

# Parameter-Efficient Fine-Tuning of LLaMA for the Clinical Domain

Anonymous ACL submission

## Abstract

Adapting pretrained language models to novel domains, such as clinical applications, traditionally involves retraining their entire set of parameters. However, this approach is increasingly proven to be impractical owing to the substantial computational requirements associated with training such large language models. To address this issue, Parameter-Efficient Fine-Tuning (PEFT) techniques offer a viable solution by selectively fine-tuning a small subset of additional parameters, significantly reducing the computational requirements for domain adaptation. In this study, we propose Clinical LLaMA-LoRA, a PEFT adapter layer built upon the open-sourced LLaMA model. Clinical LLaMA-LoRA is trained using clinical notes obtained from the MIMIC-IV database, thereby creating a specialised adapter designed for the clinical domain. Additionally, we propose a two-step PEFT framework which fuses Clinical LLaMA-LoRA with Downstream LLaMA-LoRA, another PEFT adapter specialised for downstream tasks. We evaluate this framework on multiple clinical outcome prediction datasets, comparing it to clinically trained language models. Our proposed framework achieves a state-of-the-art AUROC score averaged across all clinical downstream tasks. We observe substantial improvements of 6-9% AUROC score in the large-scale multilabel classification tasks, such as diagnoses and procedures classification.

## 1 Introduction

Large Language Models (LLMs) have consistently achieved state-of-the-art performance across various NLP tasks. However, while these models exhibit impressive generalisation abilities, they often struggle to perform in specialised domains such as clinical applications, primarily due to the absence of domain-specific knowledge. The complexity of medical terminology and the presence of incomplete sentences in clinical notes contribute to this

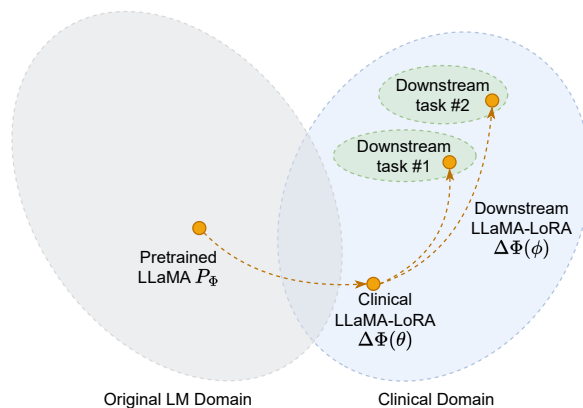


Figure 1: An illustration of the proposed two-step PEFT framework. Clinical LLaMA-LoRA fine-tunes the pretrained LLaMA to the clinical domain. Downstream LLaMA-LoRA further fine-tunes the domain-adapted model to downstream clinical tasks.

challenge (Lehman and Johnson, 2023). Unfortunately, studies have indicated that even LLMs pretrained with datasets comprising biomedical publications still exhibit suboptimal performance when applied to downstream clinical applications, particularly when compared to LLMs pretrained with clinical notes (Alsentzer et al., 2019; Li et al., 2022; Yang et al., 2022). This observation suggests that there are intrinsic nuances specific to the clinical context that can only be effectively captured if LLMs undergo pretraining using clinical datasets.

The current approach of adapting pretrained LLMs to the clinical domain typically involves fine-tuning the entire model parameters (Alsentzer et al., 2019; Peng et al., 2019; van Aken et al., 2021; Michalopoulos et al., 2021; Lehman and Johnson, 2023). However, due to the rapid increase in the size of LLMs, such a practice demands extensive computational resources, which may not be readily accessible to all researchers. Consequently, this challenge will further exacerbate the disparity between the resource-rich and resource-constrained research institutions (Ruder et al., 2022).

To address the substantial computational demands, studies have proposed various Parameter-Efficient Fine-Tuning (PEFT) techniques. These techniques present a practical solution by fine-tuning a small subset of additional parameters while keeping the remaining pretrained parameters fixed. As a result, this strategy significantly alleviates the computational burden while achieving comparable performance to that of full fine-tuning.

In this study, we propose a two-step PEFT framework (see Figure 1). Firstly, we introduce Clinical LLaMA-LoRA, a Low-Rank Adaptation (LoRA, Hu et al., 2022) PEFT adapter built upon the open-source Large Language Model Meta AI (LLaMA) (Touvron et al., 2023). Then, we introduce Downstream LLaMA-LoRA, which is trained on top of the pretrained Clinical LLaMA-LoRA. Downstream LLaMA-LoRA is specifically designed for clinical downstream tasks. The fusion of the two adapters achieves state-of-the-art performance in clinical NLP downstream tasks while considerably reducing the computational requirements. This study presents the following contributions:

- We introduce Clinical LLaMA-LoRA, a PEFT-adapted version of the LLaMA model tailored specifically for the clinical domain.
- We provide comparisons of multiple PEFT techniques in terms of language modelling performance based on perplexity score, shedding light on the optimal PEFT techniques for the clinical domain-adaptive pretraining.
- We introduce Downstream LLaMA-LoRA, built on top of Clinical LLaMA-LoRA and tailored specifically for the clinical downstream tasks.
- We evaluate the proposed mixture of Clinical LLaMA-LoRA and Downstream LLaMA-LoRA on downstream clinical datasets and tasks. Our proposed framework showcases improvements in AUROC scores over the existing clinical LLMs.

## 2 Background

### 2.1 Biomedical Large Language Models

General-domain LLMs continue to face challenges when confronted with domain-specific tasks. The complexity associated with the requisite domain knowledge is recognised as a significant factor (Ling et al., 2023), particularly within the

biomedical domain. Consequently, numerous studies have attempted to adapt LLMs specifically for the biomedical domain.

An early example of such adaptation is BioBERT (Lee et al., 2019), which was pretrained using biomedical research articles from PubMed and PubMed Central. This adaptation has shown improved performance across various biomedical NLP tasks. Recognising the significance of biomedical-specific vocabularies, Gu et al. (2022) proposed PubMedBERT, which is pretrained on biomedical data from scratch and initialised the model vocabulary with the biomedical corpus. The growing interest in biomedical NLP research has led to the adaptation of even larger models to the biomedical domain (Luo et al., 2022; Singhal et al., 2022; Wu et al., 2023; Singhal et al., 2023)

While these biomedical LLMs have demonstrated advancements in various biomedical NLP benchmarking tasks, studies have revealed that clinical LLMs still outperform their biomedical counterparts in numerous clinical downstream tasks (Alsentzer et al., 2019; Yang et al., 2022; Li et al., 2022; Lehman and Johnson, 2023). This suggests that domain-adaptive pretraining using clinical data is still the *de facto* protocol in adapting LLMs to the clinical domain.

