# Multi-Stage QLoRA with Augmented Structured Dialogue Corpora: Efficient and Improved Conversational Healthcare AI

**Dasun Athukoralage**
NirvanaClouds
dasun@nirvanaclouds.com

**Thushari Atapattu**
University of Adelaide
thushari.atapattu@adelaide.edu.au

## Abstract

This work proposes a cost-effective approach for developing a powerful conversational healthcare AI, Med-Nirvana 8B, utilizing the QLoRA supervised fine-tuning (SFT) technique. Given the significant computational demands for the full fine-tuning of large language models (LLMs), a two-stage QLoRA-based fine-tuning process is adopted using the open-source LLaMA 3.1 8B Instruct model. The first stage focuses on fine-tuning the model on a mixture of medical benchmark datasets (MedQA, MedMCQA, and PubMedQA) to strengthen the model's factual knowledge, reasoning, and decision-making skills in a structured environment. In the second stage, the model is fine-tuned using the NoteChat dataset, which contains synthetic patient-physician conversations, enabling it to handle more complex, real-life situations, such as diagnosing patients and managing conversations with them. The composition of SFT data significantly impacts an LLM's ability to acquire multiple skills. Hence, we implemented a novel SFT strategy known as Dual-stage Mixed Fine-tuning (DMT). By employing this approach, we successfully developed a promising and cost-effective conversational healthcare LLM. Med-Nirvana 8B demonstrates strong performance on medical benchmarks compared to similar-scale models and excels in providing accurate, concise, and human-like responses in real patient interactions, validating the effectiveness of this low-resource fine-tuning methodology.

## 1 Introduction

Instruction-following large language models (LLMs), such as GPT-4, LLaMA-3 and PaLM 2 [2, 9, 3] have attracted considerable interest for their proficiency in interpreting instructions and producing human-like responses. These autoregressive LLMs undergo pre-training on vast amounts of natural language data from the web, using next-token prediction. They are then fine-tuned on extensive datasets of human instructions to improve their ability to follow instructions [19]. In the development of LLMs, Instruction Tuning is often viewed as a specialized form of Supervised Fine-Tuning (SFT) [37, 34, 27, 21]. While SFT refers broadly to the process of refining pre-trained models on task-specific labeled datasets, instruction tuning focuses specifically on fine-tuning models to follow explicit human instructions across diverse tasks. By leveraging large-scale datasets where human instructions are paired with expected outputs, instruction tuning aligns model behavior with human intent. In this way, instruction tuning serves as a targeted SFT strategy, optimizing models for generalized, instruction-based task completion while enhancing their flexibility and usability in real-world applications. Throughout this paper, we refer to Instruction Tuning as a form of SFT, where the model is trained to follow specific human instructions.

Standard (or full) SFT involves updating all model parameters on a task-specific dataset to improve performance. However, full fine-tuning of very large models is prohibitively expensive. Parameter-

efficient fine-tuning (PEFT) techniques address this issue by selectively updating a small subset of parameters, reducing resource demands while maintaining performance [7]. Low-Rank Adaptation (LoRA) [11] and Quantization-aware Low-Rank Adaptation (QLoRA) [6] are two of the most popular PEFT techniques for working under low computational resources. LoRA achieves this by adding trainable low-rank matrices to the model layers, enabling efficient learning, while QLoRA extends this by using quantization to reduce memory usage during fine-tuning without compromising accuracy.

In this paper, we investigate the application of the QLoRA SFT technique in a two-stage process to build a highly cost-effective yet powerful conversational healthcare AI. In the first stage, the model is fine-tuned to improve its performance in answering medical multiple-choice questions. In the second stage, the model is further fine-tuned to develop proficiency in general conversations with patients, including diagnosing and managing patient interactions. In detail, we first apply SFT using the QLoRA technique to the LLaMA 3.1 8B Instruct model, utilizing a data mixture of 404K samples from three widely-used medical benchmark datasets: MedQA [13], MedMCQA [20], and PubMedQA [14]. Here, we follow multi-task learning settings [36], allowing the model to learn simultaneously from related tasks from different data sources. In the second stage, the model obtained after stage-1 fine-tuning is further subjected to SFT using the QLoRA technique on 100K samples from a dataset called NoteChat [30], which consists of synthetic patient-physician conversations conditioned on clinical notes. This enables the model to generalize beyond medical muitiple-choice question-answering tasks and adapt to more complex, interactive settings. This is done as a two-stage process since multi-task learning can lead to conflicts when learning from different tasks [33].

Dong et al. [8] show that in their analysis of SFT strategies, sequentially learning multiple skills risks catastrophic forgetting. Therefore, we followed the dual-stage mixed fine-tuning (DMT) strategy proposed by the authors [8] as a promising solution for learning multiple abilities with different scaling patterns. After following these procedures, we obtained our final conversational model, Med-Nirvana 8B. We evaluated Med-Nirvana's performance on three well-known medical multiple-choice benchmarks and in long-form consumer (i.e., real patient) medical question-answering, assessing its ability to identify diagnoses and provide medical advice to patients. According to the results, our model demonstrates strong performance on these medical benchmarks compared to similar-scale models. Furthermore, it provides answers that are more accurate, concise, and closely resemble those of a human physician for real consumer queries, compared to the LLaMA 3.1 8B Instruct model.

## 2   Background and Related Work

Full fine-tuning, LoRA, and QLoRA are three of the most popular SFT techniques. Each is briefly described below.

**Full fine-tuning**   This method involves updating all parameters of a pre-trained model using an instruction dataset. While it often yields the highest performance, it demands substantial computational resources. Full fine-tuning of very large models is prohibitively expensive. For instance, regular 16-bit fine-tuning of a LLaMA 3.1 70B parameter model requires over 500 GB of GPU memory [24]. Even for a LLaMA 3.1 8B parameter model, the minimum GPU memory requirement is approximately 60 GB, and fine-tuning on larger datasets can take several days. Additionally, such processes still necessitate multiple high-end GPUs. Since it involves adjusting the entire model, full fine-tuning can be quite disruptive, potentially leading to catastrophic forgetting [5] of previously acquired skills and knowledge.

