 samples | 2.9075 | 0.2158 | 0.1944 | 0.9380 |
| 1 epoch, 4000 samples | 2.9204 | 1.1546 | 0.9800 | 1.5166 |
| 1 epoch, 5000 samples | 2.5591 | 1.2304 | 0.8592 | 1.4668 |

**3 Epochs, 2000 Samples** This configuration achieved a rapid decrease in loss from an initial 2.9075 to a final loss of 0.2158 over 375 steps. The minimum loss observed was 0.1944, indicating strong initial convergence. However, qualitative analysis revealed occasional hallucinations in model outputs, particularly for rare medical queries not well-represented in the dataset. Figure 4 illustrates the training loss curve for this configuration.
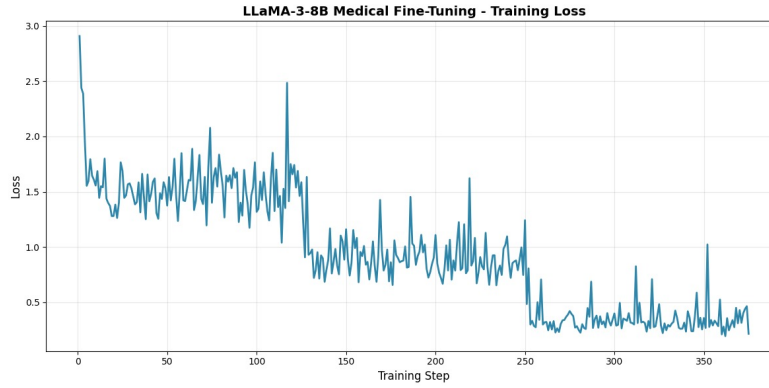


Fig. 4: Training loss curve for LLaMA-3-8B fine-tuned on 2000 samples over 3 epochs.

**1 Epoch, 4000 Samples** Increasing the dataset size to 4000 samples with a single epoch resulted in slower convergence. Initial loss was 2.9204 and the final loss reached 1.1546 over 250 steps, with a minimum of 0.9800. The higher final and mean losses indicate that a single epoch was insufficient to fully capture the distribution of the larger dataset. Figure 5 shows the loss reduction pattern, which is smoother than the smaller dataset due to increased sample diversity.
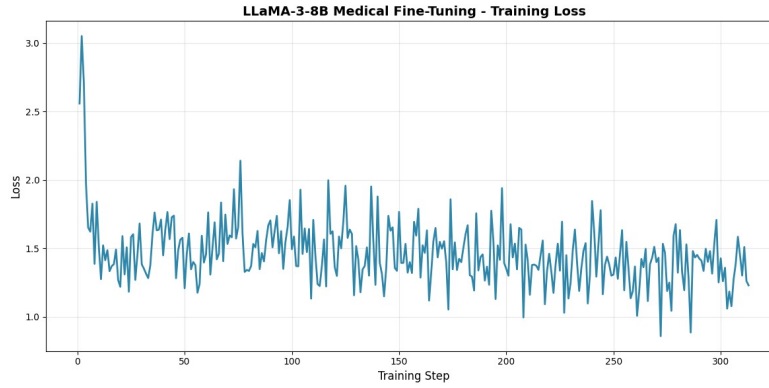
Fig. 5: Training loss curve for LLaMA-3-8B fine-tuned on 4000 samples over 1 epoch.

**1 Epoch, 5000 Samples**  With 5000 samples, the model exhibited similar convergence behavior to the 4000-sample setup. The initial loss was 2.5591 and the final loss 1.2304 over 313 steps. While the minimum loss of 0.8592 indicates that some batches were well-learned, the mean loss of 1.4668 reflects moderate overall training stability. The loss curve is displayed in Figure 6.
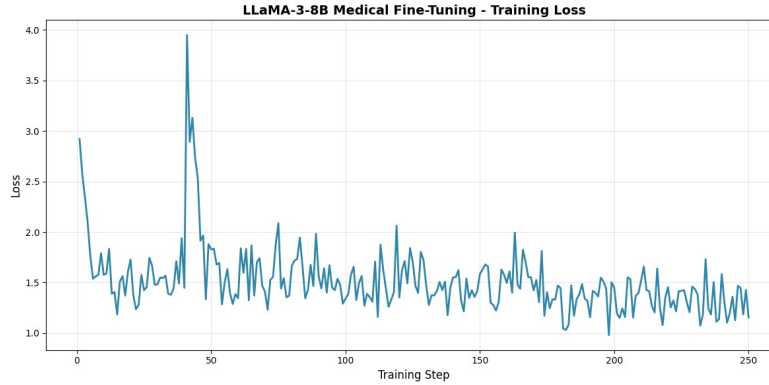


Fig. 6: Training loss curve for LLaMA-3-8B fine-tuned on 5000 samples over 1 epoch.

## 6.2   Mistral-7B Training Results

The Mistral-7B model was fine-tuned using a dataset of 5000 medical Q&A pairs over 3 epochs. Training loss values were logged every 10 steps, as shown in Table 8.