### 2.2 Clinical Large Language Models

Clinical LLMs are often fine-tuned with clinical data from an LLM that is already pretrained with datasets that encompass broader topics. For instance, Bio+ClinicalBERT (Alsentzer et al., 2019) is domain-adaptively pretrained using clinical notes from the Medical Information Mart for Intensive Care (MIMIC)-III database (Johnson et al., 2016), starting from a pretrained BioBERT (Lee et al., 2019), which itself is pretrained on biomedical articles. BlueBERT (Peng et al., 2019) is domain-adaptively pretrained using PubMed abstracts and MIMIC-III clinical notes from a BERT model (Devlin et al., 2019), that is pretrained with general-domain texts. Similarly, Clinical-T5 (Lehman and Johnson, 2023) is domain-adaptively pretrained using the union of MIMIC-III and MIMIC-IV (Johnson et al., 2023) clinical notes from T5-base (Raffel et al., 2020), another general-domain LLM.

All these studies share a common approach, which is to fine-tune the entire model parameters. With massive LLMs, this method has become cost-prohibitive and inaccessible for many researchers.

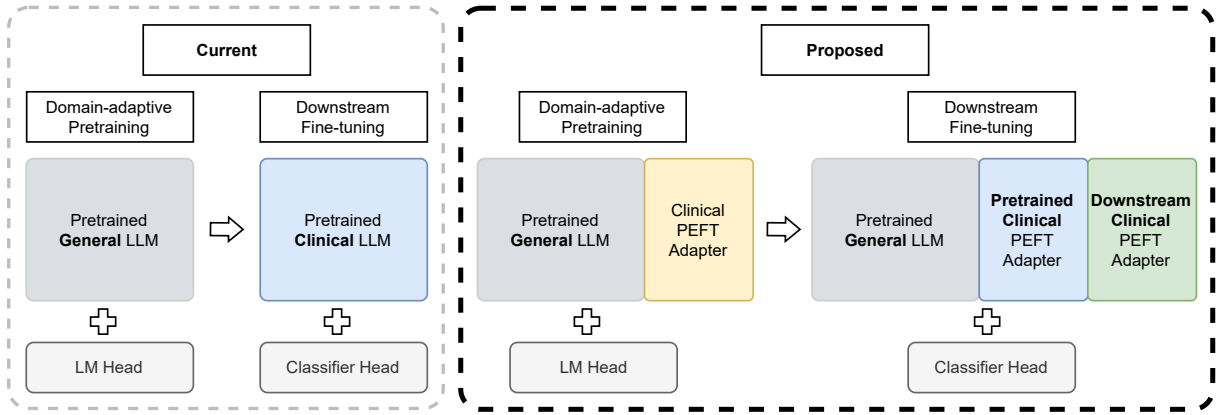


Figure 2: Frameworks of domain-adaptive and downstream fine-tuning to adapt a pretrained LLM from the general domain to the clinical domain. As opposed to a full fine-tuning process which can be prohibitively expensive (left), our approach leverages PEFT techniques to introduce a clinically-specialised adapter that is attached to a pretrained general LLM (right). Our proposed framework also introduces another clinical PEFT adapter trained on the downstream clinical tasks, such as clinical note classification.

### 2.3 Parameter-Efficient Fine-Tuning for Large Language Models

Suppose that we have a pretrained LLM  $P_{\Phi}(y|x)$ ; fine-tuning it can be effectively defined as finding the most appropriate parameter changes  $\Delta\Phi$  by optimising the fine-tuning objective. A conventional, full fine-tuning process means that the model needs to learn a  $\Delta\Phi$  whose dimension is equal to the entire parameters of the pretrained LLM  $|\Delta\Phi| = |\Phi_0|$ , which is computationally expensive. PEFT techniques address this by tuning the *delta*  $\Delta\Phi$ , which corresponds to a very small fraction of additional trainable parameters during the fine-tuning process.

Adapter tuning (Houlsby et al., 2019) is an early PEFT method that involves adding small additional parameters called *adapters* to each layer of the pretrained model and strictly fine-tuning this small set of new parameters. LoRA (Hu et al., 2022) is another PEFT approach that trains low-rank matrices to represent the attention weights update of transformer-based models.

Another group of PEFT approaches leverages the concept of prompting. Prefix Tuning (Li and Liang, 2021) optimises a sequence of continuous task-specific vectors, called a *prefix*, which are trainable parameters that do not correspond to real tokens. P-Tuning (Liu et al., 2021b) uses a similar strategy as Prefix tuning with a focus on text understanding tasks, as opposed to generative tasks. Prompt tuning (Lester et al., 2021) simplifies Prefix tuning by introducing trainable tokens, called *soft prompts*, for each downstream task. Liu et al.

(2021a) introduced P-tuning v2 which uses deep prompt tuning to address the lack of performance gain in the previous prompt tuning techniques.

By fine-tuning a small fraction of additional parameters, all PEFT approaches alleviate the issue of extensive computational resource requirements.

## 3 Methodology

### 3.1 Problem Statement

Figure 2 shows the comparison between the current and proposed problem definitions. The general problem can be decomposed into two stages:

**Domain-adaptive Pretraining.** Given a pretrained general LLM  $P_{\Phi}(y|x)$  with its parameters  $\Phi$  and a training dataset  $\mathcal{Z} = \{(x_i, y_i)\}_{i=1, \dots, N}$ . To adapt to the new domain, the model needs to update its weight iteratively from its pretrained state  $\Phi_0$  to  $\Phi = \Phi_0 + \Delta\Phi$ . This process of maximising the objective function can be defined as:

$$\operatorname{argmax}_{\Phi} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log(P_{\Phi}(y_t | x, y_{<t})) \quad 213$$

In the current paradigm, a full fine-tuning process means that the model needs to learn a  $\Delta\Phi$  whose dimension is equal to the entire pretrained parameters  $|\Delta\Phi| = |\Phi_0|$ , which is computationally expensive.

