

Personalized Transformers for Everyone

Anonymous ACL submission

Abstract

Personalized Intelligence (PI) is the problem of providing customized AI experiences tailored to each individual user. A related problem is the *compartmentalization* of intelligence that maintains a partition between the personalized and the general models. Existing personalization approaches involve fine-tuning pre-trained models to create new customized models. However, these require a significant amount of computation to train, which scales with model size and the number of users, inhibiting PI to be realized widely. A compartmentalized approach enables a small model to be specialized for each individual user, which needs to be used together with a larger model to provide personalization. By separating personalized and general models, we enable higher accuracy, scalability, and stronger privacy guarantees. In this paper, we aim to design a compartmentalized personalization approach that can scale to millions of users and beyond. We investigate the landscape of model fine-tuning techniques and construct new design adaptations based on the requirements of PI. We then introduce Personalized Head (PH), a new model training/inference framework designed for scalable PI. We explore the design space of these techniques and evaluate their efficacy under various production-level constraints. Specifically, we break down the trade-off between accuracy, scalability and production deployment limitations. We present several production-ready personalization approaches suited for various production use case scenarios.

1 Introduction

Personalized Intelligence (PI) is the problem of providing unique and customized AI experience tailored for each individual user. Today’s AI in production is often served with an unified model shared among all users and the experience remains largely homogeneous across users. However, certain problems are highly personal and the task scope can vary from user to user, which limits the effectiveness of this single model approach. For instance, in productivity software, users often use personalized category labels and tags to organize their to-do items. For the task of classifying a new unseen item to a category, a shared model is limited

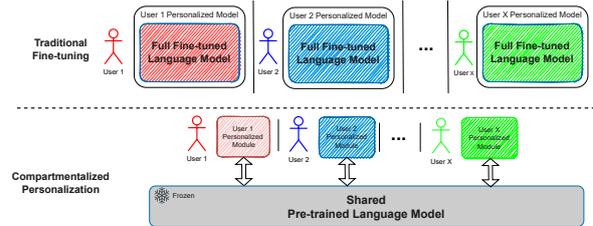


Figure 1: Full-model fine-tuning approach vs. compartmentalized personalization training approach.

as each user can have unique category labels and can interpret them differently. More PI use case examples are described in Section 2.

In order to capture the user-specific knowledge, PI tasks require personalized weights and module to learn from user-specific data. This presents several key requirements for solving PIs. First, the training of the personalized models needs to be lightweight and fast, as they often happen in an online settings. Second, the personalized models need to be parameter efficient. The amount of personalized models scales with the number of users, which can be millions or even billions, so the personalized model size needs to be small so they can be stored and served at large scale. Third, the personalized weights need to be compartmentalized from the original model, in order to protect the privacy of user data and maintain the integrity of each user’s model quality.

The goal of this paper is to design an approach for achieving personalized intelligence at the production scale of millions of users and beyond. We examine a landscape of personalization techniques and aim to answer the following research questions:

1. How scalable can we achieve with the current landscape of compartmentalized personalization training techniques?
2. What are the design knobs of each technique and what kind of trade-off do they represent in the design space?
3. How do the techniques compare under various

075	production level constraints?	
076	We investigate the current landscape of	
077	parameter-efficient fine-tuning techniques, includ-	
078	ing Adapters (Houlsby et al., 2019) and Prefix Tun-	
079	ing (Li and Liang, 2021). We identify key limita-	
080	tions of these approaches to be the compartmen-	
081	talization of the personalized weights as they are	
082	intertwined with the layers and weights of the orig-	
083	inal model. In addition, these approaches are de-	
084	signed for task-level fine-tuning and it remains an	
085	open question whether they are scalable enough	
086	for creating user-level personalized models. To ad-	
087	dress these, we make two contributions. We first	
088	construct new design adaptations specifically for	
089	personalization, namely Personalized Adapter (PA)	
090	and Personalized Prefix (PP). Secondly, we propose	
091	a new model training and inference framework, Per-	
092	sonalized Head (PH), specifically for personaliza-	
093	tion at scale. In PH, the personalized weights are	
094	completely separated from the original base model.	
095	We explore the design space of all three techni-	
096	ques (PA, PP and PH) and compare their accu-	
097	racy and scalability. We found that Personal-	
098	ized Prefix (PP) outperforms Personalized Adapter	
099	(PA) and Head (PH) when no deployment con-	
100	straints are considered and all layers of the original	
101	model are augmented with new weights. On the	
102	other hand, under production level deployment con-	
103	straints, there are several viable options across the	
104	three techniques, each with unique advantages de-	
105	pending on the use case.	
106	2 Personalized Intelligence	
107	We describe in details the concept of Personalized	
108	Intelligence and provide three concrete example	
109	use cases that can be found in production today.	
110	2.1 Definition	
111	Personalized Intelligence (PI) is the problem of	
112	creating unique and customized experience of an	
113	AI capability for each individual user. To provide	
114	experience tailored to each user, PI requires user-	
115	specific weights that are trained on user-specific	
116	data. The user-specific weights also serve as com-	
117	partmentalization for the learnt knowledge of each	
118	user to not influence and affect other users' mod-	
119	els and experience. PI requires every user-specific	
120	models to be query-able and incoming requests are	
121	routed accordingly based on their source. This cre-	
122	ates significant scalability challenges for traditional	
123	fine-tuning approach. Specifically, it is not feasible	
124	to have full size model replicas for each user due to	
125	the ever increasing size of state-of-the-art models.	
	A similar problem that has been studied recently	126
	is task-specific fine-tuning (Houlsby et al., 2019;	127
	He et al., 2021; Li and Liang, 2021; Hu et al., 2021).	128
	While both problems involve adapting pre-trained	129
	models to new data, task-specific fine-tuning gener-	130
	ates one model per task and only that model needs	131
	to be hosted and served to the users. In contrast,	132
	in personalized intelligence, unique weights and	133
	learning is required for each user and total num-	134
	ber of query-able models scales with number of	135
	users. Furthermore, the task scope is fully defined	136
	at the model creation time for task-specific fine-	137
	tuning, while for personalized intelligence, part of	138
	the problem scope is unique for each user.	139
	In summary, Personalized Intelligence tasks gener-	140
	ally have one or more of the following character-	141
	istics:	142
	• Part of the problem scope is user-specific and	143
	not fully defined at model creation and pre-	144
	training time. The scope can also evolve over-	145
	time, pre- and post-production deployment.	146
	• Unique weights are required on a user-by-user	147
	basis to capture the personalized knowledge	148
	of each user.	149
	• The knowledge learnt from the user-specific	150
	data needs to be compartmentalized as to	151
	avoid potentially contaminating the behavior	152
	of other users' models and experience.	153
	2.2 Example use-cases of PI	154
	We describe three example production use cases of	155
	Personalized Intelligence, as illustrated in Figure 2.	156
	Category Prediction in Productivity Tools - In	157
	productivity apps such as Todoist and OmniFocus,	158
	users are encouraged to create and organize their	159
	tasks into custom categories or tags (e.g., work,	160
	family, hobbies, etc). The task of classifying an	161
	item to a category/tag by learning from past user	162
	behavior is an example of Personalized Intelligence,	163
	as the category labels can be drastically different	164
	between users and are not fully defined before-	165
	hand (Figure 2a). A successful category prediction	166
	model needs to train on each user's data to capture	167
	user's personal preference. Furthermore, this per-	168
	sonalized knowledge need to be compartmentalized	169
	to avoid affecting other users.	170
	Intent Classification in Dialogue Systems - For	171
	intent classification in dialogue systems, each state	172
	has its own set of candidate intents correspond-	173
	ing to the set of next possible states, as shown in	174
	Figure 2b. Based on the user, the map of dialogue	175
	states can also be different and as a result, the scope	176

