
Clinical Prompt Learning with Frozen Language Models

Niall Taylor^{1*} Yi Zhang^{1*} Dan W Joyce^{1,2}
Alejo Nevado-Holgado¹ Andrey Kormilitzin¹

¹Department of Psychiatry, University of Oxford, Oxford, OX3 7JX, UK
²NIHR Oxford Health Biomedical Research Centre, Oxford, OX3 7JX, UK

{first_name}.{last_name}@psych.ox.ac.uk

Abstract

Prompt learning is a new paradigm in the Natural Language Processing (NLP) field which has shown impressive performance on a number of natural language tasks with common benchmarking text datasets in full, few-shot, and zero-shot train-evaluation setups. Recently, it has even been observed that large but frozen pre-trained language models (PLMs) with prompt learning outperform smaller but fine-tuned models. However, as with many recent NLP trends, the performance of even the largest PLMs such as GPT-3 do not perform well on specialized domains (e.g. medical text), and the common practice to achieve State of the Art (SoTA) results still consists of pre-training and fine-tuning the PLMs on downstream tasks. The reliance on fine-tuning large PLMs is problematic in clinical settings where data is often held in non-GPU environments, and more resource efficient methods of training specialized domain models is crucial. We investigated the viability of prompt learning on clinically meaningful decision tasks and directly compared with more traditional fine-tuning methods. Results are partially in line with the prompt learning literature, with prompt learning able to match or improve on traditional fine-tuning with substantially fewer trainable parameters and requiring less training data. We argue that prompt learning therefore provides lower computational resource costs applicable to clinical settings, that can serve as an alternative to fine-tuning ever increasing in size PLMs. Complementary code to reproduce experiments presented in this work can be found at: https://github.com/NtaylorOX/Public_Clinical_Prompt

Index terms— Prompt learning, BERT, transfer learning, clinical decision support, few-shot

1 Introduction

The field of Natural Language Processing (NLP) has seen a surge in the use of deep learning in recent years, partly due to the increased capacity and availability of powerful GPUs and cloud computing globally. Both academic and industry research have subsequently become dominated by the use of large Pretrained Language Models (PLMs), which are typically commercially produced and trained on enormous amounts of text data in a self-supervised manner through language modelling objectives such as Masked Language Modeling (MLM) and next word prediction. Two major PLM families are the bidirectional encoder representations from transformers (BERT) Devlin et al. [2019] which originally had 110 million trainable parameters, and Generative Pre-trained Transformer 3 (GPT-

*These authors contributed equally to this work.

3) Radford et al. [2019], Brown et al. [2020a] and the new Meta’s Open Pre-trained Transformer Language Model (OPT) Zhang et al. [2022], with ~ 175 billion parameters. With these PLMs one can fine-tune on new domains and design downstream tasks with relative ease, often resulting in state of the art results on a number of popular datasets and tasks Devlin et al. [2019], Lester et al. [2021]. However, "out of the box" PLMs typically do not perform well on out-of-domain texts Han et al. [2021]: Thus taking a BERT model trained on non-medical texts and applying it to a niche medical text domain often leads to a lackluster performance Lee et al. [2019], Huang et al. [2019]. Instead domain specific PLMs are often created through continued pre-training on domain specific corpora when available Alsentzer et al. [2019], Peng et al. [2019], Gururangan et al. [2020], Senior et al. [2020], Vaci et al. [2021]. Moreover, to then leverage the knowledge of these domain specific PLMs to achieve a downstream task requires further training of a task-specific module, such as a classification head, attached to the end of the PLM Devlin et al. [2019], Wolf et al. [2020]. Typically downstream task fine-tuning requires further training of all of the PLMs parameters, in addition to the attached task specific head(s).

This fine-tuning approach is suitable when the application domain has an abundance of text data, which in many situations is not feasible. For instance in clinical settings, there are major data privacy issues and consequently large open medical datasets are difficult to produce. On top of this, the written language used in clinical text can differ drastically to that of the same language found in general written texts, and even between clinical institutions Huang et al. [2019], Leaman et al. [2015], Kormilitzin et al. [2021]. Together this makes creating general purpose clinical PLMs quite difficult. Additionally, the NLP community has seen a trend of increasing model size to enhance performance; Microsoft recently produced a monolithic 530 billion parameter model named Megatron for state of the art performance on generative tasks Smith et al. [2022]. Whilst impressive, to utilise such models for specific domains of interest will likely require full or partial fine-tuning, which has the massive computational, financial investment and of course, environmental impacts Bender et al. [2021].

Regardless of the size issues of the PLMs, there is still a real benefit in their application to new domains and downstream tasks through traditional fine-tuning, including the biomedical domain Huang et al. [2019], Alsentzer et al. [2019], van Aken et al. [2021]. The persistent concern is the need to fine-tune both the entire PLM and task specific head to produce viable performance on many tasks. In the case of the recently produced super large PLMs, this can require the continual training of models that require large suites of high end GPUs, with proportional financial costs. GPUs and high-performance computing clusters are rarely available to hospitals and community clinics that hold much of the existing medical data. Further to this, traditional fine-tuning can lead to a very specific fine-tuned model that is now very far from its initial pre-trained state, which may cause catastrophic forgetting of the pretrained knowledge Chen et al. [2020]. Fine-tuning has also been reported to exploit spurious correlations of the smaller domain-specific dataset, damaging its generalizability Gururangan et al. [2018], Niven and Kao [2019]. We have also observed this lack of generalizability in medical text when fine-tuning and then validating across American and British English Hofer et al. [2018]. Considering the limitations outlined above, we recognise there is now a movement in the NLP community back towards resource efficient training regimes and models to avoid the need for full scale domain specific training. One promising strategy is known as prompt learning, which aims to close the design gap between the PLMs training objectives and downstream tasks by reformulating the downstream tasks as language modelling objectives Li and Liang [2021], Liu et al. [2021a]. Prompt learning has evolved from earlier works which have reformulated all NLP downstream tasks as text-to-text tasks Raffel et al. [2020] and more recently using task examples within the input text as a form of prompt in auto-regressive PLMs Brown et al. [2020b]. An exciting direction in the prompt learning research space has been its potential in few-shot or low resource settings, relying on frozen PLMs Tsimpoukelli et al. [2021] instead of fine-tuning them: The number of parameters to train decreases dramatically when using frozen PLMs and thus reduces computational requirements Lu et al. [2021]. The major gap in the literature is in the application of prompt learning to clinical or biomedical datasets, and in particular clinical support tasks.

We explore the suitability and performance of prompt learning applied to clinical classification tasks with a direct comparison to traditional fine-tuning methods in full and few-shot training scenarios. Our primary focus is on the performance of these approaches when using a frozen PLM, which is desirable for many reasons, but primarily the consequent reduction in training cost and computational resources required to adapt to new domains or downstream tasks. Conceptually we are not introducing a new methodology, rather exploring different applications of prompt learning to the biomedical

domain and importantly to clinical tasks, rather than simple natural language probing tasks. We observed that prompt learning strategies can outperform traditional fine-tuning on different clinical tasks in both few-shot and full training scenarios with frozen PLMs. This work can serve as a prompt learning framework for clinical tasks and as a basis for further work in this space.