In the proposed paradigm, we tune only small additional parameters  $\theta$  such that  $\Phi = \Phi_0 + \Delta\Phi(\theta)$  whose dimension is very small compared to the original parameters  $|\theta| \ll |\Phi_0|$ . Thus, the training

objective can be redefined as:

$$\operatorname{argmax}_{\theta} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log (P_{\Phi_0 + \Delta\Phi(\theta)}(y_t | x, y_{<t}))$$

In the current paradigm, the outcome of domain-adaptive pretraining would be a clinically-adapted LLM. While in the proposed paradigm, the outcome would be the clinical PEFT component, which can be combined with the untouched pre-trained general LLM for downstream applications.

**Downstream Fine-tuning.** In the current paradigm, the pretrained clinical LLM is fine-tuned to the downstream tasks, such as document classification tasks. Suppose that we have a pretrained clinical LLM  $P_{\Phi, \Theta}$  with its domain-adapted parameters  $\Phi$  and a newly initialised classifier layer  $\Theta$ , as well as a training dataset  $\mathcal{Z} = \{(x_i, y_i)\}_{i=1, \dots, N}$ . We want to maximise a specific loss function, such as a cross-entropy loss:

$$\operatorname{argmax}_{\Phi, \Theta} \frac{1}{N} \sum_{i=1}^N y_i \log (P_{\Phi, \Theta}(x_i))$$

In contrast, in the proposed paradigm, the fine-tuning process only updates the small additional parameters  $\Delta\Phi(\theta)$  and the classifier head  $\Theta$ :

$$\operatorname{argmax}_{\theta, \Theta} \frac{1}{N} \sum_{i=1}^N y_i \log (P_{\Phi + \Delta\Phi(\theta), \Theta}(x_i))$$

In fact, we can also decompose the fine-tuning into an additional "delta-updating" process:

$$\operatorname{argmax}_{\theta, \phi, \Theta} \frac{1}{N} \sum_{i=1}^N y_i \log (P_{\Phi + \Delta\Phi(\theta) + \Delta\Phi(\phi), \Theta}(x_i))$$

Similar to the Domain-adaptive Pretraining stage, the dimensions of the additional parameters  $\theta$  and  $\phi$  are very small compared to the original parameters. By updating only the additional parameters and the classifier head, the proposed paradigm reduces the computational requirements, making it more efficient and feasible, especially for clinical settings that are often resource-constrained.

### 3.2 Clinical LLaMA-LoRA

Clinical LLaMA-LoRA is a LoRA adapter built upon LLaMA (Touvron et al., 2023). Clinical LLaMA-LoRA is domain-adapted to the clinical domain and fine-tuned to the downstream tasks following the proposed procedure shown on the right-hand side of Figure 2.

Dataset	# Class	Multilabel	# Train	# Valid	# Test
LOS	4	✗	30,421	4,391	8,797
MOR	2	✗	33,954	4,908	9,822
PMV	2	✗	5,666	707	706
DIAG	1,266	✓	33,994	4,918	9,829
PROC	711	✓	30,030	4,357	8,681

Table 1: Statistics and types of downstream clinical document classification tasks: length of stay (LOS), mortality (MOR), prolonged mechanical ventilation (PMV), diagnoses (DIAG), and procedures (PROC).

**LLaMA models** In this study, we evaluate two LLaMA models; the 7 billion parameters version of LLaMA (Touvron et al., 2023) and the 7 billion parameters version of PMC-LLaMA (Wu et al., 2023). LLaMA was pretrained with an array of texts from multiple sources, such as English CommonCrawl, Wikipedia, ArXiv, and C4 (Raffel et al., 2020). While, PMC-LLaMA is a domain-adapted LLaMA model that was pretrained on 4.8 million biomedical academic papers from PubMed Central.

**Domain-adaptive Pretraining** Clinical LLaMA-LoRA is trained using a combination of MIMIC-IV de-identified discharge summaries (331,794) and radiology reports (2,321,355), resulting in a collection of 2,653,149 individual clinical notes. We evaluate five different PEFT techniques, which include *LoRA*, *Adaptation Prompt*, *Prefix Tuning*, *Prompt Tuning*, and *P-tuning*.

Our approach follows the autoregressive language modelling pretraining objective employed in the original LLaMA training. To ensure compatibility with available computational resources, we use fixed model hyperparameters that allow us to fit the LLM into a single NVIDIA A100-80GB GPU (see Appendix A.1). We optimise the hyperparameters specific to each PEFT method using Gaussian Process regression for Bayesian Optimisation (Frazier, 2018)<sup>1</sup> with a maximum of 20 trials. The detailed hyperparameters search space can be found in Appendix A.2. During this stage, we evaluate the perplexity scores of the LLM variants.

**Downstream Fine-tuning** We fine-tune the Clinical LLaMA-LoRA and Downstream LLaMA-LoRA to clinical document classification tasks:

- **Prolonged mechanical ventilation (PMV):** a binary classification task to predict whether a patient will require mechanical ventilation for

<sup>1</sup>Specifically, we use the W&B Sweep APIs: <https://docs.wandb.ai/guides/sweeps>

