
FDAPT: Federated Domain-Adaptive Pre-Training for Language Models

Lekang Jiang[†], Filip Svoboda[†], Nicholas D. Lane^{†◇}

[†]University of Cambridge [◇]Flower Labs

Abstract

Foundation models (FMs) have shown prominent success in a wide range of tasks [Bommasani et al., 2021]. Their applicability to specific domain-task pairings relies on the availability of, both, high-quality data and significant computational resources. These challenges are not new to the field and, indeed, Federated Learning (FL) has been shown to be a promising solution in similar setups [Yu et al., 2023, Zhuang et al., 2023]. This paper tackles the specific case of Domain-Adaptive Pre-Training (DAPT), a key step in the application of FMs. We conduct the first comprehensive empirical study to evaluate the performance of Federated Domain-Adaptive Pre-Training (FDAPT). We demonstrate that FDAPT can maintain competitive downstream task performance to the centralized baseline in both IID and non-IID situations. Finally, we propose a novel algorithm, Frozen Federated Domain-Adaptive Pre-Training (FFDAPT). FFDAPT improves the computational efficiency by 12.1% on average and exhibits similar downstream task performance to vanilla FDAPT, with general performance fluctuations remaining less than 1%.

1 Introduction

Foundation models (FMs) are trained on broad data and can be adapted to different downstream tasks [Bommasani et al., 2021]. Recently, these models, such as GPT-4 [OpenAI, 2023] and LLaMA [Touvron et al., 2023], have demonstrated remarkable capabilities across numerous tasks and domains, especially in the field of natural language processing (NLP) [Zhou et al., 2023]. Researchers can create high-quality models by tuning FMs to specific tasks instead of building bespoke models from scratch. This approach faces challenges due to limited availability of public data [Villalobos et al., 2022] and high demand of computation power [Bommasani et al., 2021]. To solve these problems, researchers incorporate Federated Learning (FL) [McMahan et al., 2017] into FMs [Yu et al., 2023, Zhuang et al., 2023]. FL is a decentralized approach, which allows multiple clients to collaboratively train a joint model without exchanging raw data from each client. By adopting FL when training FMs, we can leverage more distributed data and computation resources across different sources to improve model performance while preserving data privacy.

In this research, we focus on the combination of FL and Domain-Adaptive Pre-Training (DAPT) [Gururangan et al., 2020]. The aim of DAPT is to adapt original FMs to new domains by continuing the pre-training task in the target domain. Consequently, the domain-specific FMs can achieve higher performance than original models, such as Clinical BERT [Alsentzer et al., 2019] and BioBERT [Lee et al., 2020] in clinical and biomedical domains respectively. Combining DAPT with FL can offer the following advantages: 1) **Privacy guarantee**. Sensitive data can be utilized without direct sharing with other clients or the server to protect data privacy. 2) **Enhanced performance**. More enhanced and powerful domain-specific Pre-trained Language Models (PLMs) can be developed by training on extensive private and distributed data. 3) **Cost saving**. Only raw text data are needed, saving the

substantial costs and expenses of data labelling. 4) **Wide applicability.** Domain-specific PLMs can be fine-tuned on any downstream tasks within the domain to improve performance.

Despite the potential advantages, almost no studies have investigated Federated Domain-Adaptive Pre-Training (FDAPT). The previous work conducted simple experiments with a fixed number of clients and limited experimental situations, showing that pre-training in federated manners is applicable with some decline in accuracy [Liu and Miller, 2020]. This initial study leaves substantial research gaps in thoroughly assessing the performance and developing novel algorithms to improve the results.

The main contributions of this research are: 1) We derive the formal definition of non-IIDness (non-independent and identically distributed) in the context of FDAPT. 2) We design a comprehensive empirical study to evaluate the performance of standard FDAPT, and conduct extensive experiments to obtain valuable results. 3) We propose Frozen Federated Domain-Adaptive Pre-Training (FFDAPT), a straightforward but effective algorithm, which improves the computation efficiency by 12.1% and remains similar performance to the vanilla FDAPT. 4) Through a critical evaluation of our work, we identify promising future research directions for this new research area.