2 Related Work

Since the summer of 2021 there has been a steady influx of research papers concerning prompt learning for common benchmarking open-NLP datasets such as Stanford Sentiment Treebank v2 (SST2), and the General Language Understanding Evaluation (GLUE) Liu et al. [2021a], Brown et al. [2020b], Sanh et al. [2022], Lester et al. [2021], Liu et al. [2021b], Li and Liang [2021]. The datasets and tasks are standard in the field of NLP, and revolve around natural language understanding (NLU) tasks. The common finding is that prompt learning can reach the performance of traditional fine-tuning, and often outperform in few-shot settings. Although the ability of prompt learning to match performance of traditional fine-tuning seems to scale with PLM size Liu et al. [2021b]. One notable paper has investigated the use of GPT-3 for biomedical text datasets in a few-shot setting, finding a decrease in performance when compared to similar tasks in the standard NLU datasets Moradi et al. [2021]. This suggests that even the largest PLMs cannot be applied directly to specialised domains and expect good performance, and that domain specific PLMs are still sought for optimal results.

Recently, prompt learning was used to investigate the zero-shot performance on a clinical task using different PLMs and manual prompt templates Sivarajkumar and Wang [2022]. They found that biomedically trained PLMs outperformed general PLMs for one task, and we hope to extend these findings by introducing different prompt learning training strategies and clinical tasks.

3 Methods

3.1 Traditional fine-tuning

Conventional fine-tuning can be achieved by adding task-specific layer(s) or entire multi-layer perceptron (MLPs). The exact approach to processing the PLM output is dependent on the task.

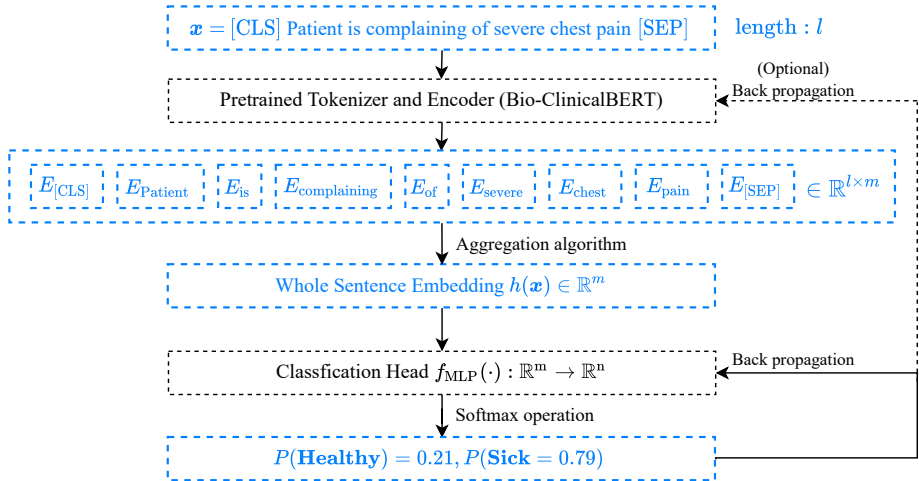


Figure 1: Illustration of conventional fine-tuning method, with an option to freeze the PLM, shown in dotted line. Here [CLS] and [SEP] tokens are special tokens for BERT-based models that are added to the beginning and end of sequences.

In the case of document classification, the downstream task head is an MLP $f_{\text{MLP}}(\cdot)$ which takes the pooled sentence embedding output by the PLM as input and generates an n -dimensional vector, where n is the number of classes. That is, given an input text \mathbf{x} , we first process the raw input with the PLM to get the m -dimensional embedding of each token. Then a pooling operation, such as

the mean, as applied to all token embeddings to produce a singular sentence embedding $h(\mathbf{x})$ of the same dimension m . Then $h(\mathbf{x})$ is fed to the MLP block in a standard feed-forward manner to get the probability across n classes with a softmax operator:

$$P(y | \mathbf{x}) = \frac{\exp((f_{\text{MLP}}(h(\mathbf{x}))_y))}{\exp(\sum_{i=1}^n f_{\text{MLP}}(h(\mathbf{x}))_i)}.$$

The MLP block can have any depth of layers $m \in \mathbb{N}$, while in our experiments, we opted for $d = 2$. Since the additional MLP block and PLMs are modular, their respective parameters are stored separately and we can opt to freeze the parameters of one or the other. An example of processing a short input text sequence using this method is shown in Fig. 1.

3.2 Prompt Learning

Generally, prompt learning can be achieved via the following steps: Given an input text \mathbf{x} , we modify it to a prompt format $\mathbf{x}' = f_p(\mathbf{x})$, where f_p , often called a template, will normally prepend, append, or insert a number of additional token embeddings to the original input along with a masked token, denoted by $\langle[\text{MASK}] \rangle$. We then feed \mathbf{x}' into the PLM to predict the masked token, which is the same as the Masked Language Modelling (MLM) pre-training objective of most BERT-based models. The result of the model will be a distribution over the fixed vocabulary \mathcal{V} of the tokenizer. A second and crucial step is to map tokens or words in the known vocabulary of the PLM to class labels in the downstream task, achieved with a mapping $g : \mathcal{V} \mapsto \mathcal{C}$, where \mathcal{C} is the set of classes. This is known as answer engineering, or verbalization (we will use the term verbalizer and verbalization throughout). The verbalizer can be seen as a mapping between single, or multiple different tokens to distinct class labels. The embedding or hidden state represented at the $\langle[\text{MASK}] \rangle$ position output by the PLM is then passed through a standard language model head, or classifier, and probabilities of the verbalizer derived class label tokens are derived.

A simple prompt-based clinical classification example could be to determine whether a patient has heart disease with class labels as sick and healthy, a prompt learning setup could be as follows: Take the template “ $\langle\text{clinical text}\rangle \langle\text{prompt}=\text{“Patient is”}\rangle \langle[\text{MASK}]\rangle$ ”, where $\langle\text{clinical text}\rangle$ represents the original input text, the $\langle[\text{MASK}]\rangle$ token is the label or class to predict. The verbalizer will map certain tokens to each class of sick and healthy separately, essentially a dictionary mapping e.g. { **“Healthy”**: ‘fine’, and **“Sick”**: ‘unwell’ }. Subsequently if the token predicted at the $\langle[\text{MASK}]\rangle$ position is ‘fine’ then this will be mapped to the **Healthy** class. Thus, the sentence “Patient is complaining of severe chest pain.” will first be wrapped by the pre-defined template as “Patient is complaining of severe chest pain. Patient is $\langle[\text{MASK}]\rangle$ ”. The wrapped sentence is then tokenized and fed into the PLM to predict the distribution over vocabulary on the $\langle[\text{MASK}]\rangle$ token position, although we just care about the probabilities of the tokens (‘fine’ and ‘unwell’) that are mapped to each of the classes that are contained in \mathcal{V} , with “unwell” hopefully having a higher probability to be predicted by the masked language model predictor and the class “sick” ultimately being predicted. We offer an illustration of the basic prompt framework in Fig. 2.