299	more than seven days (Huang et al., 2020; Naik et al., 2022).	notes from the MIMIC-III database. It is initialised from a biomedical language model called BioBERT (Lee et al., 2019), which is pretrained on biomedical research articles.	345
300			346
301	• <b>In-hospital mortality (MOR)</b> : a binary classification task to predict whether a patient will survive during their hospital stay (van Aken et al., 2021; Naik et al., 2022).	• <b>BlueBERT</b> (Peng et al., 2019): BlueBERT is pretrained on clinical notes from the MIMIC-III database and PubMed abstracts starting from the pretrained checkpoint of BERT (Devlin et al., 2019), a general-domain language model.	347
302			348
303			349
304			350
305	• <b>Length of stay (LOS)</b> : a multiclass classification task to predict the length of a patient’s hospital stay, categorised into four time-bins: less than three days, three to seven days, one to two weeks, and more than two weeks (van Aken et al., 2021; Naik et al., 2022).	• <b>CORe</b> (van Aken et al., 2021): CORe is pretrained on clinical notes from the MIMIC-III database and biomedical articles starting from the pretrained checkpoint of BioBERT (Lee et al., 2019).	351
306			352
307			353
308			354
309			355
310			356
311	• <b>Diagnoses (DIAG)</b> : a large-scale multilabel classification task to predict the differential diagnoses associated with a patient, represented by simplified ICD-9 diagnosis codes (van Aken et al., 2021).	• <b>UmlsBERT</b> (Michalopoulos et al., 2021): UmlsBERT is pretrained on clinical notes from the MIMIC-III database starting from the pretrained checkpoint of Bio+ClinicalBERT while modifying the architecture and pretraining objective by incorporating knowledge from the Unified Medical Language System (UMLS) Metathesaurus (Schuyler et al., 1993).	357
312			358
313			359
314			360
315			361
316	• <b>Procedures (PROC)</b> : a large-scale multilabel classification task to predict the diagnostics or treatments administered to a patient, represented by simplified ICD-9 procedure codes (van Aken et al., 2021).		362
317			363
318			364
319			365
320			366
321	The label and split statistics of each dataset can be found in Table 1.	These baseline models have been trained to perform specifically on clinical data, thus providing comparison points for evaluating the performance of the proposed Clinical LLaMA-LoRA in downstream clinical NLP tasks.	367
322			368
323	During this downstream fine-tuning process, we use fixed model hyperparameters to ensure compatibility with the available computational resources, a single NVIDIA A100-80GB GPU (see Appendix B.1). We optimise the hyperparameters specific to each PEFT method using Gaussian Process regression for Bayesian Optimisation with a maximum of 20 trials. The detailed hyperparameters search space of the PEFT method can be found in Appendix B.2.		369
324			370
325			371
326			372
327			373
328			374
329			375
330			376
331			377
332			378
333			379
334			380
335			381
336			382
337			383
338			384
339			385
340			386
341			387
342			388
343			389
344			390
			391
			392
			393
			394
			395
			396
			397
			398
			399
			400
			401
			402
			403
			404
			405
			406
			407
			408
			409
			410
			411
			412
			413
			414
			415
			416
			417
			418
			419
			420
			421
			422
			423
			424
			425
			426
			427
			428
			429
			430
			431
			432
			433
			434
			435
			436
			437
			438
			439
			440
			441
			442
			443
			444
			445
			446
			447
			448
			449
			450
			451
			452
			453
			454
			455
			456
			457
			458
			459
			460
			461
			462
			463
			464
			465
			466
			467
			468
			469
			470
			471
			472
			473
			474
			475
			476
			477
			478
			479
			480
			481
			482
			483
			484
			485
			486
			487
			488
			489
			490
			491
			492
			493
			494
			495
			496
			497
			498
			499
			500

Base Model	PEFT	Trainable Params	Train Perplexity	Test Perplexity	Train Time (h:m:s)
LLaMA	<b>LoRA</b>	<b>8,388,608 (0.12%)</b>	<b>1.858</b>	<b>2.244</b>	<b>21:37:42</b>
	Adaptation Prompt	1,228,830 (0.02%)	2.561	2.865	24:57:17
	Prefix Tuning	5,242,880 (0.08%)	2.815	2.748	20:11:07
	Prompt Tuning	61,440 (0.0009%)	4.846	4.007	23:27:28
	P-tuning	16,093,696 (0.24%)	2.723	3.271	23:49:31
PMC-LLaMA	<b>LoRA</b>	<b>2,097,152 (0.03%)</b>	<b>1.938</b>	<b>2.404</b>	<b>21:32:59</b>
	Adaptation Prompt	1,228,830 (0.018%)	2.374	2.867	23:33:10
	Prefix Tuning	2,621,440 (0.04%)	1.789	2.848	20:13:10
	Prompt Tuning	40,960 (0.0006%)	4.821	4.385	22:25:32
	P-tuning	2,171,392 (0.03%)	3.491	4.572	22:28:15

Table 2: Domain-adaptive Pretraining results of LLaMA and PMC-LLaMA trained on MIMIC-IV clinical notes with a language modelling objective. Lower perplexity scores indicate better language modelling performance. The **boldface row** indicates the model with the lowest perplexity score from each base model variant.

in the subsequent sections. The following experiments in downstream fine-tuning will utilise this pretrained Clinical LLaMA-LoRA.

## 4.2 Downstream results

From the downstream fine-tuning results shown in Table 3, we can decompose the analysis into multiple research questions:

**Can LoRA help fine-tune LLaMA from other domains (general and biomedical) to achieve higher AUROC scores in clinical tasks?** We compare the results obtained by LLaMA and LLaMA + LoRA, as well as PMC-LLaMA and PMC-LLaMA + LoRA, as presented in Table 3. The obtained results consistently demonstrate improved AUROC scores when utilising LoRA across all tasks. The macro-averaged AUROC score of LoRA-equipped LLaMA shows a notable 13.01% increase when compared to the LLaMA-only baseline. Similarly, LoRA-equipped PMC-LLaMA exhibits a 12.2% improvement in macro-averaged AUROC compared to the original PMC-LLaMA. Both LLaMA and PMC-LLaMA, when equipped with LoRA, exhibit significant AUROC score improvements in all tasks except the prolonged mechanical ventilation prediction task, which is proven challenging for all model variants.

Furthermore, the marginal difference in AUROC scores between PMC-LLaMA and the general-domain LLaMA can be attributed to two factors. Firstly, the original LLaMA has been exposed to biomedical concepts during its pretraining, reducing the need for domain-adaptive pretraining to the biomedical domain. Secondly, clinical NLP tasks are challenging, even for biomedical LLMs.

**Can LoRA-equipped LLaMA and PMC-LLaMA perform comparably in comparison to clinically trained LMs?** We compare the AUROC scores obtained by the baseline models, and LoRA-equipped LLaMA and PMC-LLaMA (see Table 3). Among the baseline models, BlueBERT performs the best with a macro-averaged AUROC score of 69.59%. Compared to BlueBERT, both LLaMA and PMC-LLaMA underperform with macro-averaged AUROC scores of 58.61% and 60.51%, respectively. This finding highlights the importance of clinical-specific fine-tuning.