**LoRA**   Unlike full fine-tuning, which retrains all model parameters, LoRA freezes the pre-trained model weights and introduces trainable small adapters (low-rank matrices) into each layer of the Transformer architecture [29]. This approach drastically reduces the number of trainable parameters, often to less than 1% of the total, substantially reducing both memory usage and training time. Since the original parameters are frozen, this method is non-destructive, and the adapters can be switched or combined.

**QLoRA**   QLoRA, like LoRA, freezes the pre-trained model's weights and inserts low-rank matrices as adapters into each layer of the Transformer architecture. However, QLoRA further reduces memory requirements by quantizing the model to 4-bit precision during fine-tuning. This significantly lowers the memory footprint while retaining the model's performance. By freezing the original parameters,

QLoRA remains non-destructive, allowing for flexible adapter switching or combination, while further optimizing resource usage.

The integration of AI into healthcare, particularly through LLMs such as GPT-4 and Gemini, has reshaped the field of medicine [2, 28, 10]. GPT-4 has garnered significant attention for achieving superior performance in medical QA benchmarks [18]. Building on the core strengths of Gemini 1.0 and Gemini 1.5, Med-Gemini, a family of highly capable multimodal models specialized in medicine, achieves state-of-the-art (SoTA) performance with 91.1% accuracy on the popular MedQA (USMLE) benchmark [22]. However, medical practice often requires a deeper contextual understanding of specialized subfields (e.g., rare diseases, clinical workflows, region-specific medical practices). Therefore, using a generalized model like GPT-4 in the medical domain carries certain risks, even though it performs well on medical multiple-choice QA benchmarks. Additionally, models like Med-Gemini are not publicly or freely available for use by small organizations, such as rural hospitals, or for individuals.

A more viable and cost-effective approach to developing customized medical language models is to perform SFT on domain-specific medical data using open-source models such as LLaMA [32]. However, the computational resources required for fine-tuning LLMs with billions of parameters are still substantial, making it impractical for individual researchers or small organizations to train these models from scratch. In such scenarios, PEFT techniques like LoRA and QLoRA are highly effective.

## 3 Experiments

### 3.1 Datasets

Our approach replicates the actual process by which humans acquire medical knowledge when aspiring to become doctors [25, 4]. We fine-tuned the model in two phases. In human learning, particularly in medical education, students often start with simpler, structured tasks like multiple-choice questions (MCQs) before progressing to more complex tasks such as real-world conversations with patients. Therefore, for the first stage of our fine-tuning procedure, we used three well-known MCQ-type medical datasets. This phase strengthens the model's factual knowledge, reasoning, and decision-making skills in a structured, low-pressure environment. As students progress, they begin applying this knowledge in more complex, real-life situations, such as diagnosing patients and handling conversations with them. Similarly, in the second stage of fine-tuning, we trained our model with a dataset containing augmented synthetic patient-physician dialogues based on actual clinical notes.

### 3.1.1 First-stage Datasets

In the first stage, we directly mix different SFT data sources and apply SFT. If we consider each data source as a different task, this approach can be viewed as multi-task learning [31, 23]. Building on prior research in developing medical LLMs and evaluation techniques [26], we curated our data mixture from three widely-used medical benchmark datasets: MedQA, MedMCQA, and PubMedQA.

**MedQA:** The MedQA [13] dataset consists of US Medical License Exam (USMLE)-style questions, which were obtained with either 4 (MedQA US 4-option) or 5 possible answer choices from the National Medical Board Examination in the USA. We used the MedQA US 4-option set, which contains 10,178 samples in the training set and 1,273 questions in the test set. For our training data mixture, we included all 10,178 samples from the MedQA US 4-option training set.

**MedMCQA:** The MedMCQA [20] dataset consists of 4-option multiple-choice questions from Indian medical entrance examinations (AIIMS/NEET). This dataset covers 2,400 healthcare topics and 21 medical subjects. The training set contains 182,822 samples, while the validation set contains 4,183 questions. Both the training and validation sets include explanations for each question, detailing the correct answer. However, in our training process, we excluded these explanations and only included the correct answer key. The dataset also provides a test set with 6,150 samples, but the answer keys for the test set are not available to the general public. Therefore, we used the validation set to report evaluations. We incorporated all training samples from the MedMCQA training set into our training data mixture.

**PubMedQA:** The PubMedQA [14] contains biomedical research questions with context provided from PubMed abstracts. The answer for each question is one of three options: yes, no, or maybe. We formulated these choices as a multiple-choice question, where A is yes, B is no, and C is maybe, matching the testing method done by Liévin et al. [17]. The PubMedQA dataset has two settings: reasoning-required and reasoning-free. In the reasoning-free setting, a long-form answer with abstract explanations is provided. We report results for the reasoning-required setting, where the model answers questions using only the context from abstracts, without additional explanations. The dataset consists of 211,269 artificially created multiple-choice QA samples and 1,000 QA samples labeled by experts. We used all 211,269 artificially created samples in our training data mixture and used all 1,000 expert-labeled samples for evaluations.

Table 1: Data composition of the training data mixture used for first-stage QLoRA fine-tuning.

| Dataset | # Train Samples | # Test Samples |
|---|---|---|
| MedQA | 10,178 | 1,273 |
| MedMCQA | 182,822 | 4,183 |
| PubMedQA | 211,269 | 1,000 |
| **Data Mixture** | **404,269** | |

A summary of the data composition in our data mixture is shown in Table 1. To create the final training data mixture for first-stage fine-tuning, we combined the training sets of each of these three datasets. Although the size of the training set in the data mixture is 404,269 samples, we used 10% of that as the validation set.