2 Related Work

2.1 Domain-adaptive Pre-training

Language models are often pre-trained on heterogeneous corpora, such as Wikipedia, to capture general knowledge of languages. However, these models are task-agnostic and lack domain awareness. For example, linguistic characteristics between general texts and clinical narratives are different [Alsentzer et al., 2019]. Clinical narratives contain specialized medical terminology and abbreviations that are not commonly found in general text. These differences downgrade the performance of PLMs on tasks in specific domains. Hence, researchers propose the DAPT methods to adapt language models to target domains and tasks [Guo and Yu, 2022].

There are two types of DAPT approaches for PLMs. The first method is to continue the pre-training task on abundant unlabeled domain-specific texts without changing the loss functions and model structures. Studies demonstrated the effectiveness of this straightforward approach in various domains. For example, an empirical study [Gururangan et al., 2020] conducted experiments in 4 domains, including biomedical, computer science publications, news, and reviews. It illustrated that DAPT can result in performance increases under both high- and low-resource settings. In addition, domain-specific models trained by DAPT, such as Clinical BERT [Alsentzer et al., 2019] and BioBERT [Lee et al., 2020], can significantly outperform the original model on tasks in those domains. The alternative DAPT method is to add domain-distinguish pre-training tasks to the original training objective. This approach is more complicated but can be more effective for domain adaptation. For example, studies showed that incorporating adversarial domain discrimination into DAPT can enhance domain-invariant features [Du et al., 2020], which further improves the model performance.

2.2 Federated Pre-training

Researchers trained Word2Vec [Mikolov et al., 2013] from scratch using FL on Wikipedia and compared the performance with centralized Word2Vec [Bernal et al., 2021]. The results showed that combining Word2Vec and FL can achieve enhanced word representations, with better similarity, analogy, and categorization outcomes. Both model quality and convergence time in FL settings are comparable to centralized training.

FedBERT [Tian et al., 2022] incorporated Split Learning [Gupta and Raskar, 2018] with FL to resolve computational constraints in cross-device settings. It showed that the transformer layer is computationally costly to train on edge devices. Thus, it adopted Split Learning to update the transformer layer on the powerful server and conducted other layers on edge devices. The results demonstrated that Federated Split Learning can reach same performance compared to standard FL.

While the above work focused on pre-training in the general domain, a proof-of-concept study conducted experiments in the DAPT process [Liu and Miller, 2020]. It continued pre-training BERT [Devlin et al., 2018] using FL with 5 clients on the clinical dataset and tested on 2 downstream tasks. Notably, different situations of data distributions are not considered in the experiments. The outcome

indicated that pre-training and fine-tuning of BERT are applicable to FL settings with some decline in accuracy.

3 Methodology

3.1 Experimental Setup

Training and Evaluation. In the FDAPT process, we initialize each client’s model with the weights of a PLM and continue training on domain-specific datasets using the same pre-training tasks under various federated settings. Specifically, we adopt the Federated Averaging (FedAvg) [McMahan et al., 2017] algorithm for our experiments. To evaluate domain-specific PLMs, we fine-tune them on different downstream tasks within the domain and compare their performance. Details of FDAPT system design are illustrated in Appendix A and evaluation metrics are introduced in Appendix B.

Model Architecture and Framework. As the purpose of this empirical study is to demonstrate the effectiveness of FDAPT in comprehensive experimental settings we chose the DistilBERT [Sanh et al., 2019] model due to its efficiency and representativeness. Although GPT series models [OpenAI, 2023] have achieved incredible performance recently, their enormous model sizes would limit the scale of our experiments substantially. Furthermore, we used Flower [Beutel et al., 2020], a user-friendly FL framework, for experimental simulations. Flower provides a flexible and accessible environment to easily customize FL configurations, simplifying the implementation and execution of FL experiments.