Within the broad prompt learning framework there are important decisions to make about the construction of prompt templates and verbalizers. At its infancy templates were manually created, often based on human knowledge of the task domain, with massive variance in performance with even subtle perturbations of the template and verbalizer Lester et al. [2021], Hu et al. [2021].

To enable a standardised framework for prompt learning a team have developed OpenPrompt to enable reproducible prompt based research by creating a open source and unified code-base Ding et al. [2021]. We shall first define the templates and verbalizers used in the framework and our experiments. We refer to the classical prompt learning strategy with handcrafted templates and verbalizers as manual templates and manual verbalizers respectively. This strategy was first proposed as the Pattern-Exploiting Training (PET) Schick and Schütze [2021]. We denote the set of words in the verbalizer for each class $y \in \mathcal{C}$ to be \mathcal{V}^y . The probability of each class given the input \mathbf{x} and its prompt \mathbf{x}' is thus:

$$P(y | \mathbf{x}) = \frac{\exp\left(\frac{1}{|\mathcal{V}^y|} \sum_{t \in \mathcal{V}^y} P_M(t | \mathbf{x}')\right)}{\sum_{i=1}^{|\mathcal{C}|} \exp\left(\frac{1}{|\mathcal{V}^i|} \sum_{t \in \mathcal{V}^i} P_M(t | \mathbf{x}')\right)}.$$

Manual templates and verbalizers are discrete and bounded to the PLMs vocabulary, so there are no extra parameters to train, although fine-tuning the PLM is possible.

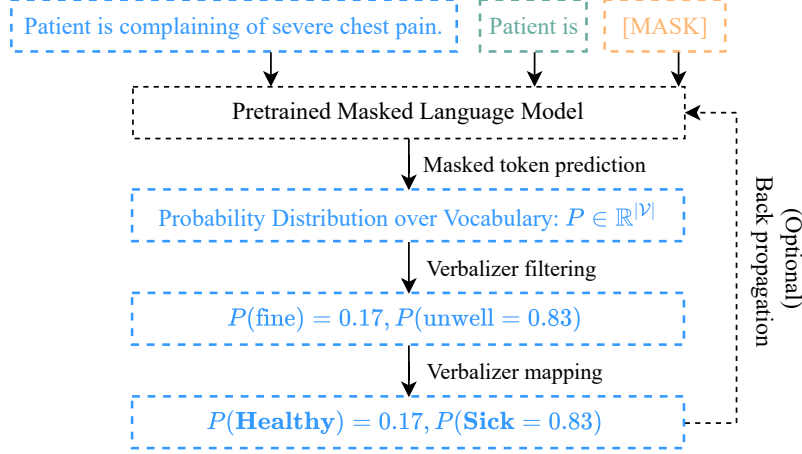


Figure 2: Illustration of manual template and verbalizer in prompt learning.

The engineering of the manual components of prompt learning is not straight forward, with large variations in performance emerging from small changes to the tokens, and typically domain expertise is required. One can however sacrifice the human interpretability of the manual components and create trainable or soft prompt components. Soft prompt learning operates in the same manner as manual approach, but replaces the fixed manual components with trainable embeddings (continuous vectors) of same dimension as the original PLM. The error from the downstream task can then be back-propagated to tune only the embeddings for the template and verbalizer Lester et al. [2021]. Normally, a manual template has the form of $\mathbf{x}' = \{[P_0, P_1, \dots, P_j], \mathbf{x}, [P_{j+1}, P_{j+2}, \dots, P_k], [\text{MASK}]\}$, where for $i \in \{0, 1, \dots, k\}$, P_i denotes the token of the template. And since \mathbf{x}' is fed to the PLM to get $h(\mathbf{x}')$, the prompt tokens P_i are also mapped to the embedding space, where we can assume $h(P_i)$ to be optimized during training and such tokens are denoted as $\langle[\text{soft}] \rangle$ in the template format. A template where all tokens are $\langle[\text{soft}] \rangle$ is called a soft template, while a template with a mixture of manual and $\langle[\text{soft}] \rangle$ tokens is called a mixed template.

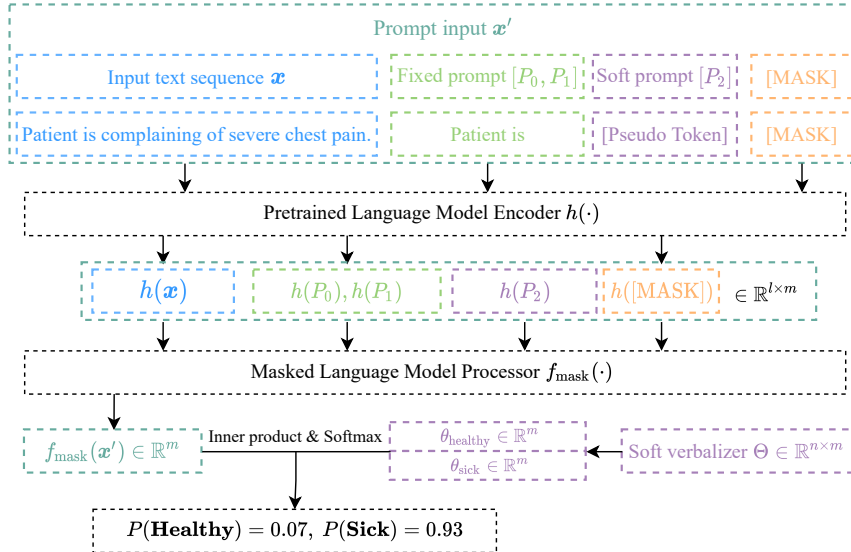


Figure 3: Illustration of soft template and verbalizer in prompt learning. If the $\langle[\text{soft}] \rangle$ token P_2 is not defined manually in advance, the embedding $h(P_2) \in \mathbb{R}^m$ will be randomly initialized in the hidden space.

Similarly, a soft verbalizer can be assumed as replacing words in verbalizer with trainable vectors for each class. Therefore, when using the soft verbalizer, there is no need to build the mapping from

vocabulary \mathcal{V} to class labels \mathcal{C} as the trainable vectors do not have semantic meaning. The resulting verbalizer then becomes a matrix operator $\Theta \in \mathbb{R}^{n \times m}$, where n represents the number of classes and m represents the dimension of generated hidden embeddings. For better understanding, we denote the i -th row of Θ as θ_i for each trainable vector of class i . To compile with the soft verbalizer which takes hidden embeddings from the PLM as input, the original decoder head of the PLM is removed. We denote the resulting mapping from $h(\mathbf{x}') \in \mathbb{R}^{l \times m}$ to the prediction of hidden representation of $\langle [\text{MASK}] \rangle$ as $f_{\text{mask}} : \mathbb{R}^{l \times m} \rightarrow \mathbb{R}^m$, where l is the sequence length of \mathbf{x}' . Therefore, the probability of class y given the input \mathbf{x} and its prompt \mathbf{x}' can be calculated by

$$P(y | \mathbf{x}) = \frac{\exp(\theta_y^\top f_{\text{mask}}(h(\mathbf{x}')))}{\sum_{i=1}^n \exp(\theta_i^\top f_{\text{mask}}(h(\mathbf{x}')))}.$$

For further details and origins of prompt learning see: P-tuning Liu et al. [2021c], prefix tuning Li and Liang [2021] and WARP Hambarzumyan et al. [2021].