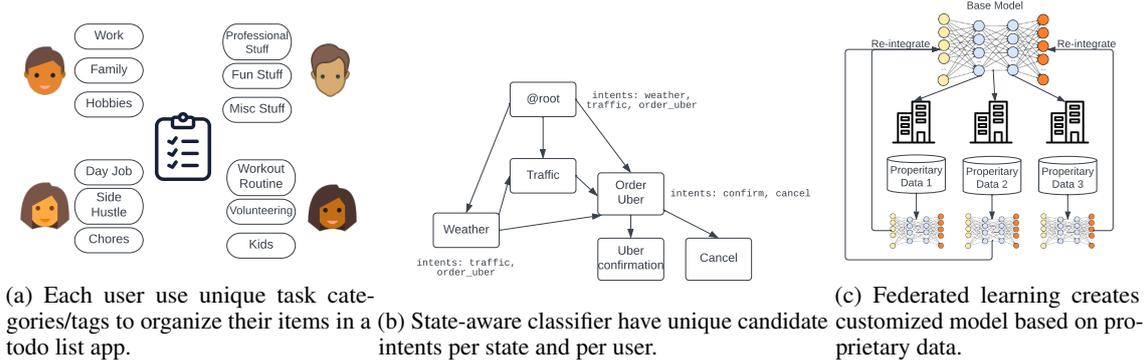


Figure 2: Example use cases of Personalized Intelligence (PI) in production today.

of the intent classification for a given user is unique. Furthermore, the set of candidate intents generally evolve and grow as new topics and capabilities are introduced. A fully built-out dialogue system can feature hundreds or more states (Larson et al., 2019), which combined with the conversation paths unique to each user, can lead to a large number of unique classification problems within one dialogue system. A personalized intent classification model is customized to the context of each state and each user and trained to focus on the state’s candidate intents and specific type of utterances encountered from a given user.

Federated Learning - The goal of federated learning is to train a central model based on data that is distributed over a large number of clients (Konečný et al., 2016) (Figure 2c). Clients independently compute model updates based on local data, and client-side updates are eventually aggregated to update the central model. The original motivation for this technique was driven by the needs of distributed systems, where learning takes place based on data at the edge (e.g. in mobile phones) which may not always be connected to the central model. Follow-on work (Gao et al., 2019) emphasized the privacy-preserving aspects of federated learning, as the training data itself does not need to leave its original location at the edge. Personalized Intelligence techniques can be used to support federated learning use-cases, as the information contained in each personalized model can be transferred to the general model, if needed, under user control. Instead of assuming that all knowledge learned at the edge is going to end up in the general model, PI approach emphasizes privacy-preservation and compartmentalization to provide user-specific inference.

3 Training for Personalization

We define the process of training personalized models. Let M be a pre-trained language model (LM) and Θ_M its trainable parameters. Consider the scenario of fine-tuning for an individual user to create a personalized model. We define $U = \{(D_1, L_1), (D_2, L_2), \dots, (D_N, L_N)\}$ as the collection of personalized tasks where $i \in [1, \dots, N]$ for N users, D_i is the unique data for user i and L_i is the loss function.

3.1 Traditional Fine-Tuning

When applying model M to a downstream task T with labelled training data D_T , an output layer K with parameters Θ_K is usually appended to M . Then M and K are trained jointly:

$$\Theta'_M, \Theta'_K \leftarrow \arg \min_{\Theta_M, \Theta_K} L_T(D_T; \Theta_M, \Theta_K) \quad (1)$$

This generates Θ'_M and Θ'_K . Θ'_M is the fine-tuned parameters of M , which are the same size as Θ_M but with distinct values. Let $\Omega(\Theta)$ be the computation complexity required to train a set of parameters Θ . We then define the training complexity of the above fine-tuning operation as:

$$\Omega(\Theta'_M) + \Omega(\Theta'_K) \quad (2)$$

The problem of fine-tuning the model to personalize for user i is defined as

$$\Theta'_{M,i}, \Theta'_{K,i} \leftarrow \arg \min_{\Theta_M, \Theta_K} L_i(D_i; \Theta_M, \Theta_K) \quad (3)$$

The aggregated training complexity scales linearly with the number of users, N :

$$\sum_1^N \Omega(\Theta'_{M,i}) + \sum_1^N \Omega(\Theta'_{K,i}) \quad (4)$$

The collection of model parameters to be stored scales linearly with N as well:

$$\sum_1^N |\Theta_{M,i}| + \sum_1^N |\Theta_{K,i}| \quad (5)$$

3.2 Compartmentalized Personalized Training

When training compartmentalized and personalized model for user i , a new set of weights W_i are introduced. During training, the base LM parameters are frozen and only W_i and K_i are updated:

$$\Theta_{W}, \Theta_{K} \leftarrow \arg \min_{\Theta_{W}, \Theta_{K}} L_T(D_T; \Theta_M, \Theta_W, \Theta_K) \quad (6)$$

Note that, compared to the traditional fine-tuning defined in Equation 1, no Θ_M is generated. During inference, Θ_{W}, Θ_{K} is combined with the original Θ_M to generate prediction. The aggregated training complexity for training the personalized model is:

$$\sum_1^N \Omega(\Theta_{W,i}) + \sum_1^N \Omega(\Theta_{K,i}) \quad (7)$$

and the total parameters is:

$$\Theta_M + \sum_1^N |\Theta_{W,i}| + \sum_1^N |\Theta_{K,i}| \quad (8)$$

The main goal for a scalable personalization approach is minimizing size of W_i , which will reduce the model size and training cost and increase the maximum users supported N given a certain amount of compute resources, and optimizing for prediction accuracy on the collection of personalized tasks U .