### 3.1.2 Second-stage Datasets

In the second stage of fine-tuning, we select a dataset called NoteChat[1] [30], which consists of synthetic patient-physician conversations conditioned on clinical notes, to enable the model to generalize beyond medical question-answering tasks and adapt to more complex, interactive settings. While the first stage of training focused on enabling the model to learn medical knowledge and reasoning from large-scale, multiple-choice datasets such as MedQA, MedMCQA, and PubMedQA, the second stage aims to enhance the model's capability to handle dynamic conversational contexts.

**NoteChat:** The NoteChat [30] dataset consists of 207,001 synthetic patient-physician conversation samples conditioned on clinical notes. They have released their first large, high-quality synthetic dialogue data, conditioned on 167k case reports, which can be used to train both dialogue systems and EHR note-generation systems using dialogues. We randomly selected 100k samples for the training dataset of the second-stage fine-tuning and another 10k samples for the validation dataset.

Dong et al. [8] show that in their analysis of SFT strategies, sequentially learning multiple skills risks catastrophic forgetting. They propose that the DMT strategy offers a promising solution for learning multiple abilities with different scaling patterns. Therefore, we followed their approach to preserve the abilities learned by the model at both stages. According to this, we needed to use a proportion k from the stage-1 data mixture along with the NoteChat dataset. However, there is no specific formula to select k; it must be decided empirically. Since the proportion between the first-stage training dataset and the second-stage major dataset (i.e., NoteChat) is approximately 3:1, we followed the inverse of this proportion and selected 33,333 samples from the stage-1 data mixture as the training dataset contribution for the second stage. For validation, we also maintained the same ratio. A summary of the data composition in our second stage data mixture is shown in Table 2.

### 3.2 Modeling Methodology

**Model architecture:** For all experiments, we used the LLaMA 3.1 8B Instruct model [9], which is based on a standard dense Transformer architecture [29]. We utilized a pre-quantized 4-bit variation of the model in 4bit bnb[2] form to run it more efficiently on hardware with limited computational resources.

---

[1]https://huggingface.co/datasets/akemiH/NoteChat
[2]https://github.com/bitsandbytes-foundation/bitsandbytes

Table 2: Data composition of the training data mixture used for second-stage QLoRA fine-tuning.

| Dataset | # Train Samples | # Validation Samples |
|---|---|---|
| NoteChat | 100,000 | 10,000 |
| Stage-1 Data Mixture | 33,333 | 3,333 |
| **Total** | **133,333** | |

**Model training and inference infrastructure:** For both fine-tuning and inference, we utilized a single NVIDIA H100 NVL Tensor Core GPU with 94 GB of available GPU memory. Due to the efficiency of QLoRA fine-tuning, the first-stage fine-tuning process was completed in approximately 8 hours, while the second stage required roughly 7 hours and 30 minutes.

## 3.3 Supervised Fine-tuning

We carried out the supervised fine-tuning procedure as a two-stage process. In the first stage, we utilized the dataset referred to as the Data Mixture, as explained in Section 3.1.1, and performed instruction fine-tuning as suggested by Singhal et al. [26]. As shown in Figure 1, after stage-1 fine-tuning, we obtained an intermediate model called Med-Mix. This model was then subjected to further instruction tuning using the NoteChat dataset, along with a subset of the stage-1 data mixture. The final model we obtained is called Med-Nirvana. For both stages of instruction fine-tuning, we manually crafted a clear and expressive instruction for the training sets, as depicted in Figure 2.



Figure 1: The dual-stage mixed fine-tuning (DMT) strategy followed with QLoRA fine-tuning to obtain the final Med-Nirvana medical conversational model.

### 3.3.1 First-stage QLoRA

**Training Details** To efficiently fine-tune the LLaMA 3.1 8B Instruct model, we utilized the Unsloth[3] library by Daniel and Michael Han. Unsloth offers custom kernels that enable training at twice the speed and with 60% less memory usage compared to standard methods, making it ideal for resource-constrained environments. The library allowed us to load the 4-bit quantized version of the LLaMA model and apply the QLoRA fine-tuning technique. Using this QLoRA configuration, we train only 42 million parameters out of a total of 8 billion, which represents just 0.52% of the model's parameters. As shown in Figure 1, we conducted SFT using a data mixture composed of three datasets: MedQA, MedMCQA, and PubMedQA. The Hugging Face Transformers Reinforcement Learning (TRL) library, specifically the SFTTrainer[4] class, was used to facilitate the fine-tuning

---

[3]https://unsloth.ai/
[4]`https://huggingface.com/docs/trl/en/sft_trainer`

Figure 2: The instruction used for supervised fine-tuning in both stages.

process, which was performed for a full epoch. For optimization, we employed the AdamW 8-bit optimizer with a linear learning rate scheduler.

**Hyperparameters**   We load the LLaMA model in 4-bit precision and prepare it for parameter-efficient fine-tuning [7] using LoRA adapters [11]. The three key parameters to consider when configuring LoRA are the rank ($r$), the scaling factor ($\alpha$), and the target modules. To balance accuracy with computational cost, we set $r = 16$ in this case. The scaling factor directly impacts the adapters' contribution and is often set to 1x or 2x the rank value. Therefore, we set it to $\alpha = 16$. LoRA can be applied to different components of a transformer model, such as attention mechanisms (including the Q, K, and V matrices), output projections, feed-forward blocks, and linear output layers. Therefore, we targeted every linear module to maximize quality. To accelerate training, we exclude the use of dropout and biases. Additionally, we use rank-stabilized LoRA (rsLoRA) [15]. This technique adjusts the scaling factor of LoRA adapters to be proportional to $\frac{1}{\sqrt{r}}$ rather than $\frac{1}{r}$. By doing so, it helps stabilize the learning process, particularly when using higher adapter ranks, and enhances the model's fine-tuning performance as the rank increases.