Significant improvements can be observed in LoRA-equipped LLaMA and PMC-LLaMA, with macro-averaged AUROC scores of 71.62% and 72.71%, respectively. We notice considerable improvements in the diagnoses and procedures prediction tasks. For example, LoRA-equipped LLaMA achieves AUROC scores of 78.37% and 87.49% in the diagnoses and procedures prediction tasks, respectively, compared to 73.81% and 77.70% for BlueBERT. This represents improvements of 4.56% in diagnoses prediction and 9.79% in procedures prediction. Improvements are also observed in the results obtained by LoRA-equipped PMC-LLaMA, outperforming BlueBERT by 5% in diagnoses prediction and 9.02% in procedures prediction.

Overall, LoRA-equipped LLaMA and PMC-LLaMA achieve higher AUROC scores than the baseline clinical LMs in various clinical prediction tasks, particularly in diagnoses, procedures, and mortality predictions, while maintaining competitive AUROC scores in length-of-stay prediction. However, LoRA-equipped LLaMA and PMC-LLaMA still underperform in prolonged mechanical ventilation prediction.

Model	PMV	MOR	LOS	DIAG	PROC	Macro Average
<i>BlueBERT</i>	53.12	76.95	66.36	73.81	77.70	69.59
UmlsBERT	55.49	75.87	66.06	64.34	74.19	67.19
Bio+ClinicalBERT	54.49	72.92	65.13	65.97	71.73	66.05
CORe	52.11	71.52	64.17	72.40	72.73	66.59
LLaMA*	51.38	66.80	57.65	60.06	63.83	58.61
+ LoRA	51.65	74.89	65.70	78.37	<b>87.49</b>	71.62
+ Clinical LLaMA-LoRA (Frozen)	51.62	65.66	58.16	63.47	69.01	61.58
+ Downstream LLaMA-LoRA	51.11	66.00	58.04	60.46	65.30	60.18
+ Clinical LLaMA-LoRA (Trainable)	55.76	74.81	64.83	76.07	82.76	70.85
+ <i>Downstream LLaMA-LoRA</i>	<b>56.72</b>	77.36	66.32	78.52	87.15	<b>73.21</b>
PMC-LLaMA*	53.06	66.77	57.94	60.17	64.63	60.51
+ LoRA	53.84	<b>78.03</b>	66.14	78.81	86.68	72.70
+ Clinical LLaMA-LoRA (Frozen)	51.33	67.19	58.13	63.59	68.26	60.06
+ Downstream LLaMA-LoRA	50.90	67.00	58.31	60.50	64.42	60.23
+ Clinical LLaMA-LoRA (Trainable)	52.88	75.86	65.89	<b>79.66</b>	86.85	72.23
+ Downstream LLaMA-LoRA	52.21	76.54	<b>68.42</b>	78.67	87.08	72.58

Table 3: AUROC scores in clinical downstream document classification tasks. The macro-averaged AUROC score is calculated by taking the average of AUROC scores across all tasks. The **boldface cell** indicates the highest AUROC score in a column, the *row in italic* indicates the model variant with the highest macro-averaged AUROC in its category. \* Due to restricted computing resources, the fine-tuning of LLaMA and PMC-LLaMA was constrained to only training the final classification layer.

Model	PMV	MOR	LOS	DIAG	PROC	Macro Average
BlueBERT	53.12	76.95	66.36	73.81	77.70	69.59
+ LoRA	55.77	<b>81.90</b>	<b>70.48</b>	70.66	78.10	71.56
UmlsBERT	55.49	75.87	66.06	64.34	74.19	67.19
+ LoRA	56.59	80.33	69.03	69.68	77.53	70.63
BioClinicalBERT	54.49	72.92	65.13	65.97	71.73	66.05
+ LoRA	56.13	78.81	68.28	68.53	75.19	69.39
CORe	52.11	71.52	64.17	72.40	72.73	66.59
+ LoRA	55.31	79.27	68.18	67.34	72.36	68.49
LLaMA + Clinical LLaMA-LoRA + Downstream LoRA	<b>56.72</b>	77.36	66.32	<b>78.52</b>	<b>87.15</b>	<b>73.21</b>

Table 4: AUROC scores of the LoRA-equipped baseline models in clinical downstream tasks. The **boldface cell** indicates the highest AUROC score in a column. The *row in italic* indicates the model variant with the highest macro-averaged AUROC in its category.

**Can LLaMA and PMC-LLaMA with Clinical LLaMA-LoRA achieve higher AUROC scores than the clinically trained LMs?** The domain-adaptive pretraining step yields the clinically-trained LoRA adapters for LLaMA and PMC-LLaMA, called Clinical LLaMA-LoRA. We compare the results of Clinical LLaMA-LoRA-equipped LLaMA and PMC-LLaMA with the baseline models. We evaluate Clinical LLaMA-LoRA with and without downstream fine-tuning, referred to as "Trainable" and "Frozen" respectively.

The results indicate that Clinical LLaMA-LoRA-equipped LLaMA and PMC-LLaMA outperform the baseline models. LLaMA with a trainable Clinical LLaMA-LoRA achieves an AUROC score of 70.85%, surpassing BlueBERT's score of 69.59%. PMC-LLaMA with a trainable Clinical LLaMA-LoRA achieves an even higher AUROC score of

72.23%. These findings demonstrate that the Clinical LLaMA-LoRA contributes to higher AUROC scores for LLaMA and PMC-LLaMA over clinically trained LLMs.

**Can LLaMA and PMC-LLaMA with Clinical LLaMA-LoRA achieve higher AUROC scores than the other fine-tuning variants?** We examine the importance of the domain-adapted LoRA by comparing the results obtained by LLaMA and PMC-LLaMA equipped with Clinical LLaMA-LoRA against the results of LLaMA and PMC-LLaMA fine-tuning, both original and with LoRA.

Firstly, we evaluate the frozen pretrained Clinical LLaMA-LoRA. Both LLaMA and PMC-LLaMA with frozen Clinical LLaMA-LoRA do not exhibit a significant increase in performance compared to the original fine-tuning. This indi-

498 cates that, despite the domain-adaptive pretraining,  
499 the limited number of trainable parameters during  
500 the downstream fine-tuning restricts the potential  
501 improvement that the model can achieve.