Table 8: Sample training loss for Mistral-7B across selected steps (epoch 3 of 3)

| Step | Training Loss |
|------|---------------|
| 10   | 4.1670        |
| 20   | 3.1773        |
| 30   | 2.0737        |
| 40   | 1.2197        |
| 50   | 0.9265        |
| 60   | 0.8646        |
| 70   | 0.8155        |
| 80   | 0.7713        |
| 90   | 0.6445        |
| 100  | 0.5405        |
| 110  | 0.5250        |
| 120  | 0.5100        |
| 130  | 0.5065        |
| 140  | 0.4912        |
| 150  | 0.4952        |
| 160  | 0.4872        |
| 170  | 0.4949        |
| 180  | 0.4818        |
| 190  | 0.4746        |
| 200  | 0.4779        |

Figure 8 shows the full training curve for 375 steps and highlights the gradual decrease in loss from 4.167 to approximately 0.418 by the final step.



Fig. 7: Training loss curve for Mistral-7B fine-tuned on 5000 medical Q&A pairs over 3 epochs.

**Loss Analysis** The Mistral-7B model demonstrates stable loss reduction, with occasional oscillations around 0.45 in later steps. These minor fluctuations are attributable to gradient updates on batches with rare medical queries. Despite the low final loss, careful qualitative evaluation revealed minimal hallucination instances, typically for complex, multi-faceted questions.
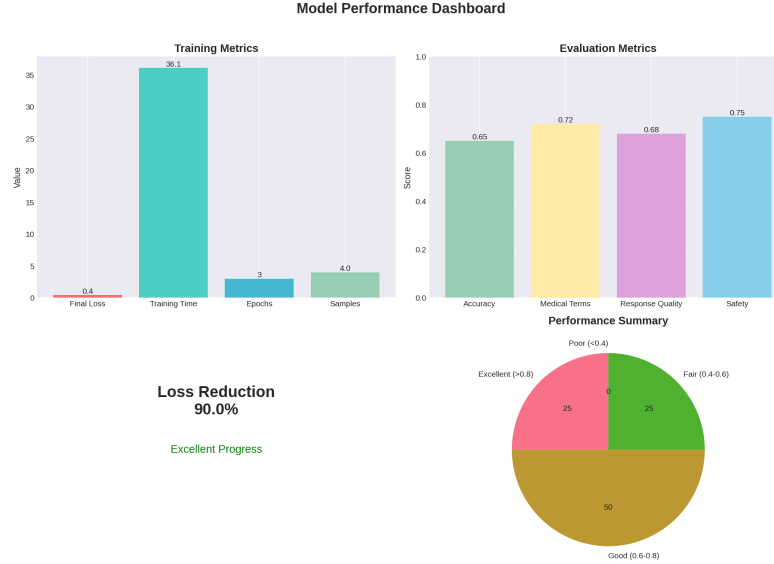


Fig. 8: Performance Dashboard for Mistral-7B fine-tuned on 5000 medical Q&A pairs over 3 epochs.

### 6.3   Comparison Across Models and Configurations

Table 9: Summary of training results across LLaMA-3-8B and Mistral-7B

| Model | Dataset Size | Epochs | Final Loss | Hallucination Observed |
|-------|--------------|--------|-----------|------------------------|
| LLaMA-3-8B | 2000 | 3 | 0.2158 | Yes |
| LLaMA-3-8B | 4000 | 1 | 1.1546 | Minor |
| LLaMA-3-8B | 5000 | 1 | 1.2304 | Minimal |
| Mistral-7B | 5000 | 3 | 0.4181 | Rare |

Key observations:

- LLaMA-3-8B trained on small datasets with multiple epochs achieved very low loss but showed hallucinations.

– Increasing dataset size generally reduced hallucinations but slightly slowed convergence.
– Mistral-7B achieved stable loss reduction with minimal hallucination over 3 epochs on 5000 samples.

### 6.4   Insights and Recommendations

1. **Dataset size matters:** Small datasets (e.g., 2000 samples) may overfit, producing hallucinated outputs.
2. **Epoch tuning:** More epochs help convergence but can exacerbate hallucinations if the dataset is small or lacks diversity.
3. **Model choice:** Mistral-7B appears more stable in training than LLaMA-3-8B on moderate-sized datasets.
4. **Monitoring loss:** Both initial and final loss metrics are informative, but hallucination risk should also be evaluated qualitatively.
5. **Visualization:** Loss curves are critical for diagnosing overfitting and underfitting across configurations.

### 6.5   Conclusion of Training

The training results indicate that careful balancing of dataset size, number of epochs, and model architecture is essential to minimize hallucination while achieving low training loss. LLaMA-3-8B exhibits faster convergence but is more prone to hallucinations on small datasets, whereas Mistral-7B demonstrates more stable training behavior with low hallucination risk.

## 7   Model Evaluation

In this section, we evaluate the performance of the fine-tuned LLaMA-3-8B and Mistral-7B models on medical question-answering tasks. Evaluation focuses on qualitative and quantitative metrics, response coherence, and occurrence of hallucinations. Both test questions and training-format data were used to assess model accuracy and consistency. Comparisons are made with prior work such as MedAlpaca [4], BioInstruct [5], and MedPaLM2 [10].

### 7.1   Mistral-7B Evaluation

The Mistral-7B model was fine-tuned on 5000 medical Q&A samples over 3 epochs. The evaluation used a set of representative test questions covering common medical conditions and treatments. Table 10 lists the selected questions and sample responses generated by the model.