4 Personalization Model Architectures

We describe the landscape of three personalization techniques studied in this work, Personalized Adapter (PA), Personalized Prefix (PP) and Personalized Head (PH), as illustrated in Figure 3.

4.1 Personalized Adapters

Personalized Adapter (PA) is constructed based on the Adapter approach, which involves inserting small trainable feedforward layers into every layer of the base transformer model (Houlsby et al., 2019; Pfeiffer et al., 2021; Hu et al., 2021). During training, the inserted adapter layers are updated while the base transformer is frozen. In the context of Personalized Intelligence, the out-of-the-box adapter

approach are not compartmentalized because the new layers are interleaved with the layers of the original model. This limits adapter’s applicability to training for personalization. Specifically, it is technically challenging to have many adapters sharing the same base model for inference because the inference execution flow switches back and forth between the original model and the adapter.

To address this limitation, we formulate Personalized Adapter (Figure 3a), where only a subset of transformer layers are augmented with adapters. This way, in production, the augmented layers can be replicated for each user to create compartmentalization and the untouched layers in the original models can be shared across users for inference. In order to understand the impact of selectively applying Adapters, we construct a range of different PA configurations and evaluate their trade-offs in Section 6.2.

4.2 Personalized Prefix

The Prefix tuning approach prepends trainable weight vectors to the keys and values weight matrices of the multi-head attention block in each transformer layer (Li and Liang, 2021). The new prefix vectors are used in the attention calculation for attention heads of every transformer layer. Similar to Adapters, Prefix tuning has the same compartmentalization limitation and we construct **Personalized Prefix**, where selected layers are augmented with the prefix vectors, as shown in Figure 3b. We experiment with a range of different PP configurations in Section 6.3.

4.3 Personalized Head

Inspired by the above approaches and considering their limitation, we propose a new framework, Personalized Head (PH). We propose to append a single layer transformer module after the base transformer. During training, only the PH is updated and the base model weights are frozen. During inference, the base model processed the input first and the output embeddings are sent to the PH to apply the personalized knowledge and generate the output. Our intuition is we can aggregate and focus the personalized knowledge into the PH module and avoid augmenting any of the base model layers. Figure 3c shows an overview of the PH architecture.

PH follows the Transformer architecture defined in the original transformer paper (Vaswani et al., 2017). Each PH has a multi-head self-attention layer and two fully connected layers, followed by

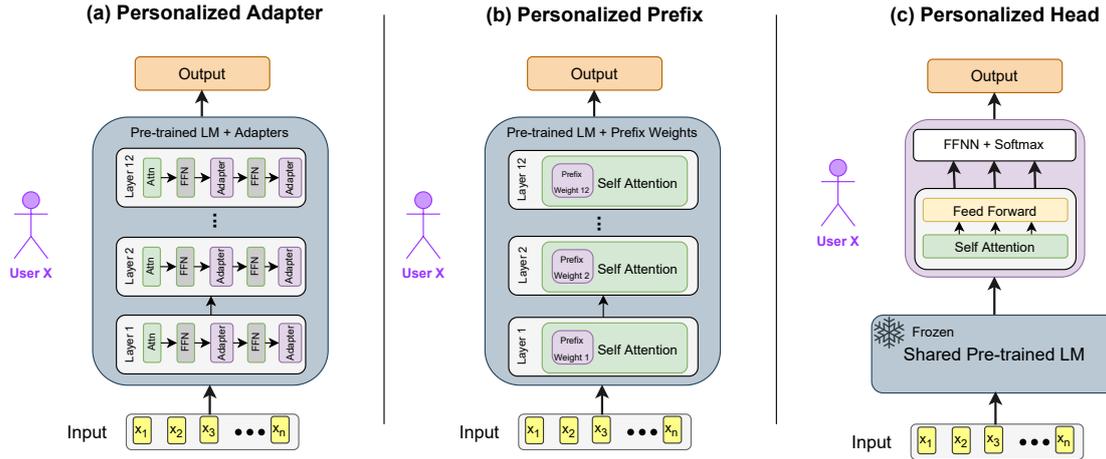


Figure 3: Three personalization approaches. The purple colored block in each approach represents the personalized module for each user. (a) Personalized Adapter insert adapter layers inside selective transformer layers; (b) Personalized Prefix augment the K,V weight matrices of the self-attention block with new personalized weights in selective transformer layers; (c) Personalized head attach a new personalized transformer layer after the original transformer.

Dataset	Description	# Classes	# Train	# Test
SNIPS	Smart assistants questions	7	13,034	1,442
Clinic150	Production VA tasks	150	15,100	1,500

Table 1: Datasets

layer normalization (Ba et al., 2016). Dropout (Srivastava et al., 2014) is applied to the output of the fully connected layers.

To help us explore the PH design space and understand its key design factors, we parameterize the size of the hidden dimension of the feed-forward network in the encoder and the number of attention heads in the attention layer. We investigate the impact of these design decisions in detail in Section 6.4.

5 Experiments

To investigate the effectiveness of these techniques, we apply them to personalize the pre-trained BERT LM to two new classification tasks as the personalized problems. We keep the BERT LM frozen during personalized training and apply each technique separately to train on the new data. In this section, we describe in details the experiments and dataset setup.

5.1 Universal Binary Classification Task

We focus on personalized classification as the PI task to evaluate the suite of personalization techniques. We aim to design a personalization framework that is generalizable to arbitrary classification tasks without requiring modification to the model architecture. To that end, we formulate the multi-

class classification problem as a series of binary classification tasks. We concatenate the class label and the text as input and the output layer generates a binary True/False prediction with a confidence score. The class with the most confident True prediction is selected as the classification prediction. We apply this binary classification across all personalization approaches.

5.2 Datasets

We use the SNIPS dataset (Coucke et al., 2018) and Clinic-150 dataset (Larson et al., 2019). We select SNIPS because its intents cover many common classification topics and it is a representative dataset widely studied in the literature. We select Clinic-150 for its focus on the complexity of production use cases. It has 150 intents and features intents and sentences inspired by real virtual assistants in production. An overview of the datasets is shown in Table 1. Specifically for Clinic, to make training and evaluation time more manageable, we randomly sample 10 examples (out of 30) per class to construct the test set and include the True examples and randomly sample 2 False examples for every True example to construct the training set. The same train and test set construction is used across all techniques and experiments.

5.3 Implementation

We use the AdapterHub framework (Pfeiffer et al., 2020) for the experiment on Adapters and Prefix Tuning. We implement the PHs using the Flair NLP framework (Akbi et al., 2019) with an underlying

pytorch runtime (Paszke et al., 2019). We use the uncased BERT encoder as the base LM (Devlin et al., 2018) for all three approaches. We train for 50 epochs (unless noted otherwise) and report F1 score and accuracy on the test set. In terms of hyper-parameters for PH, we use a batch size of 16 and a learning rate of 0.02, following the standard in (Halder et al., 2020).