Datasets. We apply FDAPT to the PubMed dataset and 9 domain-specific datasets listed in Table 1. The PubMed dataset contains abstracts and full-text research articles in the biomedical domain. The PubMed dataset adopted in our experiments [Cohan et al., 2018] is smaller than the version used for BioBERT [Lee et al., 2020], which minimizes the experimental costs while maintaining pronounced results. The domain-specific PLMs are evaluated on 9 publicly available downstream tasks in the biomedical domain, including 6 Named Entity Recognition (NER), 2 Relation Extraction (RE) and 1 Question Answering (QA) datasets, as shown in Table 1.

Table 1: List of downstream tasks.

Datasets	Task type	Entity type
NCBI Disease Doğan et al. [2014]	NER	Disease
BC5CDR Li et al. [2016]	NER	Chemical
BC4CHEMD Krallinger and et al. [2015]	NER	Chemical
BC2GM Smith et al. [2008]	NER	Gene
LINNAEUS Gerner et al. [2010]	NER	Species
Species-800 Pafilis et al. [2013]	NER	Species
GAD Bravo et al. [2015]	RE	Gene-disease
EU-ADR Van Mulligen et al. [2012]	RE	Gene-disease
BioASQ 7b-factoid Tsatsaronis et al. [2015]	QA	N/A

3.2 Non-IIDness in Federated Pre-training

Existing research in federated pre-training methods only focused on IID situations and ignored the non-IID issue, which commonly exists in practical FL applications and can cause performance drops. Non-IIDness in supervised learning is usually related to the distribution of data labels and features [Kairouz et al., 2021]. However, datasets used for pre-training only contain raw texts, which do not have pre-defined labels or features. Therefore, we define three types of non-IIDness in the federated pre-training process, including quantity skew, sentence length distribution skew and vocabulary distribution skew.

Quantity skew refers to the imbalance in the number of training data across different clients. In the context of federated pre-training, we define it as the imbalance of the number of raw texts among all clients. Additionally, sentence length and the number of vocabulary are important features of text data. Therefore, we argue that federated pre-training with imbalanced average sentence length or the

number of vocabularies may affect the final results. We define the **sentence length distribution skew** as the imbalance of the average sentence length across multiple clients, and **vocabulary distribution skew** as the imbalance of the number of unique words across all clients. When creating these skews, the aim is to maximize a single metric discrepancy among all clients, while keeping other metrics almost the same. Detailed description of non-IIDness in federated pre-training is included in Appendix C. The data distribution for experiments is reported in Appendix D.

3.3 Frozen Federated Domain-Adaptive Pre-Training

The pre-training process often requires extensive computational resources. For example, it took 23 days to train the domain-specific BioBERT model on eight V100 GPUs [Lee et al., 2020]. Hence, mitigating computation costs is important in the FDAPT process.

Freezing specific layers or parameters of a model during fine-tuning or transfer learning aims to preserve the pre-trained knowledge encoded within those layers, while allowing other parts of the model to adapt to new tasks or domains. Studies showed that fine-tuning only a fourth of model layers can achieve 90% of the original model quality [Lee et al., 2019]. Therefore, we propose FFDAPT, a simple but effective method that incorporates freezing methods into FDAPT to improve efficiency. We illustrate the FFDAPT approach in Algorithm 1.

Algorithm 1: Frozen Federated Domain-Adaptive Pre-Training (FFDAPT)

Input : Initialized weights of N -layer model $\{W_{k,1}, \dots, W_{k,N}\}_{k=1}^K$ from K clients, number of training samples for each client $\{n_1, \dots, n_K\}$, number of training round T , maximum number of frozen layers ε , scaling parameter γ

Output : Final global weights $\{W_1, \dots, W_N\}$

$start = 1, end = 1;$ // Frozen layer index

$n = \sum_{k=1}^K n_k;$

for $t \in \{1, 2, \dots, T\}$ **do**

for $k \in \{1, 2, \dots, K\}$ **do**

$N_k = \min(\varepsilon, \lceil \frac{n_k}{n} N \rceil \gamma);$

$end = start + N_k;$

if $end \leq N$ **then**

 Train $\{W_{k,1}, \dots, W_{k,N}\}$ with $\{W_{k,start}, \dots, W_{k,end}\}$ frozen;

else

$end = end - N;$

 Train $\{W_{k,1}, \dots, W_{k,N}\}$ with $\{W_{k,start}, \dots, W_{k,N}\}, \{W_{k,1}, \dots, W_{k,end}\}$ frozen;