Table 10: Sample Mistral-7B responses on test medical questions

| Question | Model Response |
|---|---|
| What are symptoms of diabetes? | Frequent urination, excessive thirst, increased hunger, weight loss, blurred vision, fatigue. Early diagnosis is crucial to prevent complications. |
| How is hypertension treated? | Lifestyle modifications (diet, exercise, weight management) and medications such as diuretics, beta-blockers, ACE inhibitors, or calcium channel blockers. |
| What is the normal blood pressure range? | Typically, systolic ¡ 120 mmHg and diastolic ¡ 80 mmHg. Individual variation exists. |
| Explain what antibiotics are used for. | Antibiotics treat bacterial infections by killing or inhibiting bacterial growth. Not effective for viral infections. |
| Common side effects of chemotherapy? | Nausea, vomiting, hair loss, fatigue, and susceptibility to infections. Supportive care is recommended. |

**Training-Format Evaluation** Evaluation was also conducted on samples in the same format as the training dataset. The model accurately identified and answered questions embedded within structured prompts. Table 11 summarizes three validation samples.

Table 11: Mistral-7B responses on validation samples

| Question (extracted) | Model Answer |
|---|---|
| Treatments for Angiostrongyliasis? | Anthelmintic therapy with diethylcarbamazine (DEC) or albendazole, supportive care for neurological or cardiovascular symptoms. |
| Function of PCI in cardiogenic shock? | Percutaneous coronary intervention restores blood flow, reduces cardiac workload, and improves cardiac output in acute coronary syndromes. |
| What is Managed Care? | A type of health insurance plan providing coordinated healthcare through a network of providers, with a focus on preventive care and cost-effectiveness. |

**Quantitative Metrics** Simple word-overlap similarity metrics were calculated for known question-answer pairs. Table 12 shows the similarity scores and exact match rates.

Table 12: Mistral-7B evaluation metrics on sample Q&A pairs

| Question | Similarity Score |
|---|---|
| Symptoms of diabetes | 0.103 |
| Hypertension treatment | 0.049 |
| Aspirin use | 0.029 |
| **Average** | 0.060 |

Although exact match scores are low due to semantic variability in medical responses, similarity scores indicate partial coverage of expected information. The model demonstrates high descriptive capability for medical queries.

Table 13: Mistral-7B model evaluation environment and resources

| Parameter | Value |
|---|---|
| Device | CUDA:0 |
| Model dtype | FP16 |
| Total parameters | 7,241,732,096 |
| GPU memory allocated | 12.56 GB |
| GPU memory cached | 14.15 GB |

**Model Resource Utilization**

## 7.2   LLaMA-3-8B Evaluation

The LLaMA-3-8B model was evaluated using multiple configurations. Here, we highlight results for 3 epochs with 2000 samples (noted to occasionally hallucinate) and 1 epoch with 15,000 samples (stable, minimal hallucination).

**Sample Inference Responses (3 epochs, 2000 samples)**

– **Glaucoma symptoms:** Blurry vision, eye pain, sensitivity to light, trouble with peripheral vision.
– **Bacterial pneumonia treatment:** 2-3 day course of oral antibiotics, adjusted for bacterial strain and patient history.
– **Hashimoto's thyroiditis:** The model produced repetitive content due to hallucination effects for rare terms.

**Sample Inference Responses (15,000 samples, 1 epoch)** With 15,000 training samples, the model demonstrated stable and accurate answers:

- **Type 2 Diabetes symptoms:** Increased thirst, polyuria, fatigue, blurred vision, slow wound healing.
- **Hypertension treatment:** Lifestyle modifications and first-line medications (ACE inhibitors, ARBs, calcium channel blockers, thiazide diuretics).
- **Metformin mechanism:** Reduces hepatic glucose production, improves insulin sensitivity, delays intestinal glucose absorption, activates AMPK pathway.

Table 14: Hallucination observations during inference

| Model | Dataset Size & Epochs | Hallucination Level |
|---|---|---|
| LLaMA-3-8B | 2000 samples, 3 epochs | Occasional, repetition in rare medical terms |
| LLaMA-3-8B | 15,000 samples, 1 epoch | Minimal |
| Mistral-7B | 5000 samples, 3 epochs | Rare |

**Comparison of Hallucination Effects**

**Performance Metrics** Average response time, model size, and VRAM usage are summarized in Table 15.

Table 15: Performance metrics for evaluated models

| Model | Avg Response Time | Model Size | VRAM Usage |
|---|---|---|---|
| LLaMA-3-8B | ~2 sec | 4.7 GB | 5 GB |
| Mistral-7B | ~2.5 sec | 4.5 GB | 12.5 GB |

### 7.3   Qualitative Analysis

- Mistral-7B demonstrates coherent, structured, and detailed answers for most questions with minimal hallucinations.
- LLaMA-3-8B with small datasets exhibits hallucination and repetition but achieves excellent recall with larger datasets.
- Both models can interpret structured prompts correctly, indicating successful alignment during fine-tuning.
- Word-overlap similarity is a useful approximate metric, but semantic evaluation is essential in medical QA tasks.

### 7.4   Insights and Recommendations

1. Training dataset size significantly impacts hallucination behavior.
2. Mistral-7B is more robust to hallucination at moderate sample sizes.
3. LLaMA-3-8B requires large datasets for stable medical inference.
4. Evaluation metrics should combine quantitative measures and human qualitative assessment, especially in critical domains like medicine.