6 Results

We evaluate the landscape of personalization techniques, including Personalized Adapters (PA), Personalized Prefix (PP) and Personalized Head (PH). We aim to understand 1) how scalable can we push these techniques to be, 2) what are the design knobs that have the most impact on performance and 3) how do these techniques compare under production-level constraints.

6.1 Naive Personalization Approaches

Table 2 shows the results for several naive personalization approaches. **Full Model Fine-tuning** - Fine-tuning both LM and linear layer achieves strong performance but it requires, for each user, the training of 109 million parameters which generates a 417MB model. The large computation cycles and storage capacity becomes untenable for production use cases that require scaling to many (millions) users.

Personalization via the linear layer only - We experiment with only training the linear classification layer. It achieves F1-score of 71.12 for SNIPS and 52.43 for Clinc, significantly underperform the full fine-tuning approach but only requiring per user 1.5K parameter.

Personalization via Zeroshot - We also investigate the efficacy of the zeroshot approach. We use the TARS zeroshot classifier from (Halder et al., 2020), which uses BERT as the underlying language model and is pre-trained on a suite of datasets including AGNews and DBPedia (Zhang et al., 2015). We test the zeroshot performance of both the trained TARS classifier and BERT out-of-the-box. BERT(OOB) achieves F1-score of 2.99 and 0.78 on SNIPS and Clinc, respectively, and the TARS classifier achieves an F1-score of 35.27 and 23.98 on SNIPS and Clinc, respectively, which are significantly lower than the reported state-of-the-art results. This shows that zeroshot approach requires significant improvement to reach the level of performance needed for production usage.

6.2 Personalized Adapter

Table 3 shows the performance of Personalized Adapters on SNIPS and Clinc. As described in Section 4, we construct Personalized Adapters by selectively insert adapters to a subset of layers in the transformer. We observe that the adapter’s ability to learn personalized knowledge depend heavily on the number and location of the inserted adapters. When only adding adapters to a selected set of layers, we see that inserting adapters to early layers achieve better performance compared to later layers, given the same number of adapters. Specifically, inserting adapters for the first half of the BERT model (1-6 layers) outperform inserting adapters to the second half (7-12). Moreover, if the situation only allows for a single transformer layer to be augmented with new weights, adding adapter to first layer outperforms the last layer.

Limitation and Trade-off for Production In Personalized Adapters, personalized layers and original layers are interleaved and the model needs to switch between them many times during an inference. This makes it challenging to have many PAs live in production to share the same base model. In addition, the tightly integrated models require the personalized weights and the original model to be co-located in the same environment, limiting the possible deployment configurations, e.g. deploying personalized weights on user’s devices and sharing the base model in the cloud. When only augmenting selective layers, those layers can be replicated across users to create a more compartmentalized approach, at the cost of accuracy. Under this approach, the most scalable design is applying PAs on the first layer or first several layers.

6.3 Personalized Prefix

Table 5 shows the result of Personalized Prefix. Similar to PA, we construct variations of Personalized Prefix by applying prefixes to different layer combinations. Overall, PP achieves similar performance as PA. However, PP requires significantly more parameters per user. In addition, at lower parameter count, PA outperforms PP with less parameters per user. We also observe that applying prefixes at early layers of the transformers yields higher accuracy than early layers, similar to PA.

Limitation and Trade-off for Production Personalized Prefix approach requires 9.87M parameters (37MB) per user when all layers are augmented. This is a substantial amount of parameters when considering potentially scaling to millions of users. However, this technique has the unique property

Naive Personalization Approach		# Params / User	Size / User	SNIPS F1 / Acc.	Clinic F1 / Acc.
ZeroShot	BERT	NA	NA	2.99 / 1.60	0.78 / 0.47
	TARS	NA	NA	35.27 / 26.70	23.98 / 23.67
Fine-tuning LM + Linear Layer	BERT	109M	417MB	98.61 / 98.61	95.74 / 95.07
	TARS	109M	417MB	98.13 / 98.06	95.22 / 94.27
Fine-tuning Linear Layer Only	BERT	1.5K	7KB	68.70 / 58.67	52.43 / 50.27
	TARS	1.5K	7KB	71.12 / 63.11	33.27 / 33.20

Table 2: F1 score and accuracy of zeroshot, fine-tuning the LM + linear head and fine-tuning the linear head only, evaluated on SNIPS and Clinc-150 datasets. We also show the number of parameters that are required to be fine-tuned for each approach, as well as the size of the personalized model to be managed for each user as a representation of the scalability of each approach.

Personalized Adapter	# Params / User	Size / User	SNIPS F1 / ACC.	Clinic F1 / ACC.
Full (1-12)	894K	3.6MB	99.02 / 99.02	89.13 / 89.13
1st half (1-6)	447K	1.8MB	98.19 / 98.19	89.26 / 89.26
2nd half (7-12)	447K	1.8MB	98.54 / 98.54	89.33 / 89.33
last layer (12)	74.5K	0.3MB	97.22 / 97.22	64.80 / 64.80
first layer (1)	74.5K	0.3MB	98.06 / 98.05	82.00 / 82.00

Table 3: Performance and size of personalized adapters.

Personalized Prefix	# Params / User	Size / User	SNIPS F1 / ACC.	Clinic F1 / ACC.
Full (1-12)	9.87M	37MB	98.95 / 98.95	91.80 / 91.80
1st half (1-6)	4.94M	19MB	98.74 / 98.75	89.40 / 89.40
2nd half (7-12)	4.94M	19MB	98.13 / 98.12	88.70 / 88.60
last layer (12)	823K	3.1MB	95.61 / 95.63	61.04 / 61.06
first layer (1)	823K	3.1MB	97.45 / 97.43	80.55 / 80.53

Table 4: Performance and size of personalized prefix.

where the prefix weights can be stored on the user side and then pass to the base model along with the input thus achieving complete compartmentalization of the personalized weights (Li and Liang, 2021). In a production settings, each user’s prefix weights can be stored in a separate database and during inference, first fetch the prefix weights and then send it to the shared base model. Note that this approach introduces an additional inference latency overhead for the step of fetching the prefix weights for the querying user.