$start = end + 1;$

if $start > N$ **then**

$start = start - N;$

$\{W_1, \dots, W_N\} \leftarrow \sum_{k=1}^K \frac{n_k}{n} \{W_{k,1}, \dots, W_{k,N}\};$

return $\{W_1, \dots, W_N\}$

For each client, a portion of the consecutive layers is frozen, and the number of frozen layers is determined by the size of the dataset. We set a scaling hyper-parameter γ to control the actual number of frozen layers and improve flexibility. Since freezing all layers are meaningless during training, ε is set to be the maximum number of frozen layers. Through our method, each client freezes some layers during training to improve the computation efficiency.

4 Results

We report all the experimental results in Table 2. We demonstrate that FDAPT achieves competitive performance on downstream tasks against the centralized baseline in both IID and non-IID situations. Moreover, FFDAPT improves the training efficiency by approximately 12.1% on average, while maintaining similar downstream tasks performance compared to vanilla FDAPT.

Table 2: Downstream tasks performance and standard deviations of the original model, the model with centralized pre-training, models with FDAPT and models with FFDAPT under IID and non-IID settings.

Settings	NER						RE		QA
	NCBI	BC5CDR	BC4CHEMD	BC2GM	LINNAEUS	Species-800	GAD	EU-ADR	BioASQ 7b
DistilBERT	84.6 (.4)	84.1 (.3)	86.9 (.1)	79.6 (.2)	83.5 (.8)	67.1 (.9)	75.8 (.4)	80.0 (.3)	24.9 (1.2)
Centralised	85.4 (.9)	86.9 (.1)	88.3 (.2)	80.8 (.4)	84.6 (1.6)	69.3 (.6)	78.0 (.1)	80.6 (.1)	27.0 (1.4)
FDAPT									
IID									
2 Clients	85.7 (.8)	86.7 (.2)	88.0 (.1)	80.7 (.2)	85.4 (1.3)	68.2 (1.2)	77.7 (.3)	80.7 (.2)	28.0 (1.2)
8 Clients	85.2 (.7)	86.1 (.2)	87.5 (.1)	80.6 (.2)	83.8 (.8)	67.8 (1.1)	77.2 (.6)	80.1 (.2)	28.5 (1.2)
Non-IID (Quantity Skew)									
2 Clients	85.8 (.9)	86.4 (.2)	88.2 (.1)	80.5 (.2)	85.1 (.7)	69.0 (.9)	78.0 (.8)	80.1 (.2)	28.3 (1.4)
8 Clients	85.1 (.5)	86.0 (.2)	87.6 (.1)	80.1 (.5)	84.7 (1.4)	69.2 (.6)	77.5 (.2)	80.8 (1.2)	29.0 (.8)
Non-IID (Sentence Length Distribution Skew)									
2 Clients	85.8 (.6)	86.5 (.2)	87.9 (.1)	80.8 (.3)	83.8 (1.4)	69.3 (.4)	78.0 (.2)	80.0 (.8)	27.3 (1.6)
8 Clients	85.2 (.2)	86.2 (.1)	87.6 (.1)	80.5 (.5)	84.5 (1.2)	68.9 (2.0)	78.2 (.1)	81.7 (1.0)	29.4 (1.0)
Non-IID (Vocabulary Distribution Skew)									
2 Clients	85.8 (1.2)	86.7 (.2)	88.2 (.2)	80.7 (.3)	84.6 (1.2)	69.4 (1.4)	77.4 (.4)	80.7 (.5)	26.7 (1.2)
8 Clients	85.3 (.4)	86.1 (.1)	87.7 (.2)	80.7 (.2)	83.7 (.8)	68.8 (1.2)	78.0 (.1)	79.8 (.9)	28.8 (1.4)
FFDAPT									
IID									
2 Clients	84.8 (.9)	86.3 (.2)	87.6 (.1)	80.3 (.1)	85.0 (1.8)	68.5 (1.3)	77.0 (.3)	79.6 (.4)	28.4 (2.3)
8 Clients	85.0 (.7)	85.9 (.2)	87.4 (.2)	80.1 (.2)	83.8 (1.2)	68.6 (.8)	77.3 (.1)	80.8 (.1)	26.9 (.9)
Non-IID (Quantity Skew)									
2 Clients	85.2 (.9)	86.2 (.2)	87.7 (.1)	80.4 (.2)	83.7 (1.0)	67.7 (.7)	78.5 (.1)	80.4 (.4)	28.2 (1.2)
8 Clients	84.5 (1.0)	86.0 (.2)	87.4 (.1)	80.2 (.5)	83.5 (1.5)	68.4 (.7)	78.2 (.2)	80.9 (.1)	28.3 (1.3)
Non-IID (Sentence Length Distribution Skew)									
2 Clients	85.6 (.4)	86.4 (.1)	87.5 (.0)	80.3 (.3)	84.1 (1.4)	68.5 (.8)	78.2 (.1)	80.7 (.3)	27.4 (1.0)
8 Clients	85.5 (.6)	85.9 (.3)	87.5 (.2)	80.6 (.3)	83.6 (2.0)	68.7 (1.2)	77.6 (.7)	80.6 (.5)	27.0 (0.9)
Non-IID (Vocabulary Distribution Skew)									
2 Clients	84.8 (1.0)	86.2 (.3)	87.6 (.1)	80.3 (.2)	84.1 (1.7)	69.7 (.8)	76.7 (.8)	79.0 (.3)	28.3 (1.7)
8 Clients	85.5 (.6)	85.9 (.3)	87.5 (.2)	80.6 (.3)	83.6 (2.0)	68.7 (1.2)	77.2 (.1)	81.2 (.0)	28.2 (.6)