3.3 Pre-trained Language Model

As we wanted to compare the performance of prompt learning and traditional fine-tuning in a best case scenario, we chose the Bio-ClinicalBERT Alsentzer et al. [2019]. Bio-ClinicalBERT was essentially pre-trained on all MIMIC-III notes and a large collection of PubMed abstracts and full articles by being initialized from weights produced by another biomedical BERT model, BioBERT Lee et al. [2019]. Whilst we appreciate this may be an overly optimized model for the dataset used in this paper, we argue the point of the experiments presented here is to compare and contrast the ability of the different modelling frameworks to leverage what has been learned by a PLM for clinical tasks. As has already been shown extensively, PLMs benefit from domain specific pre-training Gururangan et al. [2020], what is lesser known is whether current pre-prompt learning approaches are fully utilising these language models.

3.4 Clinical Dataset

We use the Medical Information Mart for Intensive Care III (MIMIC-III) [Johnson et al., 2016], an open source medical dataset developed by the MIT Lab for Computational Physiology. It comprises of de-identified health data associated with 38,597 critical care patients and 58,976 intensive care unit (ICU) admissions at the Beth Israel Deaconess Medical Center between 2001 and 2012. Data includes demographics, vital signs, laboratory tests, medications, caregiver notes, imaging reports, and mortality in and out of hospital. The number of possible tasks with this dataset is quite large and varied, but we focus on classification tasks which utilise free text notes alone. Moreover, to allow comparisons with other baselines we derive clinical task datasets used in previous research van Aken et al. [2021], Pellegrini et al. [2022], Wang et al. [2020], Boag et al. [2018] as well as deriving our own triage task, described below. An important note is that whilst some of the derived clinical tasks may benefit from utilising the multi-modal data available for each patient, we focus purely on the free text clinical notes. Full details and code for reproducing these datasets and experiments is provided by authors.²

4 Experiments - Clinical tasks

ICD-9 50 Within the MIMIC-III data and other EHRs are standardised International Classification of Diseases version 9 (ICD-9) codes, which are used to record diagnosis and procedures. A common task is to classify the ICD-9 diagnosis code based on a patients data and automate the whole process, and one can do so from the free text notes alone. There are approximately 2,000 diagnosis codes present in the MIMIC-III dataset, with a very skewed distribution, and a resulting extreme multi-class problem which is beyond the scope of this paper. Thus for our classification task we opt to subset top 50 most frequent ICD-9 diagnosis codes that have a corresponding set of clinical notes, as has been done before Yuan et al. [2022], Wang et al. [2020], van Aken et al. [2021].

²complementary code to reproduce experiments is provided at: https://github.com/Ntaylor0X/Public_Clinical_Prompt

ICD-9 Triage task A potential concern with the ICD-9 diagnosis code classification is that the codes themselves may be mentioned explicitly in the notes van Aken et al. [2021]³, and further, simply classifying patients’ ICU discharge notes by ICD-9 code lacks ecological validity as a clinical decision support task. For example, within a hospital setting, patients admitted to an ICU will be treated and then “stepped down” (discharged) to another ward or team to progress their treatment when they no longer require ICU. With assistance from clinicians, we therefore designed a novel task that aims to make the classification task more similar to the decision making process of arranging patient flow on discharge from the ICU. For example, a patient being discharged from the ICU after a cardiac event will likely be “stepped down” to a cardiology team. Similarly, a patient admitted to ICU with obstetric complications will likely be stepped-down to a maternity ward. In essence we grouped together the ICD-9 diagnosis codes into “teams” that reflect the triage or patient-flow decision making found in hospital settings.

For this task we selected the top 20 most frequent ICD-9 diagnosis codes in MIMIC-III and a clinician derived triage groups based on which team would likely continue the patient’s care on being stepped down from ICU. The training classes are therefore many-to-one mappings of ICD-9 codes to discharge teams and we derived the following seven post-ICU discharge destination teams: Cardiology, Obstetrics, Respiratory Medicine, Neurology, Gastroenterology, Acute or Internal Medicine, and Oncology. The resultant dataset consists of 15,000 clinical notes across the 7 triage categories.

In hospital mortality One of the most frequently used benchmark clinical support tasks with the MIMIC-III dataset is the prediction of whether a patient will survive their hospital episode. Within the MIMIC-III database are structured data relating to the mortality status of a patient, which paired with a date and timestamp allows for easy labelling of the data. Only notes prior to the mortality flag are considered, and some simple regular expression rules were used to filter any notes that had explicit mentions of a patients death, similar to that of previous work Boag et al. [2018], van Aken et al. [2021].

Length of stay in ICU Predicting how long a patient will require ICU is of significant value to hospitals who aim to optimise the flow of patients in resource-limited settings (that is, there are usually very few ICU beds compared to the hospital’s overall bed capacity). We model this as a three way classification task, binning length of stay in the following categories: Under 3 days, 3 to 7 days, 1 week to 2 weeks, more than 2 weeks van Aken et al. [2021].

Full and few-shot training We will be comparing the performance of models in full and few-shot training setups. Validation and test set performance is always carried out on the full validation and test sets to enable direct comparisons in performance. An important note for our few-shot experiments is that sample size will refer to the number of samples per class, i.e. $N = s \times c$ where N is the total training samples, s is the sample size per class and c is the number of unique classes. Note in some instances not all classes can fill the sample size, so for some few-shot experiments there will remain a class imbalance. All results presented are on held-out test sets for each task.

5 Results

5.1 Different prompt learning setups

The number of possible combinations of templates and verbalizers in the prompt learning framework is vast, and as such we have opted to utilise previous research to derive the most suitable for our use case. To this end we conducted an initial experiment comparing the performance of four prompt learning combinations on one clinical task to establish the best performing combination. We chose the ICD-9 Triage task as the baseline due to it being a relatively straight forward multi-class classification problem and with a reasonably balanced distribution of classes when compared to the other tasks. The prompt learning setup comprised six combinations of a manual, mixed or soft template with a manual or soft verbaliser. The results are summarised in Table 1

The performance across the different prompt combinations is very similar in the setting where the PLM is fine-tuned, however there is greater variance when the PLM is frozen. The frozen PLM setting

³it was shown samples where diagnosis was not mentioned explicitly only had a slight drop in performance

Table 1: Table comparing different prompt learning setups on ICD9 Triage task.

PLM	Prompt combination	Balanced accuracy
Fine-tuned	(manual, manual)	0.8765
	(manual, soft)	0.8818
	(mixed, manual)	0.8817
	(mixed, soft)	0.8824
	(soft, manual)	0.8860
	(soft, soft)	0.8954
Frozen	(manual, soft)	0.7524
	(mixed, manual)	0.8474
	(mixed, soft)	0.8724
	(soft, manual)	0.8591
	(soft, soft)	0.8900

is of most interest, and whilst the soft template and soft verbalizer combination performs the best overall, we opt to use the more interpretable combination of mixed template and soft verbalizer as our prompt learning benchmark going forward. The mixed template is a mixture of manual prompting and prefix tuning, whereby both discrete tokens known to the PLM and newly introduced, trainable continuous vectors of the same dimension as the PLM token embeddings are combined.