6.4 Personalized Head

For PH, we experiment with a wide range of configurations by varying the hidden dimension of the feed-forward layer and the number of attention heads in the PH. For brevity, we include results for configurations with hidden dimensions from 128 to 2048 and # attention heads of 2, 4, and 8. Results for additional configurations are included in the Appendix. We observe that PHs, across all configurations, significantly outperform fine-tuning only the linear layer on both datasets. When comparing to the baseline of fine-tuning the entire model stack of the base language model and the linear output layer, PH achieves similar results for the SNIPS dataset, while requiring orders of magnitude less

training cost. PH performs similarly as PA and PP on SNIPS, while underperforms the previous two approaches on Clinc, which is a more difficult dataset. As for scalability, PH requires less parameters than all of PP configurations except for single layer prefix and requires more parameters than PA. We have included a deep dive into how to effectively design and train a PH in the Appendix.

Limitation and Trade-off for Production PHs are completely detached from the base model and inference uses both models in combination to generate its prediction. By having a simple and clear boundary between the base model and the personalized, PH makes it trivial to share the base model among many users, using existing out-of-the-box model implementation. It also enables more flexible deployment options, such as cloud and edge collaboration which has been shown to improve latency and energy efficiency (Kang et al., 2017), where PH is deployed on user’s own mobile/edge device while keeping the large-scale base model in the cloud.

7 Discussion

We compare across the three personalization techniques and summarize their advantages and disadvantages in different production scenarios.

Personalized Adapters (PA) is the most parameter-efficient approach, achieving the highest accuracy per parameter. It also slightly outperforms the other two techniques and the full fine-tuning baseline for the simpler SNIPS dataset. The disadvantage of this approach is the difficulty to share one base model across many users when multiple layers are augmented with PA. Therefore, an option here to use the last layer PA variant and replicate it across users, when dealing with low complexity personalization tasks and compute resources is limited.

Personalized Prefix (PP) achieves the highest accuracy among all three approaches but it requires

Model			# Params / User	Size / User	SNIPS F1 / ACC.	Clinc F1 / ACC.
Personalization Head (PH) w/ frozen LM	Hidden Dim	# Attn. Heads				
	2048	8	5.52M	21MB	96.52 / 96.12	76.36 / 75.93
		4			97.18 / 96.74	76.46 / 76.47
		2			97.33 / 97.16	76.36 / 75.93
	1024	8	3.94M	15MB	97.46 / 97.09	75.07 / 74.40
		4			96.76 / 96.39	75.82 / 75.73
		2			96.77 / 96.67	76.61 / 76.60
	512	8	3.15M	12MB	97.05 / 97.02	75.95 / 76.33
		4			96.26 / 95.49	70.79 / 71.20
		2			96.94 / 96.74	75.29 / 75.27
	256	8	2.76M	11MB	95.90 / 95.70	66.99 / 67.53
		4			95.64 / 95.15	68.64 / 67.87
		2			97.06 / 96.32	67.68 / 67.13
	128	8	2.57M	9.8MB	95.32 / 95.28	63.43 / 64.53
		4			96.32 / 96.32	62.70 / 63.67
		2			96.36 / 96.36	63.82 / 64.60

Table 5: F1 score and accuracy of personalization head.

the most parameters per user. The accuracy advantage is most prominent on harder tasks, such as Clinc-150. Therefore, Personalized Prefix is the preferred solution for more sophisticated personalization tasks with less demanding scalability requirement (e.g., smaller user count and/or more compute resources available).

Personalized Head (PH) represents a unique position among the three personalization techniques. It achieves better or similar performance than PA and PP on the SNIPS dataset and under-performs them on the Clinc dataset. Because of its complete compartmentalized design, PH is the easiest to deploy for production. It also enables on-device deployment and training of the personalized module, which protects privacy of user’s data, while the other techniques require the personalized weights to co-locate with the base model.

8 Related Work

One of the most common solutions to many natural language understanding problems today is leveraging large-scale pre-trained transformer-based language models (Devlin et al., 2018; Liu et al., 2019) which are typically trained on language understanding objectives such as Masked Language Modeling and Next Sentence Prediction. These language models are then fine-tuned for a specific task. This transfer of learning to the task of interest is achieved by tuning all model weights on that new task. The performance of these LMs has shown to scale with model size (Kaplan et al., 2020), resulting in massive models consisting of billions of parameters (Brown et al., 2020; Raffel et al., 2020; Sanh et al., 2021). When applied to the online setting of personalization training, the applicability of these language models is severely constrained as it results in a dedicated model for each user.

This section explores existing works improving the applicability of transformer models at scale.

Zero-shot learning approaches aim to provide generalized models for a range of language-based tasks without needing additional training steps required by traditional transfer learning approaches. Recent approaches to this problem frame this as a text-to-text generation task (Brown et al., 2020; Raffel et al., 2020; Sanh et al., 2021) with focuses on prompt design (Perez et al., 2021; Khashabi et al., 2020). While shown to be effective in tasks such as QA and summarization, zero-shot performance is still lacking when it comes to text classification (Halder et al., 2020; Wenpeng Yin and Roth, 2019). This is further shown in our zero-shot experimental results. Halder et al. (2020) explores the shortcomings of the existing transfer learning mechanisms for text classification, proposing the formalization of text classification as a general binary classification problem.

9 Conclusion

In this work, we introduce and define a new research problem, Personalized Intelligence (PI). PI is the problem of creating customized AI experience that is tailored to each individual user. PI brings a new host of challenges for existing fine-tuning techniques, increasing the pressure on model scalability. We examine a landscape of personalization techniques and investigate their performance and trade-off and present limitations and considerations for using each technique under production level constraints. Finally, we compare across all techniques to provide concrete recommendations on the best approach for Personalized Intelligence given various production scenarios.

References

Alan Akbik, Tanja Bergmann, Duncan Blythe, Kashif Rasul, Stefan Schweter, and Roland Vollgraf. 2019. Flair: An easy-to-use framework for state-of-the-art nlp. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 54–59.

Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450*.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared 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 M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#).

Alice Coucke, Alaa Saade, Adrien Ball, Théodore Bluche, Alexandre Caulier, David Leroy, Clément Doumouro, Thibault Gisselbrecht, Francesco Caltagirone, Thibaut Lavril, et al. 2018. Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces. *arXiv preprint arXiv:1805.10190*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Dashan Gao, Yang Liu, Anbu Huang, Ce Ju, Han Yu, and Qiang Yang. 2019. [Privacy-preserving heterogeneous federated transfer learning](#). In *2019 IEEE International Conference on Big Data (Big Data)*, pages 2552–2559.

Kishaloy Halder, Alan Akbik, Josip Krapac, and Roland Vollgraf. 2020. Task-aware representation of sentences for generic text classification. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3202–3213.

Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2021. [Towards a unified view of parameter-efficient transfer learning](#). *CoRR*, abs/2110.04366.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gemundo, Mona Attariyan, and Sylvain Gelly. 2019. [Parameter-efficient transfer learning for nlp](#).

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#). *CoRR*, abs/2106.09685.