Note: Each experiment is conducted with 5 random seeds, and the average value is reported to eliminate randomness. The evaluation metrics used in the table are F1 scores for NER and RE tasks and strict accuracy for the QA dataset. We also evaluate the precision and recall for NER and RE tasks, lenient accuracy and mean reciprocal rank for the QA dataset. The results of these metrics provide similar research outcomes.

4.1 FDAPT

We make the following observations regarding FDAPT.

The performance drops of federated models compared to the centralized baseline are acceptable. While federated models exhibit a slight performance decrease compared to the centralized model on some tasks, all models pre-trained with FDAPT surpass the original model on all tasks. These datasets contain both word-level and sentence-level tasks, underscoring the effectiveness of FDAPT. On the other hand, the performance of federated models decreases by less than 1% on almost all datasets compared to the centralized baseline. Hence, we argue that this level of performance reduction is tolerable to preserve data privacy.

Federated models occasionally outperform the centralized baseline. Federated models trained in specific settings can outperform the centralized approach on NCBI, LINNAEUS, Species-800, GAD, EU-ADR and BioASQ 7b, including both token-level and sentence-level tasks. Notably, FDAPT with 2 clients under the IID setting increases the F1 scores on LINNAEUS by 0.8%. Meanwhile, on the EU-ADR dataset, FDAPT with 8 clients under non-IID settings of sentence length skew, achieves the F1 score of 81.7%, which is 1.1% higher than the centralized approach. Furthermore, federated models outperform the centralized model on BioASQ 7b with increases up to 2.4% in terms of strict accuracy.

4.2 FFDAPT

The aim of FFDAPT is to improve computational efficiency. The improvement of training efficiency is calculated by the following equation.

$$I = \frac{T - T_F}{T_F} \cdot 100\% \quad (1)$$

where T and T_F are the round time for standard FDAPT and FFDAPT respectively. We calculate the improvement of efficiency for each experimental scenario, and report the average value of 12.1% as the final result of improvement.