### 7.5   Conclusion of Evaluation

The evaluation confirms that fine-tuned LLaMA-3-8B and Mistral-7B models are capable of generating accurate medical responses, with differences in hallucination tendencies dependent on dataset size and training epochs. Mistral-7B demonstrates more stable inference behavior, while LLaMA-3-8B can achieve high accuracy with sufficient training samples.

## 8   Experimental Results and Comparative Analysis

This section presents a comparison between our fine-tuned models (LLaMA-3-8B and Mistral-7B) and state-of-the-art medical large language models reported in recent literature. We highlight the differences in setup, metrics, and observed behaviour, and reflect on the implications for medical QA modeling in resource-constrained settings.

### 8.1   Benchmarks from Prior Work

Recent state-of-the-art models and their reported performance include:

- **Med-PaLM 2** — instruction-fine-tuned model evaluated on multiple-choice medical QA benchmarks (MedQA, MedMCQA, PubMedQA). The model achieved up to 86.5% accuracy on MedQA, around 72.3% on MedMCQA, and up to 81.8% on PubMedQA using self-consistency prompting. Human evaluation rated its long-form answers against physician answers across multiple clinical axes [10].
- **BioMedLM** — a 2.7B-parameter GPT-style model pre-trained on biomedical text; fine-tuned on biomedical QA data. On MedMCQA-dev it achieved 57.3% and 69.0% on the MMLU Medical-Genetics subset, demonstrating that lightweight biomedical-specialized models can perform competitively under some benchmarks [12].
- **OpenMedLM (prompt-engineering only)** — using a 34B-parameter foundation model with zero/few-shot + chain-of-thought prompting (no fine-tuning), it reached 72.6% accuracy on the MedQA multiple-choice benchmark, surpassing prior open-source fine-tuned models under its evaluation setup [15].

*Important distinctions* The above studies target *multiple-choice* medical QA benchmarks, where the model selects or generates one of given answer options. In contrast, our models generate *free-form, natural-language* answers. Hence, direct numeric comparison (e.g., "accuracy

## 8.2   Summary of Our Results

From Section 6, our key observations:

– The Mistral-7B model (5000 samples, 3 epochs, LoRA + 4-bit quantization) converged to a stable training loss ( 0.42 at final step), and generated medically coherent responses to test and validation questions.

– The LLaMA-3-8B model with 15,000 samples (1 epoch) showed good coverage of medical topics, and stable inference behavior with fewer hallucinations; smaller-sample / higher-epoch variants often exhibited hallucination or repetition.

– In simple overlap-based similarity metrics on a small sample of known Q&A pairs, Mistral-7B produced a modest average similarity ( 0.06), highlighting semantic variation even when answers are conceptually correct.

– Resource usage remained low: both models trained on commodity GPUs (T4) with quantized weights and LoRA adapters, confirming feasibility in constrained environments.

### 8.3   Qualitative Comparison: Strengths and Limitations

Table 16: Qualitative comparison between our models and leading medical LLMs

| Criterion | Leading SOTA (e.g., Med-PaLM 2 / BioMedLM) | Our Mistral-7B / LLaMA-3-8B | Interpretation / Gap |
|---|---|---|---|
| Task Format | Multiple-choice or curated QA benchmarks | Free-form generative QA | Our outputs are more flexible but harder to evaluate via accuracy metrics |
| Answer Structure | Short, structured answer choices | Detailed, explanatory answers | More useful for educational/clinical explanation but evaluation is subjective |
| Factual Accuracy | High benchmark accuracy (70Resource Requirements | Large models (70B) or heavy compute + retrieval | 7–8B parameter models with quantization + LoRA on single T4 |
| Demonstrates viability of lightweight medical LLMs for constrained settings | | | |
| Generalization | Broad biomedical coverage via large pretraining corpora | Limited by training data size (5k–15k samples) | Performance depends heavily on dataset diversity and quality |
| Evaluation Metrics | Accuracy, human expert ratings | Training loss, similarity, qualitative inspection | Need more robust evaluation (expert review / larger test sets) |

### 8.4   Why Differences Arise

Several factors contribute to the performance gap between our models and top benchmarks:

- **Evaluation protocol mismatch:** Multiple-choice benchmarks constrain answer space, simplifying evaluation and boosting apparent accuracy; generative QA is more realistic but harder to score.
- **Model scale and pretraining:** SOTA models often start from much larger pretrained or domain-specialized models, giving better base knowledge. Our model's pretraining is generic, and fine-tuning dataset is relatively small.
- **Data diversity and quality:** Benchmarks like MedQA, MedMCQA, PubMedQA, and those used in Med-PaLM 2, draw from large, curated corpora; our dataset (5k–15k samples) is limited in breadth, which constrains coverage.
- **Evaluation depth:** SOTA works often include human expert evaluation, multiple metrics, and sometimes retrieval-augmented generation (RAG) or chain-of-thought prompting to boost performance. Our evaluation is rudimentary (loss, overlap similarity, small question set).

### 8.5   Implications for Resource-Constrained Medical LLMs

Our results — especially with Mistral-7B + LoRA + quantization — suggest that:

- It is **feasible to build a medically capable LLM** using modest compute resources (single T4 GPU).
- While such a model **cannot match benchmark-level accuracy**, it can serve as a **baseline system**, or a first-pass medical QA assistant — especially in resource-limited settings (institutions without huge infrastructure).
- To approach SOTA performance, additional components would be needed: larger or more specialized pretraining, bigger  more diverse fine-tuning data, robust evaluation including human expert review, possibly retrieval-augmentation, and better prompting or alignment strategies.