Yiping Kang, Johann Hauswald, Cao Gao, Austin Rovinski, Trevor Mudge, Jason Mars, and Lingjia Tang. 2017. Neurosurgeon: Collaborative intelligence between the cloud and mobile edge. *ACM SIGARCH Computer Architecture News*, 45(1):615–629.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. [Scaling laws for neural language models](#). *CoRR*, abs/2001.08361.

Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Clark, and Hannaneh Hajishirzi. 2020. [Unifiedqa: Crossing format boundaries with a single qa system](#).

Jakub Konečný, H. Brendan McMahan, Felix X. Yu, Peter Richtarik, Ananda Theertha Suresh, and Dave Bacon. 2016. [Federated learning: Strategies for improving communication efficiency](#). In *NIPS Workshop on Private Multi-Party Machine Learning*.

Stefan Larson, Anish Mahendran, Joseph J. Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K. Kummerfeld, Kevin Leach, Michael A. Laurenzano, Lingjia Tang, and Jason Mars. 2019. [An evaluation dataset for intent classification and out-of-scope prediction](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1311–1316, Hong Kong, China. Association for Computational Linguistics.

Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#).

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized BERT pretraining approach](#). *CoRR*, abs/1907.11692.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. [Pytorch: An imperative style, high-performance deep learning library](#). In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.

Ethan Perez, Douwe Kiela, and Kyunghyun Cho. 2021. [True few-shot learning with language models](#).

Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. 2021. [Adapterfusion: Non-destructive task composition for transfer learning](#).

- 740 Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya
741 Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun
742 Cho, and Iryna Gurevych. 2020. [Adapterhub: A
743 framework for adapting transformers](#). In *Proceedings
744 of the 2020 Conference on Empirical Methods in Nat-
745 ural Language Processing (EMNLP 2020): Systems
746 Demonstrations*, pages 46–54, Online. Association
747 for Computational Linguistics.
- 748 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine
749 Lee, Sharan Narang, Michael Matena, Yanqi Zhou,
750 Wei Li, and Peter J. Liu. 2020. [Exploring the limits
751 of transfer learning with a unified text-to-text trans-
752 former](#).
- 753 Victor Sanh, Albert Webson, Colin Raffel, Stephen H.
754 Bach, Lintang Sutawika, Zaid Alyafeai, Antoine
755 Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja,
756 Manan Dey, M Saiful Bari, Canwen Xu, Urmish
757 Thakker, Shanya Sharma Sharma, Eliza Szczechla,
758 Taewoon Kim, Gunjan Chhablani, Nihal Nayak, De-
759 bajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang,
760 Han Wang, Matteo Manica, Sheng Shen, Zheng Xin
761 Yong, Harshit Pandey, Rachel Bawden, Thomas
762 Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma,
763 Andrea Santilli, Thibault Fevry, Jason Alan Fries,
764 Ryan Teehan, Stella Biderman, Leo Gao, Tali Bers,
765 Thomas Wolf, and Alexander M. Rush. 2021. [Multi-
766 task prompted training enables zero-shot task gener-
767 alization](#).
- 768 Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky,
769 Ilya Sutskever, and Ruslan Salakhutdinov. 2014.
770 Dropout: a simple way to prevent neural networks
771 from overfitting. *The journal of machine learning
772 research*, 15(1):1929–1958.
- 773 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
774 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
775 Kaiser, and Illia Polosukhin. 2017. Attention is all
776 you need. In *Advances in neural information pro-
777 cessing systems*, pages 5998–6008.
- 778 Jamaal Hay Wenpeng Yin and Dan Roth. 2019. [Bench-
779 marking zero-shot text classification: Datasets, eval-
780 uation and entailment approach](#). In *EMNLP*.
- 781 Jack Schofield ZDNet. 2012. Is your Gmail re-
782 ally worth \$3,600? Backup now! [https://
783 www.zdnet.com/article/is-your-gmail-
784 really-worth-3600-backup-now/](https://www.zdnet.com/article/is-your-gmail-really-worth-3600-backup-now/). [On-
785 line; accessed 16-Jan-2022].
- 786 Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015.
787 Character-level convolutional networks for text clas-
788 sification. In *NIPS*.

789 A Appendix 841

790 A.1 How to effectively train a PH 842

791 We conduct experiments aimed at understanding 843
792 the learning behavior of PHs and gain insights into 844
793 how to design and deploy an effective PH for real- 845
794 world use cases. We aim to answer the following 846
795 questions: 1) How to effectively train PHs in pro- 847
796 duction? 2) How does the PH configuration affect 848
797 its learning behavior? 3) Do larger PHs achieve 849
798 better performance and does there exist a sweet 850
799 spot of PH design that is the most compute and 851
800 data efficient? 852

801 Impact of Data vs. Epoch on Training PHs 853

802 We study the impact of the scale of training data 854
803 and the number of training epochs on the PH per- 855
804 formance. In real-time training in production, there 856
805 is often limited training data and training cycles 857
806 available. Therefore, it is imperative to understand 858
807 how to train a PH in a compute and data efficient 859
808 manner. To this end, we construct a SNIPS sub- 860
809 dataset by random sampling 100 training examples 861
810 per class (1400 samples in total) and keep the full 862
811 SNIPS test set. We train the spectrum of PH de- 863
812 signs for 50 epochs and then 50 more epochs (100 864
813 epochs in total) and record the test set F1-score 865
814 at both points. We then select another 100 train- 866
815 ing samples per class to add to the training set and 867
816 repeat the same experiment. Table 6 shows the av- 868
817 erage F1 scores, as well as improvement gained by 869
818 increasing training data, training epochs, and both. 870

819 We observe that increasing the training data from 871
820 100 per class to 200 per class provides a signifi- 872
821 cantly higher F1 score increase (+19.85 on aver- 873
822 age), compared to training for more epochs (+5.98 874
823 average). This behavior is consistent across the var- 875
824 ious PH configurations. This is intuitive because 876
825 100 training samples/class represents only 5.3% of 877
826 the full SNIPS training set and does not provide ro- 878
827 bust coverage of the problem space. Increasing the 879
828 training data scale should be the priority over more 880
829 training iterations in the early stages of applying a 881
830 PH to a personalized problem. 882

831 **PH Design Analysis** Two main design choices 883
832 for PHs are the hidden dimension of the encoder 884
833 block and the number of attention heads in the 885
834 multi-attention layer. We study how these design 886
835 choices impact the learning behavior of the PH. We 887
836 conduct a set of experiments where we gradually 888
837 increase the amount of training data or epochs and 889
838 measure the F1-score at each stopping point. This 890
839 is to simulate a training setup in production, where 891
840 the model gradually gets exposed to more training 892

841 data as the applications collect more personalized 842
843 data from the users. 844