5.2 Prompt learning versus traditional fine-tuning

Next is a comparison across the different clinical tasks outlined in the methods section between prompt learning and traditional fine-tuning. Each framework utilises the exact same PLM and we present evaluation results for both fine-tuning and freezing the entire PLM. In the case of the frozen PLM, only the parameters introduced by traditional fine-tuning or prompt learning are updated during training. We found that prompt learning can match or improve on traditional fine-tuning, with a much smaller gap in performance between the frozen and fine-tuned PLM setting across few-shot and full training setups, see Fig. 4.

5.3 Hyperparameter search

There are considerable variations in any neural networks performance with changes to hyperparameters, in particular learning rates and hidden layer dimensions. With comparing the performance of two neural network frameworks as we have, one must be careful to ensure the hyperparameters are optimized for each. Our initial experiments used sensible hyperparameters based on previous research using traditional fine-tuning and prompt learning, where prompt learning and traditional fine-tuning achieved similar performance when the PLM was fine-tuned, see Fig.4. However, when freezing the PLMs, performance differences arose between the two frameworks, especially for few-shot settings in favor of prompt learning. We chose the ICD-9 Triage task as the optimal showcase task for further exploration due to its relatively stable performance. Moreover, with limited computational resources, it was impractical to run hyperparameter searches for all tasks and frameworks. The hyperparameter search space is provided below in Table 2, with results of the subsequent optimized training runs for the ICD-9 Triage task presented in Table 3. Further details of the hyperparameter search and results are presented in supplementary materials, see Appendix A.

5.4 Sensitivity analyses

Results suggested that on certain tasks prompt learning outperformed the traditional fine-tuning model when using a frozen PLM Fig.4. We will focus on the triage task again, for which we optimized each of the frameworks. There is a risk that the performance drop for the traditional fine-tuning classification head is due to over or under fitting with its larger number of trainable parameters in the original setting. We manipulated the number of trainable parameters in each framework and compared the effects on performance, for results see Fig.5. Adjusting the number of trainable parameters for traditional fine-tuning involves adjusting the number of layers and hidden dimension size of the classification head, whilst adjusting number of trainable parameters for prompt learning requires

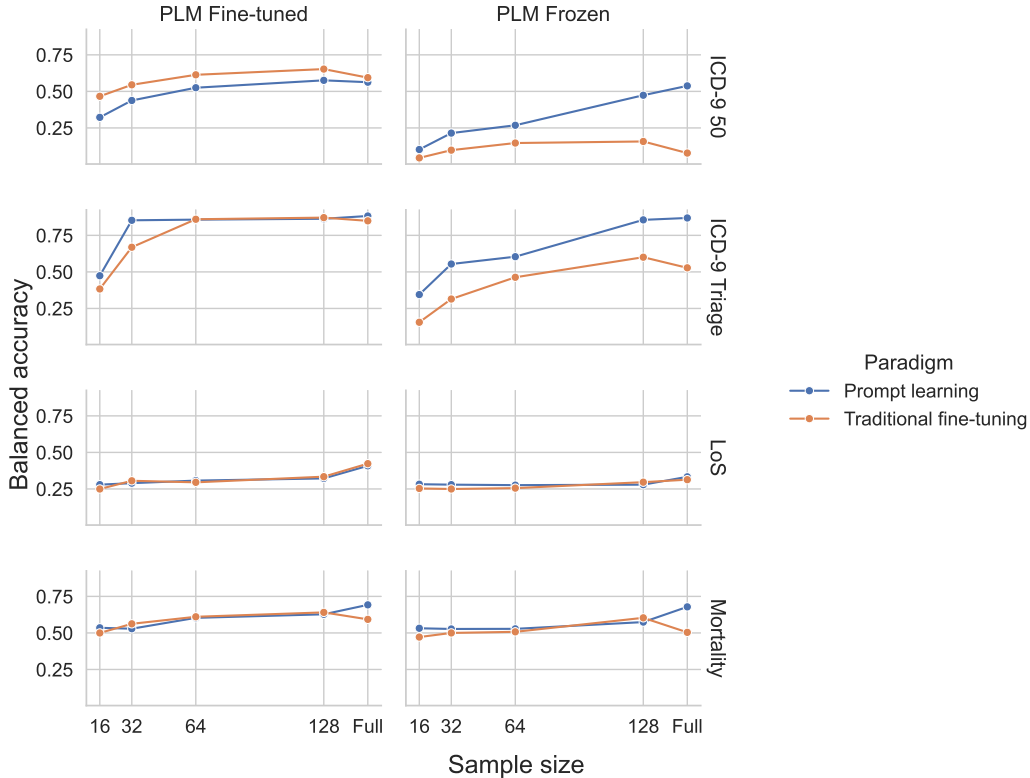


Figure 4: Balanced accuracy for prompt learning and traditional fine-tuning frameworks across the four clinical tasks. “LoS” refers to length of stay and “Full” refers to a full data set size which varies from task to task.

Table 2: Hyperparameter search space used for optimization

Parameter	Search space
classifier learning rate	$\log.\text{uniform}[1 \times 10^{-5}, 3 \times 10^{-1}]$
batch size	[4]
gradient accumulation steps	$\text{range}[2, 10]$
dropout	$\text{range}[0.1, 0.5]$
optimizer	categorical[adamw, adafactor]
prompt learning rate	$\log.\text{uniform}[1 \times 10^{-5}, 3 \times 10^{-1}]$
verbalizer learning rate	$\log.\text{uniform}[1 \times 10^{-5}, 1 \times 10^{-1}]$

Table 3: Hyperparameter optimized model comparison with frozen PLM for ICD9 triage.

Paradigm	Balanced accuracy	F1 weighted	AUC
Traditional fine-tuning	0.8162	0.8919	0.9811
Prompt learning	0.8698	0.9246	0.9889

just changing the number of soft template tokens and whether to include a soft verbalizer (manual templates and verbalizers have no trainable parameters). Training used 128 samples per class as this approached peak performance without requiring a full training run. Note that prompt learning with the *fewest* trainable parameters (N params = 1,536) achieves comparable performance to the traditional fine-tuning model with 1000 times the number of trainable parameters (N params = 1,552,007).