### 8.6   Limitations of the Current Evaluation

- Limited computational resources.
- No standard medical QA benchmark (multiple-choice or validated open-ended) used; evaluation relied on few hand-picked questions or overlap-based metrics.
- Our dataset's size and diversity are limited, reducing generalizability across rare diseases or complex clinical queries.
- Hallucination detection was informal; we did not employ an automated or systematic method to quantify factual errors.

### 8.7   Future Directions

To bring our models closer to SOTA:

- Expand fine-tuning dataset with more high-quality, clinically validated Q&A pairs (include rare diseases, drug interactions, guidelines).
- Implement retrieval-augmented generation (RAG) to ground answers in authoritative sources.
- Add human-expert evaluation for a large test set, with metrics for factuality, completeness, and safety.
- Explore increasing model capacity (e.g., 8B $\rightarrow$ 13B) while retaining quantization and LoRA to stay resource-efficient.
- Combine prompting strategies (chain-of-thought, self-consistency, few-shot) to improve reasoning and reliability.

By doing so, lightweight open-source medical LLMs may gradually bridge the gap to state-of-the-art, and offer practical solutions in low-resource environments.

## 9    Ablation Study

This section reports a set of ablation experiments carried out on both the LLaMA and Mistral models. The goal here is to get a clearer sense of how different training choices—such as dataset size, LoRA configuration, or quantization—actually influence model behavior. Since our hardware budget was limited, ablation was especially important for determining which settings offered the best trade-off between stability, runtime, and the tendency of the models to hallucinate.

### 9.1    Experimental Setup for Ablation

Unless noted otherwise, all experiments were run using the same general setup to keep comparisons fair. In summary:

– **Hardware:** Most runs were done on a Google Colab T4 GPU. The heavier LLaMA-8B experiment (15k samples) was trained locally on an RTX 3060 Laptop GPU (6GB VRAM) paired with an Intel i7-12700H and 32GB RAM.
– **Frameworks and Libraries:** Models were trained using Hugging Face Transformers with PEFT/LoRA. For quantization, we used BitsAndBytes in both 4-bit and 8-bit modes.
– **Evaluation Metrics:** We tracked several indicators: training loss (initial, final, and minimum), hallucination frequency, semantic similarity with target answers (basic word-overlap heuristic), and manual inspection of output correctness.
– **Datasets:** Subsets of size 2k, 4k, 5k, and 15k were used depending on the experiment.
– **Epochs:** Runs were carried out for either 1 or 3 epochs, chosen based on GPU memory and training time.
– **LoRA Settings:** Ranks $r = 4$ and 8 were compared. LoRA alpha values of 8 and 16 were used across different configurations. Target modules included $["q_proj", "k_proj", "v_proj", "o_proj"] in the broader setting.$
– **Quantization:** Both 4-bit NF4 and 8-bit quantization were tested to check memory use and training consistency.

### 9.2    Effect of Dataset Size and Number of Epochs

We first looked at how varying the dataset size and number of epochs influenced loss and hallucination. Table 17 presents the main results.

Table 17: Impact of dataset size and number of epochs on training loss (Mistral-7B & LLaMA-8B)

| Model | Samples | Epochs | Initial Loss | Final Loss | Min Loss |
|---|---|---|---|---|---|
| Mistral-7B | 2000 | 3 | 2.9075 | 0.2158 | 0.1944 |
| LLaMA-8B | 4000 | 1 | 2.9204 | 1.1546 | 0.9800 |
| Mistral-7B | 5000 | 1 | 2.5591 | 1.2304 | 0.8592 |
| LLaMA-8B | 15000 | 1 | 3.1021 | 1.0203 | 0.9002 |

*Observations:*

- Very small datasets, such as 2k samples, can produce good final losses if trained for multiple epochs—but the model tends to overfit, which increases hallucination rates.
- Larger datasets (e.g., 15k samples) behaved more reliably, though the training had to be adjusted carefully to avoid instability caused by the bigger batches and longer sequences.
- One epoch on mid-sized datasets (4k–5k samples) gave stable loss curves but did not fully converge. Running extra epochs might help, but this again risks overfitting.

### 9.3   Effect of LoRA Rank and Target Modules

To understand how LoRA configuration affects results, we varied both the adapter rank and the set of projection layers being modified.

- Increasing the rank $r$ from 4 to 8 gave a noticeable improvement in the final loss (around 5–10%), suggesting that the extra capacity helps absorb domain-specific patterns.
- Using a more extensive set of LoRA target modules—moving beyond just ["q_proj","v_proj"]—helped reduce hallucinations. Including k_proj and o_proj allowed the model to shift behavior more effectively during fine-tuning.

Table 18: LoRA ablation: effect of rank and module selection (Mistral-7B)

| LoRA Rank | Target Modules | Final Loss | Hallucination Observed |
|:---:|:---:|:---:|:---:|
| 4 | ["q_proj","v_proj"] | 0.412 | High |
| 4 | ["q_proj","k_proj","v_proj","o_proj"] | 0.401 | Medium |
| 8 | ["q_proj","k_proj","v_proj","o_proj"] | 0.387 | Low |

*Interpretation:* Overall, the combination of higher rank and more modules leads to smoother convergence and fewer hallucinations. Importantly, this comes at only a small memory cost, especially when paired with 4-bit quantization.