845 *Hidden Dimension Size* Figure 4 shows the F1 846
847 score of PHs with different hidden dimensions as 848
849 they are trained with more epochs on the same 850
851 amount of training data. We experiment with 50, 852
853 100, 150, and 200 training examples per class for 854
855 25, 50, 75, and 100 epochs. Conversely, Figure 5 856
857 shows the F1 score of the same suite of PHs as they 858
859 are trained with more training examples for the 860
861 same number of training epochs. We make several 862
863 observations. 864

865 First, we observe that when exposure to training 866
867 data is limited in the early stages of training, PH 868
869 training can exhibit unpredictable behavior. This 870
871 is shown in the leftmost graphs of Figure 4 and 5, 872
873 where the model performance is not improved with 874
875 additional training data or more training epochs. 876

877 Furthermore, larger PHs perform better than 878
879 smaller PHs but with diminishing returns at higher 879
880 ends, indicating a sweet spot of PH design. We 881
882 observe 512 to be the sweet spot of PH design 882
883 for SNIPS as it performs better or similar to the 883
884 other configurations across all experiments. This 884
885 finding is corroborated with results on the Clinc- 885
886 150 dataset. Figure 6 shows the F1 score of PHs 886
887 with varying hidden dimensions on both SNIPS and 887
888 Clinc-150. We observe similar trends for dimini- 888
889 shing return in performance for Clinc as the PH 889
890 design gets larger. Similarly, 512 is the inflection 890
891 point of F1 score improvement, making it the sweet 891
892 spot PH design for Clinc. Furthermore, we observe 892
893 a slower rate of F1 score improvement with respect 893
894 to hidden dimension size for Clinc than SNIPS. 894
895 This can be explained with the observation that 895
896 Clinc is a more diverse and challenging task with 896
897 significantly more classes than SNIPS, as described 897
898 in Section 5.2, 898

899 *# of Attention Heads* We study the impact of at- 899
900 tention heads on PHs performance. Figure 7 shows 900
901 F1 score w.r.t hidden dimensions per attention head. 901
902 We follow the design in (Vaswani et al., 2017) 902
903 where the hidden dimensions of the feed-forward 903
904 layer are effectively distributed evenly among the 904
905 available attention heads. We observe that PHs 905
906 achieve better performance with higher hidden di- 906
907 mensions per head but eventually see diminishing 907
908 returns. 908

909 A.2 Scalability 889

910 We study the scalability of the PH approach and 890
911 its impact on production deployment. To help us 891
912 holistically evaluate a personalization approach, 892

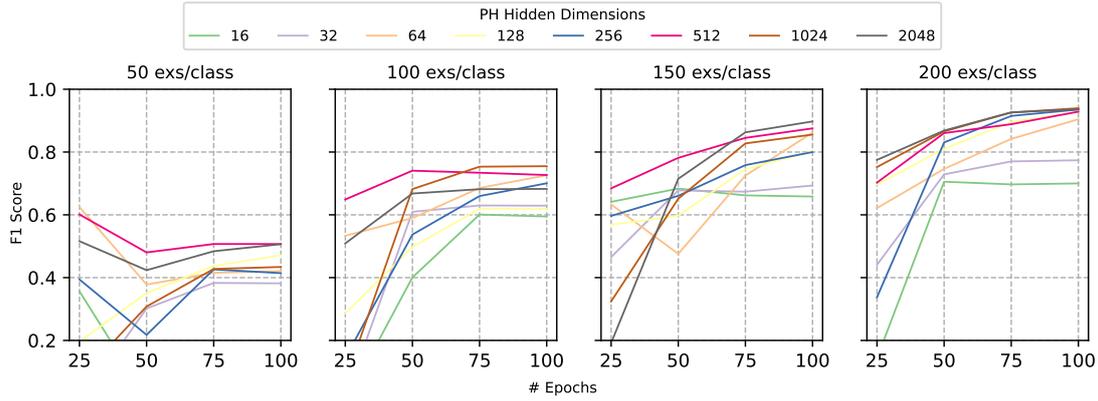


Figure 4: F1 Score w.r.t # training epochs, for fixed amounts of data.

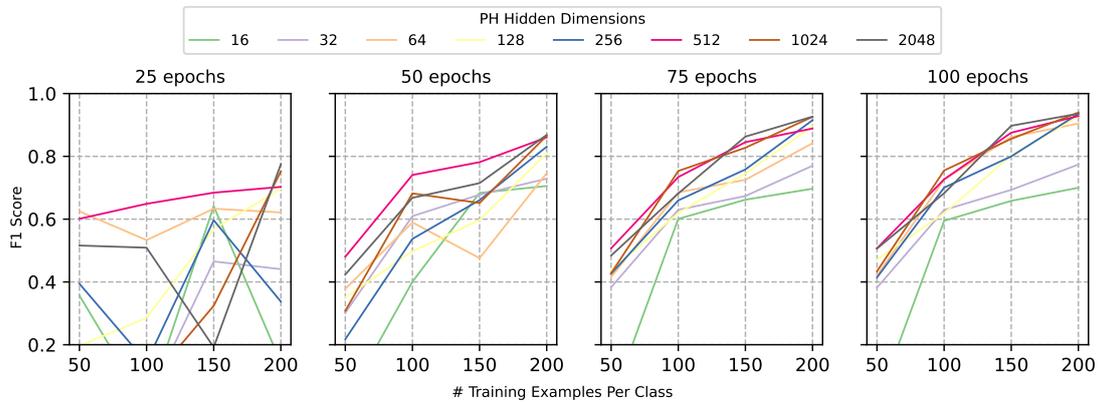


Figure 5: F1 Score w.r.t # training data, for fixed amounts of epochs.

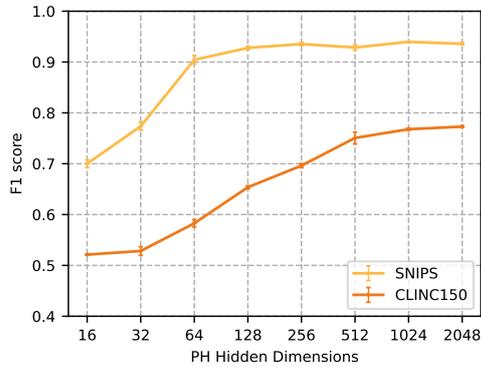


Figure 6: F1 score of PH w.r.t. hidden dimension size

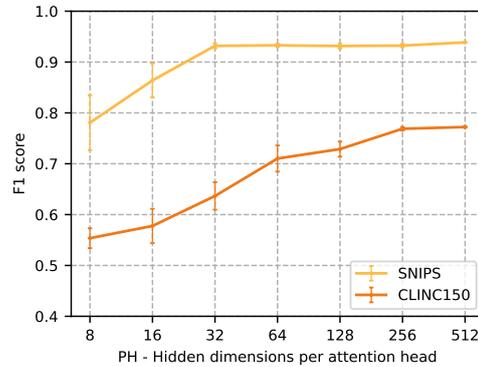


Figure 7: F1 score of PH w.r.t. hidden dims per head

we first introduce a new metric, Personalization Efficiency (PE):

$$PE = \frac{F - score^2}{Training Cost \times Model Size} \quad (9)$$

This new metric considers both the model performance and the computation requirements of training and inference. We use the number of trainable parameters as an approximation for training cost in

this study. Figure 8 shows the efficiency of 4 PH configurations normalized to the fine-tuning BERT baseline. We show that PHs achieve efficiency up to 155X compared to the fine-tuning baseline.