The variability in prompt learning performance based on the template and verbalizer has been well established Liu et al. [2021a], Li and Liang [2021], Ding et al. [2021]. We opted to focus on the use

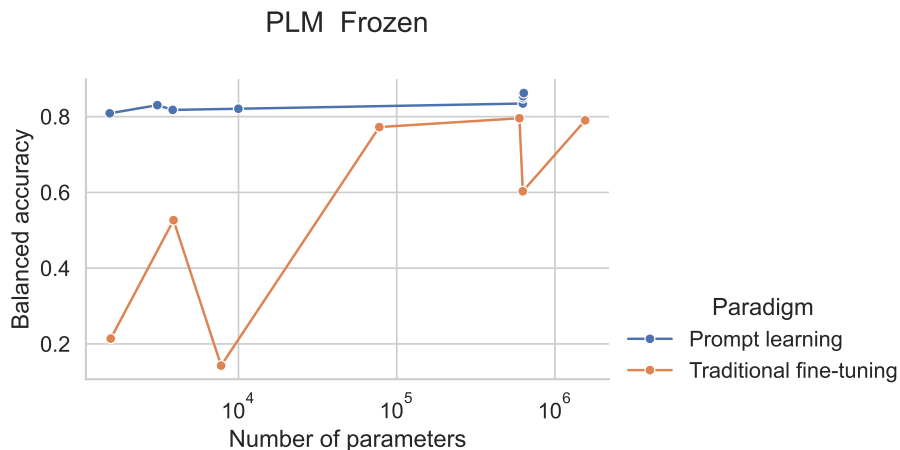


Figure 5: Balanced accuracy for prompt learning versus traditional fine-tuning across increasing number of trainable parameters with frozen PLM. For readability, logarithmic scale is used for x -axis.

of a mixed template format which is based around designing a common sense manual template for the task alongside soft and trainable tokens or embeddings. Moreover these soft tokens can be initialised from a known token of the PLM’s vocabulary. To determine whether mixed templates benefit from a common sense or domain specific manual template, we compared performance of different templates including one with a mix of unrelated and random tokens. Results are shown in Table 4 and we can see that having just one soft token or a set of random and unrelated manual tokens leads to a drop in performance. The `<[soft]>` token represents the trainable continuous vector or embedding of the mixed template that has been initialised from the PLMs vocabulary. Thus `<[soft]>:"This"` indicates a soft embedding initialised from the PLMs representation of the token "This".

Table 4: Performance of the classification model on a test set for different mixed templates for the ICD9 triage task.

Prompt text	Balanced accuracy
<code><[soft]>: "This" <[MASK]></code>	0.8195
<code><[soft]>: "This" patient <[soft]>:"should go to" <[MASK]>.</code>	0.8539
<code><[soft]>: "This" patient should <[soft]>:"go to" <[MASK]>.</code>	0.8491
<code><[soft]>: "This" patient should <[soft]>:"go to this medical team based on symptoms of their illness" <[MASK]>.</code>	0.8624
random words here <code><[soft]>:"random" <[MASK]>.</code>	0.8346

6 Discussion

The experiments presented here have attempted to directly compare the prompt learning paradigm with the traditional fine-tuning paradigm across a number of clinical tasks that frame classification as a clinical decision support task. The objective was to ascertain whether the literature describing promising performance for prompt learning in general domain text datasets can be leveraged on a more niche biomedical domain. We present four clinical decision tasks of varying complexity, in both full training and few-shot setups. In the full training scenario, prompt learning can typically match the performance of traditional fine-tuning, and prompt learning outperforms traditional fine-tuning in the few-shot setting. Of particular interest was the performance of each model with frozen the PLM, where only parameters added to the PLM after pre-training are tuned for downstream classification

tasks. This is where prompt learning appears to prove superior, out-performing traditional fine-tuning with considerably fewer trainable parameters, see Figure .5. Moreover, the use of a mixed template appears to allow the intuitive common sense approach to domain derived prompts, whilst maintaining a trainable soft embedding that can reduce the difficulty in finding optimal manual prompts. We argue that mixed templates achieve similar performance to entirely soft templates, whilst retaining a level of transparency and interpretability. Understanding how models arrive at a decision is especially important in high-stake applications, such as medicine Taylor et al. [2021], Rajpurkar et al. [2022]. Future work should focus on the utility of interpretable prompts for helping clinicians understand a model’s decision making.

6.1 Limitations

Pre-training data leakage A notable limitation was the choice of PLM, which is arguably too well suited to the clinical tasks presented, with probable data leakage from initial pre-training and the subsequent downstream tasks. Although it must be stated that this would have benefited both paradigms, but there is the possibility that the reformulation of the downstream tasks as a masked language modelling style objective may allow easier "remembering" for prompt learning when compared to traditional fine-tuning. However, we include results for the ICD-9 Triage task using biomedical BERT (trained only on biomedical literature) and this yielded a similar pattern of results, see Appendix D.

Task performance variance We presented four clinical tasks derived from MIMIC-III notes data, and whilst we achieved results in line with previous research, the relative performance on the length of stay and mortality prediction tasks were quite poor regardless of the framework. This limits the interpretability of framework differences in performance, and whether one is more suitable to some tasks than others. Similarly we did find that using hyperparameter search for the ICD-9 Triage task improved the frozen PLM performance of the traditional fine-tuning approach by a reasonable margin and a more extensive hyperparameter search may shift this further. However, this was also true for the prompt learning approach, but these models appeared far more robust to changes in hyperparameters. Future work would benefit from exploring this more extensively, given adequate computing resource.

6.2 Conclusion

The key finding was that prompt learning outperforms the traditional fine-tuning approach when PLMs are frozen during training on the downstream task. Most striking is the relatively few trainable parameters required for prompt learning to converge and match or even outperform traditional fine-tuning. This is in line with previous prompt learning research and may offer a useful framework for building clinical support tools in low compute resource settings, as well as enabling a faster, flexible, modular training pipeline for new downstream tasks and novel data. The ability to utilise a single, frozen PLM and share or reuse these embeddings across a number of task specific modules, each with their own trainable prompt is very desirable for specialised domains. Whilst using smaller PLMs and prompts may not achieve the state of the art performance on certain tasks, it can approach similar levels of performance with a fraction of the model size and training time. In the field of clinical support tools, a computationally efficient and interpretable model with good enough performance that can run on a CPU is arguably more desirable than a trillion parameter model that requires high-performance computing clusters with arrays of GPUs. The prompt learning framework is an evolving paradigm with variants being introduced regularly, thus we cannot claim to have fully covered prompt learning in this work. We have opted to use the most readily available, and arguably resource efficient prompt approach to achieve our results. This work can act as a basis for further clinical prompt learning work, and may encourage the use of relatively small domain specific PLMs rather than relying on the giant PLMs produced by commercial enterprises. We suggest that it is more efficient to train a small BERT model on a specialised domain and applying prompt learning, than attempting to apply prompt learning directly to models such as GPT-3 which often lack the domain knowledge required.

Acknowledgement

NT is supported by the EPSRC Center for Doctoral Training in Health Data Science (EP/S02428X/1). AK, ANH, YZ and DWJ were supported in part by the NIHR AI Award for Health and Social Care (NIHR-AI-AWARD0-2183); AK and ANH declare a research grant from GlaxoSmithKline. DWJ is supported by the NIHR Oxford Health Biomedical Research Centre (grant BRC-1215-20005). The views expressed are those of the authors and not necessarily those of the UK National Health Service, the NIHR, the UK Department of Health, or the University of Oxford.