502 This reasoning is further supported by the sig-  
503 nificant improvement observed in the AUROC  
504 scores of LLaMA and PMC-LLaMA with train-  
505 able Clinical LLaMA-LoRA. LLaMA and PMC-  
506 LLaMA with trainable Clinical LLaMA-LoRA  
507 achieve 70.85% and 72.23% macro-averaged AU-  
508 ROC scores, respectively, massive improvements  
509 from the vanilla fine-tuning performance (58.61%  
510 and 60.51% AUROC scores respectively).

511 However, Clinical LLaMA-LoRA does not  
512 yield significant improvements when compared  
513 to LLaMA and PMC-LLaMA, which are directly  
514 equipped with LoRA without pretraining. For in-  
515 stance, we can observe that LLaMA with LoRA  
516 achieves a slightly higher macro-averaged AUROC  
517 score of 71.62% compared to LLaMA with Clinical  
518 LLaMA-LoRA, which achieves 70.85%.

519 **Can a downstream LoRA adapter improve the**  
520 **AUROC scores of LLaMA and PMC-LLaMA**  
521 **equipped with Clinical LLaMA-LoRA?** By  
522 considering Clinical LLaMA-LoRA as the "delta-  
523 updating" outcome of the domain-adaptive pre-  
524 training, we can view the downstream fine-tuning  
525 process as an additional "delta-updating" step.  
526 To investigate the impact of this approach, we  
527 conduct experiments by adding a Downstream  
528 LLaMA-LoRA to LLaMA and PMC-LLaMA  
529 models that were already equipped with Clinical  
530 LLaMA-LoRA. From Table 3, we can observe  
531 that Downstream LLaMA-LoRA fails to improve  
532 the performance of LLaMA and PMC-LLaMA  
533 with frozen Clinical LLaMA-LoRA. On the other  
534 hand, improvement can be observed when adding  
535 Downstream LLaMA-LoRA to LLaMA with train-  
536 able Clinical LLaMA-LoRA. This combination of  
537 LLaMA with trainable Clinical LLaMA-LoRA and  
538 Downstream LLaMA-LoRA achieves the highest  
539 macro-averaged AUROC score of 72.81%. The  
540 macro-averaged AUROC score of Clinical LLaMA-  
541 LoRA was almost similar to that of PMC-LLaMA  
542 with LoRA, suggesting similar efficacy between  
543 Clinical LLaMA-LoRA and the full fine-tuning  
544 process that PMC-LLaMA has undergone. More-  
545 over, Clinical LLaMA-LoRA offers the advantage  
546 of reduced computational resources and training  
547 time, which is aligned with the requirements of  
548 practical implementation in clinical settings.

**Can LoRA help better fine-tune clinically-**  
549 **trained LMs?** The baseline models are relatively  
550 smaller in size compared to the LLaMA-based mod-  
551 els, which may be a better fit to care providers with  
552 limited access to computing resources. To that  
553 end, we experimented with fine-tuning the baseline  
554 models with LoRA. 555

556 Table 4 shows the obtained results. All base-  
557 line models see improvements in AUROC scores  
558 in all tasks. For instance, the LoRA-equipped Blue-  
559 BERT achieves an improved macro-averaged AU-  
560 ROC score of 71.56% compared to the conven-  
561 tional fine-tuning with 69.59%. 562

563 This finding highlights the possibility of using  
564 LoRA to efficiently fine-tune clinically trained  
565 LMs, such as BlueBERT, to downstream use cases.

## 566 5 Conclusions

567 In this study, we propose a two-step PEFT frame-  
568 work. We introduce Clinical LLaMA-LoRA,  
569 a LoRA (Hu et al., 2022) adapter built upon  
570 LLaMA (Touvron et al., 2023). Then, we intro-  
571 duce Downstream LLaMA-LoRA, a task-specific  
572 adapter that is trained on top of the pretrained Clin-  
573 ical LLaMA-LoRA. The fusion of the two adapters  
574 achieves state-of-the-art performance with an AU-  
575 ROC score of 72.81% macro-averaged across  
576 all clinical NLP downstream tasks, which rep-  
577 represents a 3.22% improvement over the previous  
578 best-performing model. Our proposed framework  
579 achieves improvement in performance while reduc-  
580 ing the computational requirements, which is suited  
581 for clinical settings that are often constrained by  
582 their computational power.

583 We also find that the LoRA-equipped BlueBERT  
584 model achieves a considerable improvement of  
585 macro-averaged AUROC score over the full fine-  
586 tuning (71.56% compared to 69.59%), with no-  
587 table improvements in mortality and length-of-stay  
588 prediction. These findings further highlight the  
589 potential to achieve strong performance without  
590 extensive computational resources.

591 Future works may explore developing a schema  
592 to address various real-world use cases, building  
593 upon the findings of this study. Such a schema  
594 would use multiple Downstream LLaMA-LoRA  
595 adapters tailored for different use cases while lever-  
596 aging the pretrained LLM and Clinical LLaMA-  
597 LoRA as the foundation. This solution would also  
598 be suited for use cases that rely on private data  
599 commonly encountered in care provider settings.



## 599 Limitations

600 This study presents a two-step PEFT framework  
601 aimed at effectively adapting LLMs to diverse clinical  
602 downstream applications. However, the evaluation  
603 of our model was restricted to MIMIC-based  
604 datasets, which are constrained to English and obtained  
605 exclusively within the Commonwealth of Massachusetts,  
606 United States of America. Consequently, despite the  
607 promising efficacy demonstrated by our proposed  
608 method, it would have been advantageous to directly  
609 assess its performance across diverse hospital systems  
610 spanning various geographical locations and languages.  
611 This would enable a more comprehensive understanding  
612 of its applicability and generalizability. However, it  
613 is essential to acknowledge that conducting such an  
614 analysis would require working within a trusted  
615 research environment and obtaining the necessary  
616 permissions to access the relevant datasets.  
617

618 It is crucial to recognise the restrictions imposed  
619 on accessing internal clinical datasets, as they limit  
620 our ability to evaluate the effectiveness of our  
621 approach across different care provider systems.  
622 Therefore, we encourage care providers to conduct  
623 internal experiments within their trusted research  
624 environment to ensure the efficacy of our proposed  
625 method within their specific use cases should they  
626 adopt this approach.