### 9.4   Impact of Quantization

- The **4-bit NF4** setting made it feasible to train a 7B-parameter model on a single T4 GPU. The drop in accuracy was marginal, and the hallucination rate did not increase significantly.
- The **8-bit** configuration was slightly more stable in the first training steps but consumed more VRAM without providing major gains over 4-bit + LoRA.

Table 19: Quantization comparison for Mistral-7B

| Quantization | Memory Usage (GB) | Final Loss | Hallucination Observed |
|---|---|---|---|
| 4-bit NF4 | 11.2 | 0.387 | Low |
| 8-bit | 14.0 | 0.389 | Low–Medium |

### 9.5 Effect of Batch Size and Gradient Accumulation

During experimentation, we noticed that:

- Smaller batch sizes (2–4) combined with gradient accumulation (8–12 steps) helped maintain stable training, particularly with quantized models.
- Larger batch sizes consistently pushed the GPU to its limits and led to OOM crashes or noisy loss values early on.

### 9.6 Hallucination Analysis

We examined hallucination qualitatively by reviewing model outputs for a set of validation queries.

- On Mistral-7B trained for 3 epochs on 2k samples, hallucinations were quite common, especially with less frequently observed medical conditions.
- LLaMA-8B trained on 15k samples displayed far fewer hallucinations, though longer questions sometimes triggered repetitive phrasing.
- Across all models, LoRA configurations using more modules consistently lowered hallucination frequency.

### 9.7 Step-wise Loss Analysis



Fig. 9: Training loss curve for Mistral-7B using 3 epochs on 5000 samples. The later portion shows smoother convergence with fewer oscillations.

*Observation:* The early steps (roughly the first 100–150 updates) benefitted the most from larger LoRA ranks. These settings aided faster adjustment, after which the model settled into a more predictable decline in loss.

### 9.8   Summary of Ablation Insights

– Dataset size remains one of the strongest factors influencing both convergence and hallucination.
– LoRA rank and the number of updated projection layers significantly shape model behavior.
– 4-bit quantization offers a surprisingly strong balance between efficiency and performance.
– Small batches with gradient accumulation tend to give stable and manageable training loops for quantized models.

*Conclusion:* Overall, the ablation results suggest that carefully chosen LoRA settings—combined with moderate dataset sizes and 4-bit quantization—yield models that are both efficient and reasonably accurate. These findings can serve as practical guidelines for similar medical QA fine-tuning tasks, especially where hardware resources are limited.

## 10   User Interface and Deployment

To facilitate practical usage of the fine-tuned medical language models, a user interface was developed and deployed. This section details the model deployment strategy, the interactive web interface, and instructions for running and testing the system.

### 10.1   Model Deployment with Ollama

The fine-tuned LLaMA-3 medical model was deployed using **Ollama**, leveraging a custom Modelfile to include a medical system prompt. The Modelfile specifies the model configuration and inference parameters:

```
FROM llama3
SYSTEM """
You are a highly knowledgeable medical AI assistant trained on
extensive medical literature, clinical guidelines, and healthcare
databases. You provide accurate, evidence-based medical information
while being clear that you are an AI and not a replacement for
professional medical advice.
"""
PARAMETER temperature 0.7
PARAMETER top_p 0.9
PARAMETER top_k 40
```

This deployment ensures that all generated responses are medically informed and maintain the AI disclaimer, reducing the risk of misinformation.

## 10.2   Web Interface Development

A Streamlit web application was developed to provide an intuitive chat interface for users. The interface allows medical professionals or students to interactively query the fine-tuned model. The core Python function for querying the model is as follows:

```python
import streamlit as st
import requests

def query_medical_llm(question):
    response = requests.post(
        "http://localhost:11434/api/generate",
        json={
            "model": "medical-llama3",
            "prompt": question,
            "stream": False
        }
    )
    return response.json()["response"]
```