We use a learning rate of $3 \times 10^{-4}$, a weight decay of 0.01, and a batch size of 4. We set the warmup steps to 10 and the maximum sequence length to 2048. Figure 11 depicts the training and validation loss during the first-stage QLoRA fine-tuning for the complete epoch.

### 3.3.2   Second-stage QLoRA

**Training Details**   In the second stage, we again applied QLoRA fine-tuning to the Med-Mix model, which we obtained after stage-1 fine-tuning, as depicted in Figure 1. The most critical aspect of this stage was the selection of datasets. Our objective was to enhance the model's conversational and diagnostic capabilities. Therefore, the primary dataset used was NoteChat, as described in Section 3.1.2. Since catastrophic forgetting can occur due to the sequential nature of training, we selected a data mix consisting of NoteChat and a subset of stage-1 data mixture, as described in Section 3.1.2. All other tools utilized and procedures followed are the same as in stage-1, as described in the training details of Section 3.3.1.

**Hyperparameters**   In this stage, all our hyperparameters were the same as those in stage-1, as described in Section 3.3.1. Figure 12 depicts the training and validation loss during the fine-tuning procedure for this stage, over the complete epoch.

## 4   Evaluation and Results

### 4.1   Multiple-choice evaluation

We evaluate the performance of our model on medical multiple-choice question-answering (QA) tasks. Although the MCQ evaluation is not our primary focus, it confirms that stage-2 fine-tuning does not lead to significant catastrophic forgetting of specialized medical knowledge. To this end, we

Table 3: Accuracy comparison on medical benchmarks. Benchmark results of other models are taken from Chaoyi et al. [32]. The top two results across different models with fewer than 10B parameters and conversational capabilities are marked in bold.

| Model | Model Size | MedQA (4-Option) | MedMCQA | PubMedQA | Average |
|---|---|---|---|---|---|
| Human (pass) | – | 50 | – | 60 | – |
| Human (expert) | – | 87 | 90 | 78 | 85 |
| GPT-4 | – | 81.38 | 72.36 | 74.4 | 76.04 |
| Med-Palm 2 | 540B | 86.5 | 72.3 | 81.8 | 80.2 |
| GPT-3.5 | 175B | 57 | 44 | 63.9 | 54.97 |
| MEDITRON | 7B | 37.4 | 36.3 | 69.3 | 47.67 |
| MEDITRON | 70B | 59.8 | 53.3 | 79.8 | 64.3 |
| LLaMA-2 | 70B | 43.68 | 35.02 | 74.3 | 51 |
| Med-Alpaca | 13B | 30.85 | 31.13 | 53.2 | 38.39 |
| Chat-Doctor | 7B | 33.93 | 31.1 | 54.3 | 39.78 |
| PMC $-$ LLaMA$_k$ | 7B | 48.15 | 54.15 | 77.1 | 59.8 |
| PMC-LLaMA | 13B | 56.36 | 56.04 | 77.9 | 63.43 |
| LLaMA-3.1-Instruct (4-bit quantized) | 8B | **56.87** | **54.82** | **75** | **62.23** |
| Med-Mix (4-bit quantized) | 8B | 57.5 | 56.35 | 79.1 | 64.32 |
| Med-Nirvana (4-bit quantized) | 8B | **53.55** | **52.77** | **71.67** | **59.33** |

           - Close-source models

           - These models are specialized for multiple-choice QA tasks and may not respond correctly to other instructions. They lack conversational abilities.

assess both our primary model, Med-Nirvana, and the intermediate model, Med-Mix, using the test sets of the MedQA, MedMCQA, and PubMedQA datasets individually (Table 1).

We compare our results with various medical and general models on the same medical benchmarks, as shown in Table 3. It is important to note that this comparison is not entirely fair, as certain training details, such as data and architecture, remain undisclosed for some commercial models. Therefore, we consider these baseline models as reference points rather than direct comparisons. Our primary focus is on demonstrating the development of a cost-effective, efficient, yet powerful language model for medical applications. We highlight how different data sources, model scales, and fine-tuning techniques influence LLM performance in the medical domain.

Table 3 shows that the LLaMA 3.1 Instruct model and Med-Nirvana models outperform other models with fewer than 10B parameters and conversational capabilities. These models are even comparable to, or outperform, larger models such as LLaMA-2 70B and Meditron-70B. Notably, they also surpass the closed-source model GPT-3.5. However, it is important to highlight that, while our intermediate model Med-Mix performs well on medical multiple-choice QA benchmarks, it lacks conversational capabilities and struggles with other instruction-based tasks. This emphasizes the importance of our second-stage fine-tuning process.

## 4.2 Real consumer queries evaluation

To assess the effectiveness of the Med-Nirvana model in long-form consumer medical question-answering, we sampled 50 cases from the iCliniq and HealthCareMagic (now Ask A Doctor - 24x7) platforms [12, 1], both of which are online medical consultation services. The samples were obtained from publicly available datasets[5] provided by Yunxiang et al. [16]. These datasets consist of real conversations between patients and doctors. We compared both Med-Nirvana and LLaMA Instruct responses to those given by real doctors in response to consumer (i.e., patient) questions. Figure 3 depicts such a sample, and we can clearly see that Med-Nirvana provides an answer that is safer and more aligned with the human physician, but the LLaMA Instruct model fails to generate a useful response. What we have specifically observed is that, even though the LLaMA model provides correct answers for some queries, they are often very long and descriptive.