627 Despite the demonstrated performance improvements,  
628 the proposed model may still be susceptible to spurious  
629 correlations. Predicting patient outcomes solely  
630 based on clinical notes presents significant challenges  
631 due to the other factors that may not be captured  
632 within those notes. For instance, the length of a  
633 patient’s in-hospital stay is not solely correlated  
634 with their diagnoses and disease progression. Factors  
635 such as the patient’s insurance status, which is not  
636 typically mentioned in clinical notes, can severely  
637 impact the duration of a patient’s stay. Therefore,  
638 we encourage end users of such clinical LLMs to  
639 consider additional measures to ensure predictions  
640 that reflect a holistic view of the patient’s situation,  
641 instead of relying solely on the predictions of LLMs.  
642

## 643 Ethics Statement

644 In this study, we use MIMIC-based datasets obtained  
645 after completing the necessary training. These  
646 datasets comply with de-identification standards  
647 set by the Health Insurance Portability and  
648 Accountability Act (HIPAA) through data cleans-

649 ing. Due to privacy concerns, we refrain from including  
650 direct excerpts of the data in the paper. We also  
651 refrain from publicly sharing the pretrained  
652 checkpoints.

653 While our model demonstrates effectiveness, it is  
654 important to acknowledge the risks associated with  
655 relying solely on clinical outcome prediction models.  
656 There are crucial pieces of information that can be  
657 found beyond the scope of clinical notes. Considering  
658 the potential impact on patient health outcomes,  
659 it is crucial to exercise caution when utilising these  
660 clinical LLMs. Therefore, we propose that the PEFT  
661 adapter generated by our framework, in conjunction  
662 with the pretrained LLM, should be used as an aid  
663 rather than a replacement for trained clinical  
664 professionals.

## 665 References

- 666 Emily Alsentzer, John Murphy, William Boag, Wei-  
667 Hung Weng, Di Jindi, Tristan Naumann, and  
668 Matthew McDermott. 2019. [Publicly available clinical BERT embeddings](#). In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 72–78, Minneapolis, Minnesota, USA. Association for Computational Linguistics. 670
- 671 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
672 Kristina Toutanova. 2019. [BERT: pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics. 680
- 681 Peter I. Frazier. 2018. A tutorial on bayesian optimization. *CoRR*, abs/1807.02811. 684
- 685 Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto  
686 Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng  
687 Gao, and Hoifung Poon. 2022. [Domain-Specific Language Model Pretraining for Biomedical Natural Language Processing](#). *ACM Transactions on Computing for Healthcare*, 3(1):1–23. 690
- 691 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,  
692 Bruna Morrone, Quentin De Laroussilhe, Andrea  
693 Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. [Parameter-efficient transfer learning for nlp](#). In *Proceedings of the 36th International Conference on Machine Learning*, page 2790–2799. PMLR. 696
- 697 Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-  
698 Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu  
699 Chen. 2022. [LoRA: Low-rank adaptation of large language models](#). In *International Conference on Learning Representations*. 700

702	Kexin Huang, Abhishek Singh, Sitong Chen, Edward Moseley, Chih-Ying Deng, Naomi George, and Charolotta Lindvall. 2020. <a href="#">Clinical XLNet: Modeling Sequential Clinical Notes and Predicting Prolonged Mechanical Ventilation</a> . In <i>Proceedings of the 3rd Clinical Natural Language Processing Workshop</i> , pages 94–100, Online. Association for Computational Linguistics.	758
703		759
704		760
705		761
706		
707		762
708		763
709		764
710	Alistair E. W. Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J. Pollard, Benjamin Moody, Brian Gow, Li-wei H. Lehman, Leo A. Celi, and Roger G. Mark. 2023. <a href="#">MIMIC-IV, a freely accessible electronic health record dataset</a> . <i>Scientific Data</i> , 10(1):1.	765
711		766
712		767
713		768
714		769
715		
716	Alistair E.W. Johnson, Tom J. Pollard, Lu Shen, Li-wei H. Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. 2016. <a href="#">MIMIC-III, a freely accessible critical care database</a> . <i>Scientific Data</i> , 3(1):160035.	770
717		771
718		772
719		773
720		
721		
722	Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2019. <a href="#">BioBERT: a pre-trained biomedical language representation model for biomedical text mining</a> . <i>Bioinformatics</i> , 36(4):1234–1240.	774
723		775
724		776
725		777
726		778
727		779
728		780
729		781
730	Eric Lehman and Alistair Johnson. 2023. <a href="#">Clinical-T5: Large Language Models Built Using MIMIC Clinical Text</a> .	782
731		783
732		784
733		785
734		786
735		787
736		788
737	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. <a href="#">The power of scale for parameter-efficient prompt tuning</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	789
738		790
739		791
740		792
741		793
742		794
743		795
744		
745	Xiang Lisa Li and Percy Liang. 2021. <a href="#">Prefix-Tuning: Optimizing Continuous Prompts for Generation</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 4582–4597, Online. Association for Computational Linguistics.	796
746		797
747		798
748		799
749		800
750		801
751		
752		802
753		803
754		804
755		805
756		806
757		807
		808
		809
		810
		811
		812
		813