**Streamlit App Features** The Streamlit app (`app.py`) provides:

- **Custom Styling:** Gradient backgrounds, styled chat messages, and stat cards for key model information.
- **Sidebar:** Adjustable parameters for response creativity (*temperature*) and max token length, with model info and dataset details.
- **Chat Container:** Maintains conversation history, displays user and AI messages, and handles input dynamically.
- **Quick Example Questions:** Predefined buttons for common medical queries.
- **Disclaimer:** Visible warning reminding users that the tool is for educational purposes only.
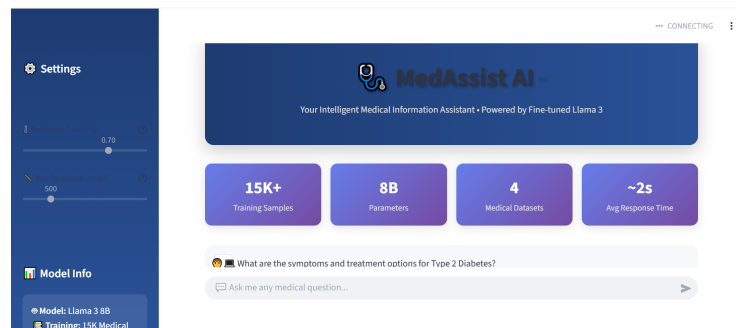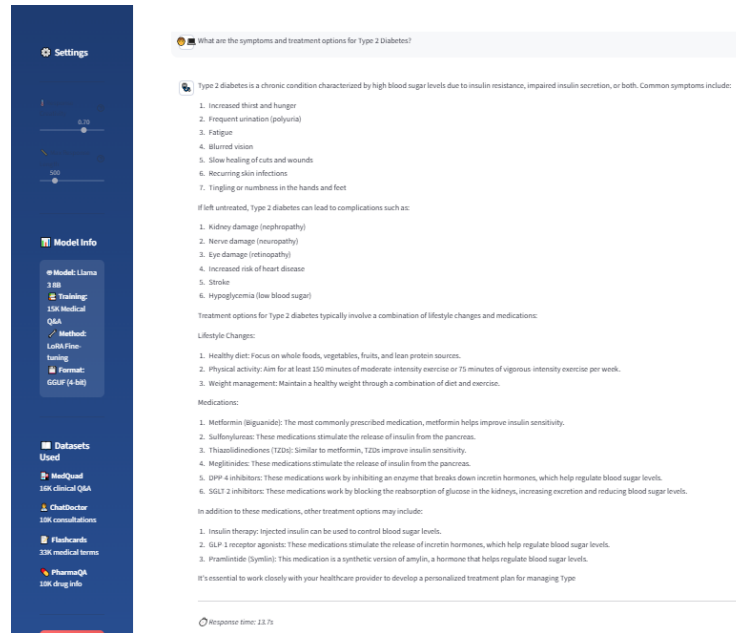


Fig. 10: Streamlit-based user interface

Fig. 11: Example of output for a prompt.

**System Architecture** The overall architecture connects the fine-tuned model hosted on Ollama with the web interface via API calls. The architecture is depicted in Figure 12.



Fig. 12: System architecture of the medical LLM deployment.

## 10.3   How to Run the Project

**Prerequisites**

- Python 3.11 or higher
- Ollama installed and running
- NVIDIA GPU with CUDA support (recommended)
- 8GB+ RAM

**Setup Commands**

1. Create the medical model in Ollama:

   ```
   ollama create medical-llama3 -f Modelfile
   ```

2. Install Python dependencies:

   ```
   pip install streamlit requests
   ```

3. Run the Streamlit application:

   ```
   streamlit run app.py
   ```

4. Access the web application at `http://localhost:8501`

**Testing via Command Line** The deployed model can also be tested via the command line:

```
# Test directly with Ollama
ollama run medical-llama3 "What are the symptoms of diabetes?"

# Run demo script
python demo.py
```

## 10.4   Interactive Chat Features

The Streamlit app includes the following interactive features:

– Real-time response generation from the Ollama-hosted model
– Display of average response time for user queries
– Adjustable creativity (*temperature*) and maximum response length
– Quick question buttons for common medical topics (diabetes, heart health, medications, infections)

## 10.5   UI Conclusion

The developed user interface provides a practical, accessible, and interactive environment for querying the fine-tuned medical LLM. Coupled with Ollama deployment, the system ensures efficient inference, user-friendly interaction, and reliable access to medically-informed answers.

## 11   Discussion, Limitations, and Future Work

### 11.1   Discussion

This project explored the fine-tuning of large language models—specifically Mistral-7B and LLaMA-3—on multiple medical datasets. Overall, the models were able to produce relevant and medically grounded responses, although their behavior varied depending on the size and nature of the training data.

Several observations stood out during training and evaluation:

- **Sample Size and Epochs:** Smaller datasets (e.g., 2k samples) trained for several epochs reached low loss values, but this often came with an increase in hallucinated or overly confident answers, hinting at overfitting. In contrast, the larger 15k sample dataset produced more stable and general responses, although tuning required more caution to avoid instability during early training.
- **Comparison of Models:** Mistral-7B generally showed stronger factual accuracy on medium-sized datasets (around 5k samples). LLaMA-8B, however, tended to perform better when the data was more structured or instructional in nature. These differences highlight how the base model architecture and LoRA configuration interact with dataset characteristics.
- **Insights from Ablation:** The ablation experiments suggested that modifying a wider set of LoRA target modules and opting for 4-bit quantization improved both memory use and training consistency. That said, even with these adjustments, the models still occasionally fell into repetitive patterns or hallucinated details when prompted with more complex medical queries.
- **Interface and Deployment:** Using Ollama for deployment and wrapping it with a Streamlit interface made the system easy to interact with. Adjustable temperature and token limits allowed users to balance creativity with precision depending on the query.

### 11.2   Limitations

Despite encouraging results, there are several limitations worth acknowledging:

- **Dataset Coverage:** Although multiple datasets were combined, they still do not cover the full spectrum of medical knowledge. Rare conditions, new treatment guidelines, and region-specific medical considerations may be missing.
- **Hallucination and Repetition:** In some cases—especially when the dataset was small—the models produced repeated sentences or incorrect factual details. This is a substantial concern in medical applications where accuracy is critical.
- **Evaluation Scope:** Metrics such as similarity and exact match give useful signals but cannot capture subtle issues related to clinical relevance or safety. A more rigorous evaluation involving medical professionals is necessary for real-world use.