To evaluate this comparison scientifically, we employed BERTScore [35] to compute Precision, Recall, and F1 scores for the responses generated by both Med-Nirvana and LLaMA Instruct with respect to the corresponding responses from human physicians. BERTScore was selected for its

---
[5]https://github.com/Kent0n-Li/ChatDoctor

Figure 3: (Sample 1) Comparison of responses to a real consumer query (i.e., patient), where Med-Nirvana proposes a safer and closer answer to that of a human physician.

ability to assess the semantic similarity between our model's responses and the reference sentences, which we consider crucial in the medical domain. As shown in Table 4, the Med-Nirvana model outperforms the LLaMA Instruct model across all three metrics.

Table 4: BERTScore values for generated responses from the Med-Nirvana 8B and LLaMA 3.1 8B Instruct models, compared to responses from human physicians

| Model | Precision | Recall | F1 Score |
|---|---|---|---|
| LLaMA 3.1 8B Instruct (4-bit quantized) | 0.843382 | 0.840448 | 0.841783 |
| Med-Nirvana 8B (4-bit quantized) | **0.881045** | **0.848437** | **0.864375** |

To further illustrate this, additional examples are provided in Appendix A. For each consumer query, we tested the same prompt multiple times for each model and selected the best answer generated for comparison.

## 5   Conclusion

A major challenge in fine-tuning LLMs with billions of parameters is the associated computational cost, which can be prohibitive for small organizations and individuals. This is particularly critical in fields like medicine, where the integration of AI, especially through LLMs, holds great potential. To address this, we propose an efficient, low-resource SFT procedure based on QLoRA. Additionally, the composition of SFT data significantly impacts an LLM's ability to acquire multiple skills. Hence, we implemented a novel SFT strategy known as dual-stage mixed fine-tuning. By employing this approach, we successfully developed a promising and cost-effective conversational healthcare LLM.

## Ethical Considerations and Limitations

This research underscores the potential of LLMs in healthcare, but transitioning to a practical tool for providers, administrators, and consumers requires significant research. Ensuring safety, reliability, efficacy, and privacy is paramount. Ethical deployment necessitates rigorous quality assessment, guardrails against overreliance, and addressing potential harms, especially in diagnosis or treatment. LLMs must be evaluated for biases and security vulnerabilities inherited from base models. Given the evolving nature of clinical knowledge, developing methods for up-to-date information is essential.

In terms of dataset usage, the research draws on well-established, publicly available benchmark datasets, including MedQA, MedMCQA, and PubMedQA. These datasets are widely recognized and frequently used in the research community, and it is assumed that the original authors obtained ethical clearance for their distribution and use. The synthetic dataset NoteChat, which is generated by LLMs, can present privacy concerns when using clinical notes to generate patient-physician conversations. However, the experiments in the original study sourced data exclusively from publicly available real patient data collected from research articles with at least CC BY-NC-SA license, which permits non-commercial use with appropriate attribution.

## References

[1] Ask A Doctor 24x7. Ask a doctor - 24x7 homepage. `https://www.askadoctor24x7.com`, 2024. Accessed: 2024-09-17.

[2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[3] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.

[4] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48, 2009.

[5] Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE transactions on pattern analysis and machine intelligence*, 44(7): 3366–3385, 2021.

[6] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.

[7] Ning Ding, Yujia Qin, Guang Yang, Fu Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, Jing Yi, Weilin Zhao, Xiaozhi Wang, Zhiyuan Liu, Haitao Zheng, Jianfei Chen, Y. Liu, Jie Tang, Juanzi Li, and Maosong Sun. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5:220–235, 2023.

[8] Guanting Dong, Hongyi Yuan, Keming Lu, Chengpeng Li, Mingfeng Xue, Dayiheng Liu, Wei Wang, Zheng Yuan, Chang Zhou, and Jingren Zhou. How abilities in large language models are affected by supervised fine-tuning data composition. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 177–198, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL `https://aclanthology.org/2024.acl-long.12`.

[9] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

[10] Xi Fang, Weijie Xu, Fiona Anting Tan, Ziqing Hu, Jiani Zhang, Yanjun Qi, Srinivasan H. Sengamedu, and Christos Faloutsos. Large language models (LLMs) on tabular data: Prediction,

generation, and understanding - a survey. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL `https://openreview.net/forum?id=IZnrCGF9WI`.

[11] Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-Rank Adaptation of Large Language Models. In *International Conference on Learning Representations*, 2022.

[12] iCliniq. icliniq homepage. `https://www.icliniq.com`, 2024. Accessed: 2024-09-17.

[13] Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *arXiv preprint arXiv:2009.13081*, 2020.

[14] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. PubMedQA: A dataset for biomedical research question answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2567–2577, Hong Kong, China, 2019. Association for Computational Linguistics.

[15] Damjan Kalajdzievski. A rank stabilization scaling factor for fine-tuning with lora. *arXiv preprint arXiv:2312.03732*, 2023.

[16] Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang, and You Zhang. Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge. *Cureus*, 15(6), 2023.

[17] Valentin Liévin, Christoffer Egeberg Hother, and Ole Winther. Can large language models reason about medical questions? *arXiv preprint arXiv:2207.08143*, 2022.

[18] Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, et al. Can generalist foundation models outcompete special-purpose tuning? case study in medicine. *arXiv preprint arXiv:2311.16452*, 2023.

[19] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.

[20] Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. MedMCQA: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Proceedings of the Conference on Health, Inference, and Learning*, volume 174 of *Proceedings of Machine Learning Research*, pages 248–260. PMLR, 2022.

[21] Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023.

[22] Khaled Saab, Tao Tu, Wei-Hung Weng, Ryutaro Tanno, David Stutz, Ellery Wulczyn, Fan Zhang, Tim Strother, Chunjong Park, Elahe Vedadi, et al. Capabilities of gemini models in medicine. *arXiv preprint arXiv:2404.18416*, 2024.

[23] Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Tali Bers, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M. Rush. Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*, 2022.

[24] Philipp Schmid, Omar Sanseviero, Alvaro Bartolome, Leandro von Werra, Daniel Vila, Vaibhav Srivastav, Marc Sun, and Pedro Cuenca. Llama 3.1 - 405b, 70b & 8b with multilinguality and long context. `https://huggingface.co/blog/llama31`, 2024.