Furthermore, for SNIPS we observe smaller PHs generally measure higher in efficiency than larger PHs but see a diminishing return. For Clinc, 512 achieves the highest efficiency of the PHs tested. This corroborates our recommendation earlier that

Hidden Dims	# Params	Size	epoch=50	epoch=100	epoch=50	epoch=100
			# data/class=100	# data/class=100	# data/class=200	# data/class=200
2048	2.7M	10.5MB	65.92	70.43 (+4.50)	78.12 (+12.20)	92.99 (+27.70)
1024	3.2M	12.0MB	57.42	65.39 (+7.97)	80.17 (+22.75)	93.72 (+36.30)
512	3.9M	15.0MB	61.75	67.01 (+5.26)	83.72 (+21.97)	93.99 (+32.24)
256	5.5M	21.0MB	57.40	63.58 (+6.18)	79.85 (+22.45)	94.58 (+37.18)

Table 6: F1 score (and differential) with increasing training data and/or training for more epochs

512 is the sweet spot for PH design.

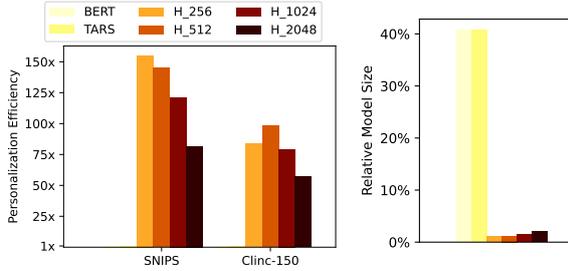


Figure 8: Personalization Efficiency

Figure 9: Data Scalability

We also quantify the potential storage overhead required for personalized models. Figure 9 shows the additional storage overhead required by the personalized models per individual user relative to the existing user data in production. We use Gmail as an example application. We calculate approximately the current per-user data usage based on a report that Gmail user creates up to 1.4MB of data per day and 3 years is the average account lifetime (ZDNet, 2012). Figure 9 shows that the proposed PHs constitute 1% - 1.5% of additional storage overhead across all 4 sizes, while the fine-tuning baselines would incur around 40% additional storage overhead per user.

A.3 More PH configurations

hidden dim	# attn. heads	SNIPS	Clinic
64	8	95.89	57.85
	4	96.32	54.41
	2	96.03	56.60
32	8	96.32	51.37
	4	96.24	50.83
	2	95.15	46.92
16	8	94.87	49.92
	4	94.43	49.99
	2	94.94	50.28
8	8	83.09	52.58
	4	80.86	52.66
	2	86.34	53.29

Table 7: F1 score of PHs with smaller sizes.

Table 7 shows the F1 score of PHs with 64, 32, 16, and 8 hidden dimensions, on SNIPS and Clinc datasets. This is an extension to the result shown in Table 5. We observe similar trends carry over to this set of even smaller PHs. This shows that even

a tiny PH can adapt LM well to the SNIPS task. On the other hand, the smaller PHs are not as effective for the more challenging Clinc datasets.

A.4 Training for more epochs on Clinc

Table 8 shows the F1 scores of PHs of 4 different sizes on Clinc when training for an additional 50 epochs (100 epochs in total). This shows that PHs performance continues to improve with more training iterations, indicating that continuing training for more iterations are beneficial in improving PH performance.

hidden dim	# attn. heads	Clinic
2048	2	78.29
	4	77.25
	8	77.78
512	2	75.82
	4	75.05
	8	74.98
128	2	65.29
	4	65.72
	8	63.63
32	2	51.35
	4	53.24
	8	52.84

Table 8: F1 score of PHs trained on Clinc-150 dataset for an additional 50 epochs (100 epochs in total)

A.5 Analyzing # of attention heads

Figure 10 and 11 analyze the impact of # attention heads on the performance of PHs. We conduct experiments similar to that in Section A.1. We gradually increase the amount of training data while holding the training epochs fixed and measure the F1 score at each stopping point, and vice versa. We observe that, compared to hidden dimension sizes, # of attention heads has less effect on the learning behavior and capacity of the PHs.

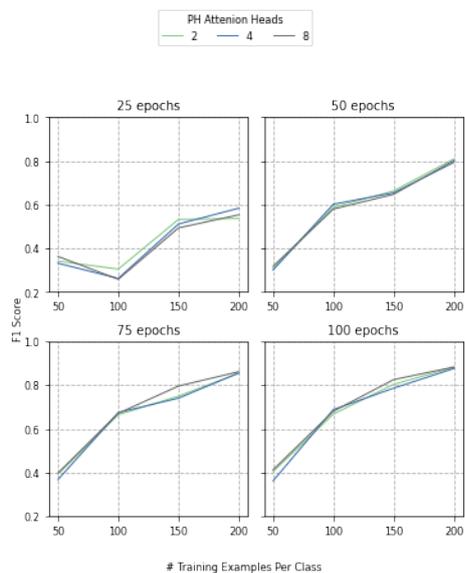


Figure 10: With the same number of training epochs, the impact of training data on performance.

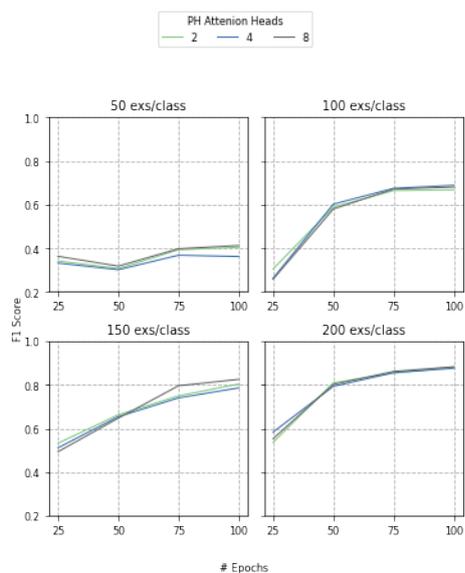


Figure 11: With the same amount of training examples per class, the impact of epochs on performance.