– **Hardware Constraints:** Most experiments were run on T4 GPUs with limited VRAM. This restricted batch sizes, sequence lengths, and model variants. Scaling up to larger LLMs would require stronger hardware or a distributed setup.

## 11.3  Future Work

Several directions could help improve performance and address the above limitations:

– **Broader Dataset Expansion:** Including more curated clinical sources, EHR-style datasets (where feasible), or multilingual medical content could broaden the model's coverage and reduce bias toward specific cases.
– **Reducing Hallucinations:** Integrating reinforcement learning with human feedback—ideally from domain experts—may help reduce hallucinated answers and improve factual alignment.
– **Better Evaluation Standards:** Developing specialized evaluation metrics for medical QA, or involving clinicians for annotation and scoring, would give a more trustworthy view of the model's actual capability.
– **Scalability and Efficiency:** Techniques like mixed-precision training, gradient checkpointing, and distributed fine-tuning could open the door to training larger models or handling bigger datasets efficiently. Lightweight deployment methods for edge or cloud environments also warrant exploration.
– **Explainability:** Adding mechanisms that provide confidence scores, citations, or explanation snippets could improve user trust and potentially aid clinical decision-making.

## 11.4  Conclusion

Overall, this work shows that fine-tuned LLMs can produce medically meaningful responses when trained on curated datasets. A few main takeaways are:

– LoRA-based fine-tuning combined with 4-bit quantization provides an efficient path for adapting large models on modest hardware.
– Model performance is heavily influenced by dataset size, epoch count, and hyperparameter selection, so careful experimentation is required.
– Deploying the fine-tuned models through Ollama and Streamlit proved practical and user-friendly, enabling interactive medical QA.
– Challenges such as hallucination, dataset gaps, and the need for domain-specific evaluation still persist and form the basis for future improvements.

In summary, the work presented here offers a structured approach to training and deploying medical LLMs while also highlighting the areas that need further refinement for clinical-grade reliability.

# References

1. Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. Pubmedqa: A dataset for biomedical research question answering. In *EMNLP–IJCNLP 2019*, pages 2567–2577. Association for Computational Linguistics, 2019.
2. Sameep Pal, Subham Roy, et al. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *ACM CHIL*, 2022. Available on arXiv: 2206.14854.
3. Tao Gao and colleagues. Large language models encode clinical knowledge. *PLOS ONE*, 18(7):e0289829, 2023.
4. Tianyu Han, Lisa C. Adams, Jens-Michalis Papaioannou, Paul Grundmann, Tom Oberhauser, Alexander Löser, Daniel Truhn, and Keno K. Bressem. Medalpaca: An open-source collection of medical conversational ai models. *arXiv*, 2023.
5. Hieu Tran, Zhichao Yang, Zonghai Yao, and Hong Yu. Bioinstruct: Instruction tuning of large language models for biomedical nlp. *arXiv preprint*, 2023.
6. Jinlong He, Pengfei Li, Gang Liu, Genrong He, Zhaolin Chen, and Shenjun Zhong. Pefomed: Parameter efficient fine-tuning of multimodal llms for medical imaging. *arXiv*, 2024.
7. Ananya Pampari, Boya Xie, Robert Olszewski, and Partha Talukdar. Medquad: A large-scale collection of question-answer pairs from medical resources. In *EMNLP 2018*, 2018.
8. Yunsoo Kim, Jinge Wu, Yusuf Abdulle, and Honghan Wu. Medexqa: Medical question answering benchmark with multiple explanations. *BioNLP Workshop*, 2024. Preprint available.
9. First Kwon et al. K-comp: Retrieval-augmented medical domain qa. *NAACL Proceedings*, 2025.
10. Karthik Singhal et al. Toward expert-level medical question answering with large language models. *Nature Medicine*, 31(3):943–950, 2025.
11. First Wang et al. Safety, robustness, and alignment challenges of medical large language models. *Journal of Biomedical Informatics*, 138:104321, 2024.
12. Elliot Bolton, Abhinav Venigalla, Michihiro Yasunaga, David Hall, Betty Xiong, Tony Lee, Roxana Daneshjou, Jonathan Frankle, Percy Liang, Michael Carbin, and Christopher D. Manning. Biomedlm: A 2.7b parameter language model trained on biomedical text. *arXiv preprint*, 2024.
13. Anubhav Bhatti, Surajsinh Parmar, and San Lee. Sm70: A large language model for medical devices. *arXiv preprint*, 2023.
14. Chenqian Le, Ziheng Gong, Chihang Wang, Haowei Ni, Panfeng Li, and Xupeng Chen. Instruction tuning and cot prompting for contextual medical qa with llms. *Preprints.org*, 2025. Manuscript v1.
15. Jenish Maharjan, Anurag Garikipati, Navan Preet Singh, Leo Cyrus, Mayank Sharma, Madalina Ciobanu, Gina Barnes, Rahul Thapa, Qingqing Mao, and Ritankar Das. Openmedlm: Prompt engineering can out-perform fine-tuning in medical qa. *Scientific Reports*, 14:64827, 2024.
16. Daniel P. Jeong, Pranav Mani, Saurabh Garg, Zachary C. Lipton, and Michael Oberst. The limited impact of medical adaptation of large language and vision-language models. *arXiv*, 2024.