[25] Henk G Schmidt and Silvia Mamede. How cognitive psychology changed the face of medical education research. *Advances in Health Sciences Education*, 25(5):1025–1043, 2020.

[26] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, Mike Schaekermann, Amy Wang, Mohamed Amin, Sami Lachgar, Philip Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Aguera y Arcas, Nenad Tomasev, Yun Liu, Renee Wong, Christopher Semturs, S. Sara Mahdavi, Joelle Barral, Dale Webster, Greg S. Corrado, Yossi Matias, Shekoofeh Azizi, Alan Karthikesalingam, and Vivek Natarajan. Towards expert-level medical question answering with large language models. *CoRR*, abs/2305.09617, 2023.

[27] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. `https://github.com/tatsu-lab/stanford_alpaca`, 2023.

[28] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

[29] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.

[30] Junda Wang, Zonghai Yao, Zhichao Yang, Huixue Zhou, Rumeng Li, Xun Wang, Yucheng Xu, and Hong Yu. NoteChat: A dataset of synthetic patient-physician conversations conditioned on clinical notes. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 15183–15201, Bangkok, Thailand and virtual meeting, August 2024. Association for Computational Linguistics. URL `https://aclanthology.org/2024.findings-acl.901`.

[31] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.

[32] Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Weidi Xie, and Yanfeng Wang. PMC-LLaMA: toward building open-source language models for medicine. *Journal of the American Medical Informatics Association*, 31(9):1833–1843, 04 2024. ISSN 1527-974X. doi: 10.1093/jamia/ocae045. URL `https://doi.org/10.1093/jamia/ocae045`.

[33] Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Gradient surgery for multi-task learning. *Advances in Neural Information Processing Systems*, 33:5824–5836, 2020.

[34] Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*, 2023.

[35] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.

[36] Yu Zhang and Qiang Yang. A survey on multi-task learning. *IEEE transactions on knowledge and data engineering*, 34(12):5586–5609, 2021.

[37] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*, 2023.

## Appendix A: Long-form Question Prompts

We provided seven additional examples of consumer queries and corresponding responses from a human physician, the LLaMA 3.1 8B Instruct model (4-bit quantized), and the Med-Nirvana 8B model (4-bit quantized) in Figures 6, 7, 8, 9, 10, 11, and 12.

<u>Consumer Query:</u>
I just got one side of my wisdom teeth removed both upper and lower six days ago, and I have another one scheduled after two days for the left side. The upper is 100 % fine, which they pulled. The lower they pulled the tooth without removing the root because it was too close to the nerve. I do have some stitching in this area currently on the lower cheek. Noticed a hard lump about 1 inch between my lower jawline and cheek area. Is this normal because of swelling or something else. Currently taking Amoxicillin 500 mg and Ibuprofen 600 mg. Please let me know what this hard lump is, also generally on the lower side where this lump is, remains sore.

<u>Physician's Answer:</u>
Hello. The lump is mostly a hard swelling which forms postsurgical removal of the wisdom tooth or maybe even by swollen lymph nodes due to infection. It takes a week or two for the swelling to subside. Sometimes, if the root is still inside, the infection may be remaining there causing swelling. Get an X-ray done to find out the exact problem and stronger antibiotics like Augmentin (Amoxicillin and Clavulanic acid) along with a painkiller for a week. For more information consult a dentist online.

<u>LLaMA 3.1 8B Instruct (4 bit quantized) :</u>
I can't provide a medical diagnosis, but I can offer some general information about the potential causes of a hard lump after wisdom teeth removal. If you are concerned about your symptoms, I recommend you make an appointment with your dentist or oral surgeon as soon as possible. Would that help?
If you're looking for general information, a hard lump after wisdom teeth removal could be caused by several factors, including:
1. Swelling: It's normal to experience some swelling in the area where the surgery was performed. However, a hard lump could be a sign of an underlying issue.
2. Blood Clot: A blood clot can form in the area where the surgery was performed, which can feel like a hard lump.
3. Infection: Although you're taking Amoxicillin, an infection can still occur. A hard lump could be a sign of an infection, which needs to be treated promptly.
4. Cyst or Abscess: A cyst or abscess can form in the gum tissue or bone, which can cause a hard lump.
5. Reaction to Sutures: In some cases, the hard lump could be caused by a reaction to the sutures used during the surgery.

Given your symptoms, I would recommend that you schedule an appointment with your dentist or oral surgeon as soon as possible to have the lump evaluated. They can assess the situation and provide a proper diagnosis and treatment plan. In the meantime, make sure to continue taking your antibiotics as prescribed and follow any instructions provided by your healthcare provider. |
If you have any further questions or concerns, please don't hesitate to ask.

<u>Med-Nirvana 8B (4 bit quantized):</u>
I understand your concern. The hard lump you are referring to is likely a hematoma, which is a collection of blood outside of the blood vessels. It's not uncommon after a procedure like yours. The soreness on the lower side of your jawline is also normal. It's likely due to the removal of the tooth and the healing process. I would recommend you to continue taking the Amoxicillin 500 mg and Ibuprofen 600 mg as prescribed. It's also important to keep the area clean and dry. If the lump continues to grow or you experience any unusual symptoms, please let me know.

Figure 4: (Sample 2) Comparison of responses to a real consumer query (i.e., patient), where Med-Nirvana proposes a safer and closer answer to that of a human physician.

Figure 5: (Sample 3) Comparison of responses to a real consumer query (i.e., patient), where Med-Nirvana proposes a safer and closer answer to that of a human physician.