In terms of downstream tasks performance, **FFDAPT can achieve similar outcomes compared to standard FDAPT in all situations.** The performance variations are no greater than 1% for most of

the cases. Remarkably, FFDAPT leads to higher performance than vanilla FDAPT in some specific scenarios. For example, the F1 scores on the Species-800 dataset increase by 0.8% for models trained with 8 clients under the IID setting. Additionally, FFDAPT with 8 clients improves the F1 scores by 1.4% on the EU-ADR dataset in non-IID settings of vocabulary distribution skew. Moreover, FFDAPT with 2 clients results in a 1.6% performance enhancement on the QA task in the situation of vocabulary distribution skew. Although there are performance decreases in some situations, they still remain better than the original model.

5 Discussion and Future Work

This study makes significant contributions to FDAPT, but we acknowledge some limitations and point out promising research directions to stimulate future studies.

More real-world simulations. Due to limited data and computational resources, our experiments focused on a particular set of common settings to obtain meaningful results while minimizing experimental costs. As it is possible that interesting connections could be discovered between model-domain pairs, and the process of federation, future research could expand on ours by attempting larger-scale simulations on an expanded selection of alternative model backbones and domain datasets. In addition, more complex non-IID simulations can be conducted, such as the skewness in the distribution of sentence embeddings, which is crucial in determining model quality.

Improve computation and communication efficiency. We demonstrate the efficacy of our FFDAPT algorithm, but further efficiency improvements may be possible. For instance, it is likely that module adapters [Cai et al., 2022] could be used to mitigate computational costs, and communication-efficient algorithms, such as FedPCL [Tan et al., 2022], could improve communication efficiency.

Other challenges in FL. Beyond the concerns of data privacy, addressed in this paper, Federated Learning can be used to address other concerns of distributed training. For example, it can address system heterogeneity, develop specific federated strategies and integrate privacy-enhancing approaches. All of these are likely to be fruitful extensions to this work.

Domain-related challenges. This research is based on the biomedical domain and suggests that FDAPT can have similar performance compared to centralized pre-training. In contrast, a previous study [Liu and Miller, 2020] showed that FDAPT is worse than centralized pre-training in the clinical domain. Our results are based on 9 diverse downstream datasets, in contrast to the alternative assessment, which was based on 2, highly similar, tasks. The difference in results is worth investigating, as, in the most interesting case, this could point to a special case identified by [Liu and Miller, 2020], in which the specific task they chose happens to be harder to federate.

6 Conclusion

In this paper, we presented the first comprehensive empirical investigation of FDAPT for NLP tasks in cross-silo settings. We trained the DistilBERT model on biomedical corpora using FDAPT, and tested the model performance on 9 different task-specific datasets, including both token-level and sentence-level tasks. We formalized three types of non-IIDness in the context of FDAPT, and simulated both IID and non-IID data distributions. The results show that FDAPT can retain competitive downstream task performance to the centralized baseline across all tested scenarios. The performance drops of FDAPT are less than 1%, most of the time, and it can occasionally outperform the centralized approach. In all experimental situations, models trained by FDAPT surpass the original, non-adapted, DistilBERT model, demonstrating the effectiveness of FDAPT.

Furthermore, we propose our own efficiency-improving algorithm, the FFDAPT. It can increase training speed by approximately 12.1%, on average. Additionally, FFDAPT exhibits similar downstream task performance to vanilla FDAPT, with general performance fluctuations remaining less than 1%. In specific circumstances, FFDAPT can improve the performance by up to 1.6%.

Finally, we point out promising future research directions for this new research area. It is crucial to continue investigating this field, in order to advance our understanding of how Federated Learning can aid in domain adaptation, in particular, and the applicability of Foundation Models, in general.