References

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 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, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020a.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.243. URL <https://aclanthology.org/2021.emnlp-main.243>.
- Wenjuan Han, Bo Pang, and Ying Nian Wu. Robust transfer learning with pretrained language models through adapters. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 854–861, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-short.108. URL <https://aclanthology.org/2021.acl-short.108>.
- Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *CoRR*, abs/1901.08746, 2019. URL <http://arxiv.org/abs/1901.08746>.
- Kexin Huang, Jaan Altosaar, and Rajesh Ranganath. Clinicalbert: Modeling clinical notes and predicting hospital readmission, 2019. URL <https://arxiv.org/abs/1904.05342>.
- Emily Alsentzer, John Murphy, William Boag, Wei-Hung Weng, Di Jindi, Tristan Naumann, and Matthew McDermott. Publicly available clinical BERT embeddings. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 72–78, Minneapolis, Minnesota, USA, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-1909. URL <https://aclanthology.org/W19-1909>.
- Yifan Peng, Shankai Yan, and Zhiyong Lu. Transfer learning in biomedical natural language processing: An evaluation of BERT and ELMo on ten benchmarking datasets. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, pages 58–65, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-5006. URL <https://aclanthology.org/W19-5006>.

- Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. Don't stop pretraining: Adapt language models to domains and tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.740. URL <https://aclanthology.org/2020.acl-main.740>.
- Morwenna Senior, Matthias Burghart, Rongqin Yu, Andrey Kormilitzin, Qiang Liu, Nemanja Vaci, Alejo Nevado-Holgado, Smita Pandit, Jakov Zlodre, and Seena Fazel. Identifying predictors of suicide in severe mental illness: a feasibility study of a clinical prediction rule (oxford mental illness and suicide tool or oxmis). *Frontiers in psychiatry*, 11:268, 2020. doi: 10.3389/fpsy.2020.00268. URL <https://www.frontiersin.org/article/10.3389/fpsy.2020.00268>.
- Nemanja Vaci, Ivan Koychev, Chi-Hun Kim, Andrey Kormilitzin, Qiang Liu, Christopher Lucas, Azad Dehghan, Goran Nenadic, and Alejo Nevado-Holgado. Real-world effectiveness, its predictors and onset of action of cholinesterase inhibitors and memantine in dementia: retrospective health record study. *The British Journal of Psychiatry*, 218(5):261–267, 2021. doi: 10.1192/bjp.2020.136.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-demos.6. URL <https://aclanthology.org/2020.emnlp-demos.6>.
- Robert Leaman, Ritu Khare, and Zhiyong Lu. Challenges in clinical natural language processing for automated disorder normalization. *Journal of Biomedical Informatics*, 57:28–37, 2015. ISSN 1532-0464. doi: <https://doi.org/10.1016/j.jbi.2015.07.010>. URL <https://www.sciencedirect.com/science/article/pii/S1532046415001501>.
- Andrey Kormilitzin, Nemanja Vaci, Qiang Liu, and Alejo Nevado-Holgado. Med7: a transferable clinical natural language processing model for electronic health records. *Artificial Intelligence in Medicine*, 118:102086, 2021. ISSN 0933-3657. doi: <https://doi.org/10.1016/j.artmed.2021.102086>. URL <https://www.sciencedirect.com/science/article/pii/S0933365721000798>.
- Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhun Liu, Shrimai Prabhumoye, George Zerveas, Vijay Korthikanti, Elton Zhang, Rewon Child, Reza Yazdani Aminabadi, Julie Bernauer, Xia Song, Mohammad Shoeybi, Yuxiong He, Michael Houston, Saurabh Tiwary, and Bryan Catanzaro. Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model, 2022. URL <https://arxiv.org/abs/2201.11990>.
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623, 2021.
- Betty van Aken, Jens-Michalis Papaioannou, Manuel Mayrdorfer, Klemens Budde, Felix Gers, and Alexander Loeser. Clinical outcome prediction from admission notes using self-supervised knowledge integration. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 881–893, Online, April 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.eacl-main.75. URL <https://aclanthology.org/2021.eacl-main.75>.
- Sanyuan Chen, Yutai Hou, Yiming Cui, Wanxiang Che, Ting Liu, and Xiangzhan Yu. Recall and learn: Fine-tuning deep pretrained language models with less forgetting. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7870–7881, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.634. URL <https://aclanthology.org/2020.emnlp-main.634>.
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A. Smith. Annotation artifacts in natural language inference data. In *Proceedings of the*

- 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pages 107–112, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2017. URL <https://aclanthology.org/N18-2017>.
- Timothy Niven and Hung-Yu Kao. Probing neural network comprehension of natural language arguments. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4658–4664, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1459. URL <https://aclanthology.org/P19-1459>.
- Maximilian Hofer, Andrey Kormilitzin, Paul Goldberg, and Alejo Nevado-Holgado. Few-shot learning for named entity recognition in medical text. *arXiv preprint arXiv:1811.05468*, 2018.
- Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *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)*, pages 4582–4597, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.353. URL <https://aclanthology.org/2021.acl-long.353>.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing, 2021a. URL <https://arxiv.org/abs/2107.13586>.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67, 2020. URL <http://jmlr.org/papers/v21/20-074.html>.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020b. URL <https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf>.
- Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Multimodal few-shot learning with frozen language models. *Advances in Neural Information Processing Systems*, 34:200–212, 2021.
- Kevin Lu, Aditya Grover, Pieter Abbeel, and Igor Mordatch. Pretrained transformers as universal computation engines. *arXiv preprint arXiv:2103.05247*, 2021.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, 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, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=9Vrb9D0WI4>.
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks, 2021b. URL <https://arxiv.org/abs/2110.07602>.
- Milad Moradi, Kathrin Blagec, Florian Haberl, and Matthias Samwald. Gpt-3 models are poor few-shot learners in the biomedical domain, 2021. URL <https://arxiv.org/abs/2109.02555>.

- Sonish Sivarajkumar and Yanshan Wang. Healthprompt: A zero-shot learning paradigm for clinical natural language processing, 2022. URL <https://arxiv.org/abs/2203.05061>.
- Shengding Hu, Ning Ding, Huadong Wang, Zhiyuan Liu, Jingang Wang, Juanzi Li, Wei Wu, and Maosong Sun. Knowledgeable prompt-tuning: Incorporating knowledge into prompt verbalizer for text classification, 2021. URL <https://arxiv.org/abs/2108.02035>.
- Ning Ding, Shengding Hu, Weilin Zhao, Yulin Chen, Zhiyuan Liu, Hai-Tao Zheng, and Maosong Sun. Openprompt: An open-source framework for prompt-learning, 2021. URL <https://arxiv.org/abs/2111.01998>.
- Timo Schick and Hinrich Schütze. Exploiting cloze-questions for few-shot text classification and natural language inference. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 255–269, Online, April 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.eacl-main.20. URL <https://aclanthology.org/2021.eacl-main.20>.
- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt understands, too, 2021c. URL <https://arxiv.org/abs/2103.10385>.
- Karen Hambardzumyan, Hrant Khachatryan, and Jonathan May. WARP: Word-level Adversarial ReProgramming. In *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)*, pages 4921–4933, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.381. URL <https://aclanthology.org/2021.acl-long.381>.
- 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. Mimic-iii, a freely accessible critical care database. *Scientific Data*, 3, 5 2016. ISSN 20524463. doi: 10.1038/sdata.2016.35.
- Chantal Pellegrini, Anees Kazi, and Nassir Navab. Unsupervised pre-training on patient population graphs for patient-level predictions, 2022. URL <https://arxiv.org/abs/2203.12616>.
- Shirly Wang, Matthew B. A. McDermott, Geeticka Chauhan, Marzyeh Ghassemi, Michael C. Hughes, and Tristan Naumann. Mimic-extract: A data extraction, preprocessing, and representation pipeline for mimic-iii. In *Proceedings of the ACM Conference on Health, Inference, and Learning, CHIL '20*, page 222–235, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450370462. doi: 10.1145/3368555.3384469. URL <https://doi.org/10.1145/3368555.3384469>.
- Willie Boag, Dustin Doss, Tristan Naumann, and Peter Szolovits. What’s in a note? unpacking predictive value in clinical note representations. *AMIA Summits on Translational Science Proceedings*, 2018:26, 2018.
- Zheng Yuan, Chuanqi Tan, and Songfang Huang. Code synonyms do matter: Multiple synonyms matching network for automatic icd coding, 2022. URL <https://arxiv.org/abs/2203.01515>.
- Niall Taylor, Lei Sha, Dan W Joyce, Thomas Lukasiewicz, Alejo Nevado-Holgado, and Andrey Kormilitzin. Rationale production to support clinical decision-making. *arXiv preprint arXiv:2111.07611*, 2021.
- Pranav Rajpurkar, Emma Chen, Oishi Banerjee, and Eric J Topol. Ai in health and medicine. *Nature Medicine*, pages 1–8, 2022.
- Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 4596–4604. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/shazeer18a.html>.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing, 2020.