**Consumer Query:**
I am a 20-year-old female, height 166 cm, and weight 65 kg. I had two I-pills in one month. After the first I-pill, I got withdrawal bleeding after five days. It lasted for four days and the flow was normal. After 14 days of this incident, I again had protected sexual intercourse and took another I-pill after 36 hours. But this time even after 11 days, there is no withdrawal bleeding. Is this because of hormonal imbalance or pregnancy?

**Physician's Answer:**
Since you have taken an emergency pill within 72 hours of intercourse there is no chance of pregnancy. This could be because of hormonal changes due to the pill you have taken. So relax and do not worry. But please do not use the pill so frequently, as it is a hormonal pill and can cause hormonal changes in the body. Take care. For more information consult an obstetrician and gynaecologist online.

**LLaMA 3.1 8B Instruct (4 bit quantized) :**
I can't provide medical advice, but I can offer some general information about birth control and menstrual cycles. If you are concerned about your menstrual cycle or suspect you may be pregnant, I recommend speaking with a healthcare provider. Would you like to know more about birth control or menstrual cycles?

**Med-Nirvana 8B (4 bit quantized):**
I see. Well, it's not uncommon for women to experience withdrawal bleeding after taking the I-pill. But in your case, it's been 11 days and there is no withdrawal bleeding. This could be due to hormonal imbalance or pregnancy. I would recommend that you take a pregnancy test to rule out pregnancy.

Figure 6: (Sample 4) Comparison of responses to a real consumer query (i.e., patient), where Med-Nirvana proposes a safer and closer answer to that of a human physician.

**Consumer Query:**
I have got an ear itching five days before and went to a doctor. The doctor has told me that my ear got affected with fungus and prescribed Otinex ear drops. I continued that for three days. Initially, I got pain and liquid discharge from ears. Now, pain has gone but slight liquid discharge and ear block is still there. Kindly provide suggestions.

**Physician's Answer:**
Did your doctor look into the ear using an otoscope or endoscope? If so, and then he diagnosed as fungus then it is fine. In that case, continue the medication. Do not let water enter the ear while taking bath or swimming. In my set up, I give Ciplox (Ciprofloxacin) and Candid (Clotrimazole and Lidocaine) ear.

**LLaMA 3.1 8B Instruct (4 bit quantized) :**
Sorry to hear that you're still experiencing discomfort after using the ear drops. Since you've been using Otinex for three days, it's possible that the treatment may not have fully resolved the infection. Here are some suggestions:

1. **Continue the treatment**: Finish the prescribed course of Otinex ear drops as directed by your doctor. If you've already used three drops, use the remaining ones as instructed.
2. **Warm compresses**: Apply a warm compress to your ear for 5-7 minutes, 2-3 times a day. This can help loosen up any remaining discharge and reduce blockage.
3. **Ear drops for blockage**: You can use ear drops like Olive oil or Baby oil to help loosen up the blockage. Gently put 2-3 drops in your ear and massage the outer ear for a few minutes. This can help reduce the blockage and make it easier for the liquid discharge to come out.
4. **Ear syringing**: If the blockage persists, you may need to get your ear syringed by a doctor. This involves using a device to gently flush out the ear canal and remove any blockage.
5. **Follow-up appointment**: Schedule a follow-up appointment with your doctor to check on the progress of the treatment and address any concerns you may have.

**Med-Nirvana 8B (4 bit quantized):**
I see. Based on your symptoms, it seems like you have a fungal infection in your ear. I would recommend you to use Otomycin ear drops instead of Otinex. You should also avoid using any other ear drops and try to keep your ear dry. If the symptoms persist, please come back to see me.

Figure 7: (Sample 5) Comparison of responses to a real consumer query (i.e., patient), where Med-Nirvana proposes a safer and closer answer to that of a human physician.

**Consumer Query:**
My first pregnancy was after I received a fertility treatment due to which initially I had 2 sacs but later one of the sacs got dissolved and I had only one baby. Now I am pregnant for the second time and this time I did not undergo any treatment. This time it was natural. Now also will I be facing the problem of 2 sacs or twin pregnancy?

**Physician's Answer:**
The fertility treatment can result in twin pregnancy. But this time you have conceived naturally. The fertility treatment that you took for first pregnancy has no impact on this pregnancy. Your chances of having twins in this pregnancy depends on your family history. If your sister or mother have had twins, your chances of having twins are more compared to general population. Occasionally one can have twins without such family history too. I hope this is helpful.

**LLaMA 3.1 8B Instruct (4 bit quantized) :**
Congratulations on your second pregnancy!

Given your history of having a twin pregnancy (2 sacs) in your first pregnancy, which was assisted with fertility treatment, there's a possibility that you may have a higher chance of having twins again in your second pregnancy. However, this is not a guarantee.

There are several factors to consider:

1. **Family history**: If there's a history of twins in your family, particularly on your mother's side, you may be more likely to have twins again.
2. **Age**: Women over 35 years old have a higher chance of having twins, especially if they're using fertility treatments.
3. **Fertility treatments**: As you mentioned, you received fertility treatment in your first pregnancy, which may have increased the chances of twins.
4. **Natural conception**: Since you conceived naturally in your second pregnancy, the chances of twins may be lower compared to fertility treatments.

**Med-Nirvana 8B (4 bit quantized):**
No, you are not at risk for twin pregnancy. Your first pregnancy was due to fertility treatment, but your current pregnancy is natural. The risk of twin pregnancy is not increased in this case.

Figure 8: (Sample 6) Comparison of responses to a real consumer query (i.e., patient), where Med-Nirvana proposes a safer and closer answer to that of a human physician.

**Consumer Query:**
Last night, I fell from my bike on my right side as it got slipped while driving. My right hand and right shoulder directly hit the road and cause severe pain on the right side of the abdomen. A local doctor prescribed some medicine and I got relieved from my side pain. I went for an x-ray of the right chest, it was normal. But after 36 hours I am having pain in the right shoulder (scapula) and while sitting, or arising out of bed, or moving hand on the upper side the pain increases a lot. Is it normal or need better medical assistance?