- 814 Ajay Kumar Tanwani, Heather Cole-Lewis, Stephen  
815 Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble,  
816 Chris Kelly, Nathaneal Schärli, Aakanksha Chowdh-  
817 ery, Philip Andrew Mansfield, Blaise Agüera y Arcas,  
818 Dale R. Webster, Gregory S. Corrado, Yossi Matias,  
819 Katherine Chou, Juraj Gottweis, Nenad Tomasev,  
820 Yun Liu, Alvin Rajkomar, Joelle K. Barral, Christo-  
821 pher Semturs, Alan Karthikesalingam, and Vivek  
822 Natarajan. 2022. [Large language models encode  
823 clinical knowledge](#). *CoRR*, abs/2212.13138.
- 824 Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres,  
825 Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl,  
826 Heather Cole-Lewis, Darlene Neal, Mike Schaecker-  
827 mann, Amy Wang, Mohamed Amin, Sami Lachgar,  
828 Philip Andrew Mansfield, Sushant Prakash, Bradley  
829 Green, Ewa Dominowska, Blaise Agüera y Arcas,  
830 Nenad Tomasev, Yun Liu, Renee Wong, Christo-  
831 pher Semturs, S. Sara Mahdavi, Joelle K. Barral,  
832 Dale R. Webster, Gregory S. Corrado, Yossi Matias,  
833 Shekoofeh Azizi, Alan Karthikesalingam, and Vivek  
834 Natarajan. 2023. [Towards expert-level medical ques-  
835 tion answering with large language models](#). *CoRR*,  
836 abs/2305.09617.
- 837 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier  
838 Martinet, Marie-Anne Lachaux, Timothée Lacroix,  
839 Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal  
840 Azhar, Aurélien Rodriguez, Armand Joulin, Edouard  
841 Grave, and Guillaume Lample. 2023. [Llama: Open  
842 and efficient foundation language models](#). *CoRR*,  
843 abs/2302.13971.
- 844 Betty van Aken, Jens-Michalis Papaioannou, Manuel  
845 Mayrdorfer, Klemens Budde, Felix Gers, and Alexan-  
846 der Loeser. 2021. [Clinical Outcome Prediction from  
847 Admission Notes using Self-Supervised Knowledge  
848 Integration](#). In *Proceedings of the 16th Conference  
849 of the European Chapter of the Association for Com-  
850 putational Linguistics: Main Volume*, pages 881–893,  
851 Online. Association for Computational Linguistics.
- 852 Alex Wang, Amanpreet Singh, Julian Michael, Felix  
853 Hill, Omer Levy, and Samuel R. Bowman. 2019.  
854 [GLUE: A multi-task benchmark and analysis plat-  
855 form for natural language understanding](#). In *7th In-  
856 ternational Conference on Learning Representations,  
857 ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*.  
858 OpenReview.net.
- 859 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien  
860 Chaumond, Clement Delangue, Anthony Moi, Pier-  
861 ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,  
862 Joe Davison, Sam Shleifer, Patrick von Platen, Clara  
863 Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le  
864 Scao, Sylvain Gugger, Mariama Drame, Quentin  
865 Lhoest, and Alexander M. Rush. 2020. [Transform-  
866 ers: State-of-the-art natural language processing](#). In  
867 *Proceedings of the 2020 Conference on Empirical  
868 Methods in Natural Language Processing: System  
869 Demonstrations*, pages 38–45, Online. Association  
870 for Computational Linguistics.
- 871 Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang,  
872 and Weidi Xie. 2023. [Pmc-llama: Further finetuning  
873 llama on medical papers](#). *CoRR*, abs/2304.14454.
- Xi Yang, Aokun Chen, Nima M. Pournajatian,  
Hoo Chang Shin, Kaleb E. Smith, Christopher  
Parisien, Colin Compas, Cheryl Martin, Anthony B.  
Costa, Mona G. Flores, Ying Zhang, Tanja Magoc,  
Christopher A. Harle, Gloria P. Lipori, Duane A.  
Mitchell, William R. Hogan, Elizabeth A. Shenkman,  
Jiang Bian, and Yonghui Wu. 2022. [A large lan-  
guage model for electronic health records](#). *npj Digit.  
Medicine*, 5.

## A Hyperparameters for the Domain-adaptive Pretraining

### A.1 Fixed Model Hyperparameters

Hyperparameter	Value
Learning rate	3e-4
Warmup steps ratio	0.06
Maximum sequence length	128
Gradient accumulation step	4
Batch size	10

Table 5: Fixed model hyperparameters for language modelling pretraining. These hyperparameters remain unchanged to fit LLaMA into a single GPU.

### A.2 PEFT Hyperparameters Optimisation Search Space

PEFT	Hyperparameter	Search space
LoRA	r	[2, 4, 8, 16]
	alpha	[4, 8, 16, 32]
	dropout	[0.0, 0.1, 0.2]
Prefix Tuning	num virtual tokens	[1, 5, 10, 15, 20]
	prefix projection	[true, false]
Prompt Tuning	num virtual tokens	[1, 5, 10, 15, 20]
	prompt init	[text, random]
P-Tuning	num virtual tokens	[1, 5, 10, 15, 20]
	reparameterisation	["MLP", "LSTM"]
	hidden size	[64, 128, 256, 768]
	num layers	[1, 2, 4, 8, 12]
	dropout	[0.0, 0.1, 0.2]
Adaptation Prompt	adapter length	[5, 10]
	adapter layers	[10, 20, 30]

Table 6: The search space for PEFT Hyperparameters optimisation runs during the domain adaptation fine-tuning with language modelling objective. Each PEFT technique has a specific set of hyperparameters to tune, we selected the combination of hyperparameters which has the lowest perplexity score.

Specifically for Prompt Tuning, we use a common prompt initialisation text "Finish this clinical note:".

## B Hyperparameters for the Downstream Fine-tuning

### B.1 Fixed Model Hyperparameters

Hyperparameter	Value
Learning rate	5e-5
Warmup steps ratio	0.06
Maximum sequence length	128
Gradient accumulation step	10
Batch size	10

Table 7: Fixed model hyperparameters for the clinical downstream fine-tuning. These hyperparameters remain unchanged to fit LLaMA into a single GPU.

## B.2 PEFT Hyperparameters Optimisation Search Space

PEFT	Hyperparameter	Search space
LoRA	r	[2, 4, 8, 16]
	alpha	[4, 8, 16, 32]
	dropout	[0.0, 0.1, 0.2]

Table 8: The search space for PEFT Hyperparameters optimisation runs during the downstream fine-tuning. Each PEFT technique has a specific set of hyperparameters to tune, we selected the combination of hyperparameters which has the highest AUROC score.

## C Training Configurations

We use HuggingFace’s Transformers (Wolf et al., 2020) and PEFT (Mangrulkar et al., 2022) libraries for the experiments. All LLaMA-based models are trained on one NVIDIA A100-80GB GPU, while the baseline models are trained on a single NVIDIA GeForce GTX 1080 Ti-16GB GPU.

## D Artefacts

The pretrained baseline models including BioClinicalBERT (Alsentzer et al., 2019), BlueBERT (Peng et al., 2019), and CORE (van Aken et al., 2021) were released under the Creative Commons designation CC0 1.0 Universal license, whereas UmlsBERT (Michalopoulos et al., 2021) was released under the MIT license. LLaMA (Touvron et al., 2023) was released under a noncommercial license.

MIMIC-III and MIMIC-IV dataset was released under the PhysioNet Credentialed Health Data License 1.5.0 and can only be accessed after one finishes the CITI Data or Specimens Only Research training<sup>2</sup>.

<sup>2</sup><https://physionet.org/about/citi-course/>