Appendix

A Training details

We implement our experiments using a combination of the OpenPrompt framework Ding et al. [2021] and the Pytorch packages. For prompt learning, we use Adafactor Shazeer and Stern [2018] optimizer for soft and mixed templates, and AdamW Loshchilov and Hutter [2019] optimizer for language models and soft verbalizers. For traditional fine-tuning, we use AdamW optimizer for MLP heads and language models. We train the model on a Nvidia RTX 1080 Ti GPU, with a batch size of 4 due to the memory limitation. To overcome this, we use gradient accumulation for 10 steps during training. Further details of training and hyperparameters can be in the complimentary code repository.

Table A.1 shows the derived optimal hyperparameters for each training paradigm based on the hyperparameter random search. The search consisted of 100 training runs using randomly generated hyperparameters from the search space shown in Table 2. Due to relatively limited computational resource, this was only performed for the ICD-9 Triage task and a sub-sample of the training data was used, similar to that of our few-shot experiments with 128 samples per class.

Table A.1: Optimized hyperparameters for each training paradigm

hp	Traditional fine-tuning	Prompt learning
learning rate	0.0048	0.0121
batch size	4	4
gradient accumulation steps	4	3
dropout	0.382	0.1536
optimizer	adamw	adafactor
verbalizer learning rate	n/a	0.007

B Dataset details

Mortality and Length of Stay For all clinical tasks a combination of available clinical notes pertaining to the outcome of interest were used, including admission and discharge summaries. Each task dataset was created separately and a 70-10-20 split of training-validation-test sets was used. We followed the data engineering steps outlined in the clinical outcomes paper van Aken et al. [2021].

ICD-9 50 and ICD-9 Triage The ICD-9 50 task was simply all clinical notes data corresponding to the top 50 most frequently occurring ICD-9 diagnosis codes. The production of the ICD-9 Triage task was derived from taking the top 20 ICD-9 diagnosis codes. From this subsample, a clinician derived suitable groups representing the destination team on discharge from ICU: Cardiology, Obstetrics, Respiratory Medicine, Neurology, Gastroenterology, Acute or Internal Medicine, and Oncology.

See Fig.B.1 showing class distributions for each of the clinical tasks presented in this paper.

C Prompt examples

Examples of different prompt methods are shown. For each task we show one manual prompt template and one mixed template. The <[soft]> token represents the trainable continuous vector or embedding of the mixed template that has been initialised from the PLMs vocabulary. Thus <[soft]>:"This" indicates a soft embedding initialised from the PLMs representation of the token "This".

ICD-9 diagnosis code triage

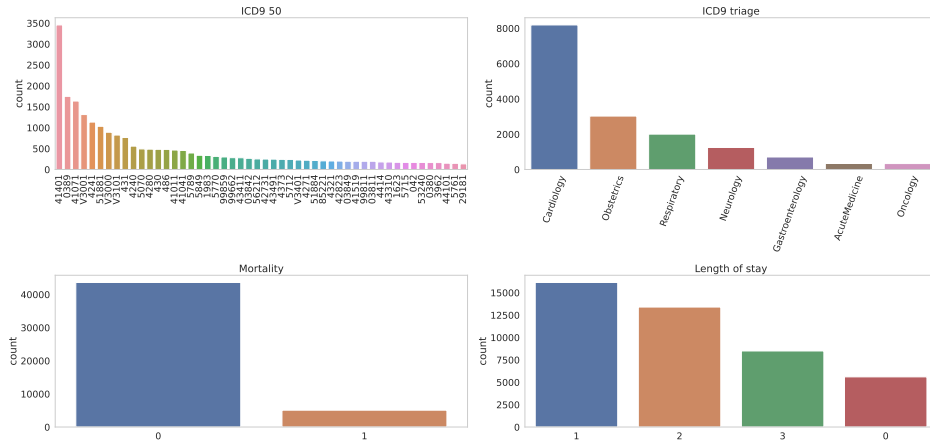


Figure B.1: Distribution of classes for each clinical task

- <clinical note> Best department is <[MASK]>.
- <clinical note> <[soft]>: "This" patient should <[soft]>:"go to this medical team based on symptoms of their illness" <[MASK]>.

Mortality prediction

- <clinical note> Patient is on the path to <[MASK]>.
- <clinical note> <[soft]>: "This" patient <[soft]>:"on path to" <[MASK]>.

ICD-9 diagnosis code classification - top 50

- <clinical note> Patient has diagnosis <[MASK]>
- <clinical note> <[soft]>: "This" patient <[soft]>:"has diagnosis" <[MASK]>.

Length of stay prediction

- <clinical note> The patient will be at hospital with a <[MASK]> length.
- <clinical note> <[soft]>: "This" patient <[soft]>:"will be in hospital for a " <[MASK]> length.

D Prompt learning versus Traditional fine-tuning with PubMed BERT

The PLM used for all presented results in the main body of the paper was the Bio-ClinicalBERT Alsentzer et al. [2019], which we have observed was trained using Mimic-III notes. Whilst this was arguably advantageous for both traditional fine-tuning and prompt learning, it may have overly favoured prompt learning due to the reformulation of the classification task as a Masked Language Modelling (MLM) objective. Therefore we present results of another biomedical BERT model from Microsoft, the PubMedBERT, which was pre-trained from scratch using abstracts from PubMed Gu et al. [2020] in Table D.1. It can be seen that prompt learning still outperforms traditional fine-tuning by a large margin on the ICD-9 Triage task, in line with our other results.

Table D.1: Balanced accuracy results for prompt learning and traditional fine-tuning using Microsoft’s PubMedBert

Sample size	Balanced Accuracy	
	Traditional fine-tuning	Prompt learning
16	0.1554	0.2249
32	0.1521	0.3749
64	0.4048	0.4621
128	0.5621	0.7814