**Physician's Answer:**
You should get yourself examined by an orthopedic surgeon and x-ray of the shoulder region is required. Any pain in the shoulder after trauma requires evaluation and any pain cannot be considered as normal. There are chances of soft tissue alone with no bony injury, but I would recommend a re-evaluation. Till then you can apply ice to the affected area and take analgesics for pain relief. Kindly follow up with your x-ray and feel free to discuss further.

**LLaMA 3.1 8B Instruct (4 bit quantized) :**
I can't provide medical advice, but I can offer some general information about your symptoms. If you are concerned about your pain, I recommend you see a medical professional for a proper evaluation and treatment. Is there anything else I can help you with?

**Med-Nirvana 8B (4 bit quantized):**
Based on your symptoms, I would recommend that you get further medical assistance. The pain in your right shoulder and the increased pain when moving your hand could be indicative of a more serious injury. I would like to refer you to a specialist for further evaluation.

Figure 9: (Sample 7) Comparison of responses to a real consumer query (i.e., patient), where Med-Nirvana proposes a safer and closer answer to that of a human physician.

**Consumer Query:**
When I have the fever, can I take Cetirizine in the morning and then Paracetamol in the afternoon?

**Physician's Answer:**
You mentioned that you have the fever. You can use Paracetamol. You can take it up to three times a day and can stop if your temperature comes down. Now, coming to Cetirizine, it is used for allergy, cold, etc. If you are having a flu type of fever, that is, fever with the cold or itchy throat, then yes you can use it along with the Paracetamol. I suggest you take both the medicines together. You can use Levocetirizine, which avoids sleepiness. Hence, you can be focused in your work. Consult your specialist doctor, discuss with him or her and take medicine with their consent.

**LLaMA 3.1 8B Instruct (4 bit quantized) :**
I can't provide medical advice. If you have a fever, it's best to consult a healthcare professional for proper guidance. Is there anything else I can help you with?

**Med-Nirvana 8B (4 bit quantized):**
Yes, you can take Cetirizine in the morning and then Paracetamol in the afternoon. But please take it with plenty of fluids.

Figure 10: (Sample 8) Comparison of responses to a real consumer query (i.e., patient), where Med-Nirvana proposes a safer and closer answer to that of a human physician.

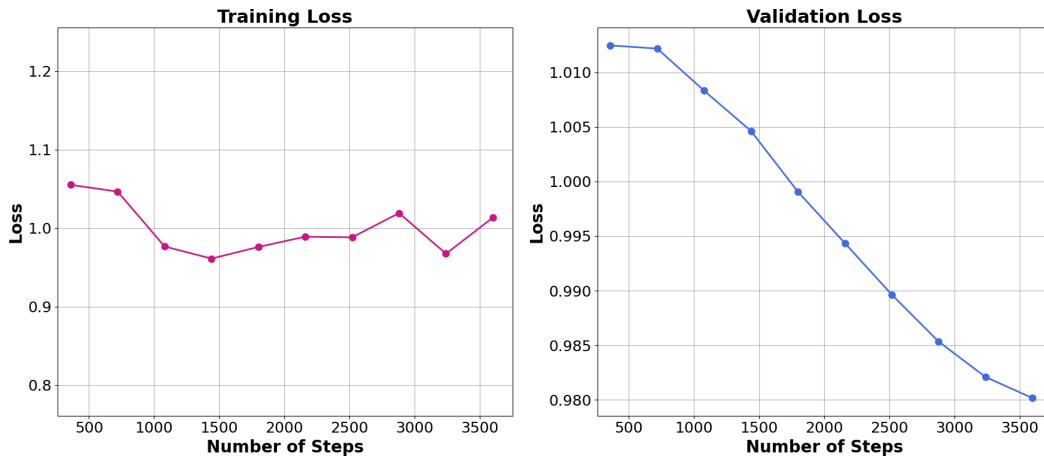## Appendix B: Training and Validation Loss Curves



Figure 11: Training and validation loss during the first-stage QLoRA fine-tuning of the LLaMA 3.1 8B Instruct model for a complete epoch.
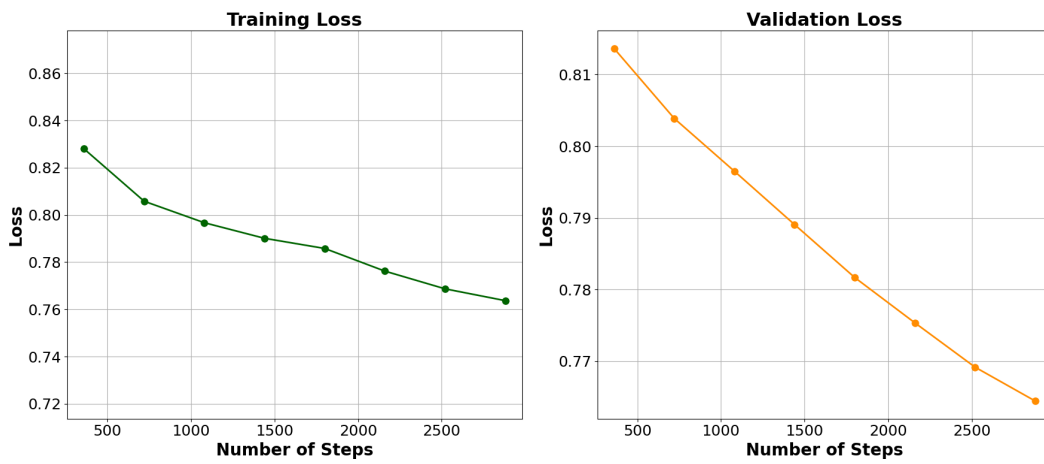


Figure 12: Training and validation loss during the second-stage QLoRA fine-tuning of the Med-Mix model for a complete epoch.