References

- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- Sixing Yu, J Pablo Muñoz, and Ali Jannesari. Federated foundation models: Privacy-preserving and collaborative learning for large models. *arXiv preprint arXiv:2305.11414*, 2023.
- Weiming Zhuang, Chen Chen, and Lingjuan Lyu. When foundation model meets federated learning: Motivations, challenges, and future directions. *arXiv preprint arXiv:2306.15546*, 2023.
- OpenAI. GPT-4 Technical Report. *arXiv e-prints*, art. arXiv:2303.08774, March 2023. doi: 10.48550/arXiv.2303.08774.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben Yan, Lifang He, et al. A comprehensive survey on pretrained foundation models: A history from bert to chatgpt. *arXiv preprint arXiv:2302.09419*, 2023.
- Pablo Villalobos, Jaime Sevilla, Lennart Heim, Tamay Besiroglu, Marius Hobbhahn, and Anson Ho. Will we run out of data? an analysis of the limits of scaling datasets in machine learning. *arXiv preprint arXiv:2211.04325*, 2022.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- 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. *arXiv preprint arXiv:2004.10964*, 2020.
- Emily Alsentzer, John R Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew McDermott. Publicly available clinical bert embeddings. *arXiv preprint arXiv:1904.03323*, 2019.
- 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. *Bioinformatics*, 36(4):1234–1240, 2020.
- Dianbo Liu and Tim Miller. Federated pretraining and fine tuning of bert using clinical notes from multiple silos. *arXiv preprint arXiv:2002.08562*, 2020.
- Xu Guo and Han Yu. On the domain adaptation and generalization of pretrained language models: A survey. *arXiv preprint arXiv:2211.03154*, 2022.
- Chunning Du, Haifeng Sun, Jingyu Wang, Qi Qi, and Jianxin Liao. Adversarial and domain-aware bert for cross-domain sentiment analysis. In *Proceedings of the 58th annual meeting of the Association for Computational Linguistics*, pages 4019–4028, 2020.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- Daniel Garcia Bernal, Lodovico Giarretta, Sarunas Girdzijauskas, and Magnus Sahlgren. Federated word2vec: Leveraging federated learning to encourage collaborative representation learning. *arXiv preprint arXiv:2105.00831*, 2021.
- Yuanyishu Tian, Yao Wan, Lingjuan Lyu, Dezhong Yao, Hai Jin, and Lichao Sun. Fedbert: When federated learning meets pre-training. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2022.

- Otkrist Gupta and Ramesh Raskar. Distributed learning of deep neural network over multiple agents. *Journal of Network and Computer Applications*, 116:1–8, 2018.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *ArXiv*, abs/1910.01108, 2019.
- Daniel J Beutel, Taner Topal, Akhil Mathur, Xinchu Qiu, Javier Fernandez-Marques, Yan Gao, Lorenzo Sani, Kwong Hei Li, Titouan Parcollet, Pedro Porto Buarque de Gusmão, et al. Flower: A friendly federated learning research framework. *arXiv preprint arXiv:2007.14390*, 2020.
- Arman Cohan, Franck Dernoncourt, Doo Soon Kim, Trung Bui, Seokhwan Kim, Walter Chang, and Nazli Goharian. A discourse-aware attention model for abstractive summarization of long documents. *arXiv preprint arXiv:1804.05685*, 2018.
- Rezarta Islamaj Doğan, Robert Leaman, and Zhiyong Lu. Ncbi disease corpus: a resource for disease name recognition and concept normalization. *Journal of biomedical informatics*, 47:1–10, 2014.
- Jiao Li, Yueping Sun, Robin J. Johnson, Daniela Sciaky, Chih-Hsuan Wei, Robert Leaman, Allan Peter Davis, Carolyn J. Mattingly, Thomas C. Wieggers, and Zhiyong Lu. Biocreative V CDR task corpus: a resource for chemical disease relation extraction. *Database J. Biol. Databases Curation*, 2016, 2016. doi: 10.1093/database/baw068. URL <https://doi.org/10.1093/database/baw068>.
- Martin Krallinger and et al. The chemdner corpus of chemicals and drugs and its annotation principles. *Journal of Cheminformatics*, 7:S2 – S2, 2015.
- Larry Smith, Lorraine K Tanabe, Cheng-Ju Kuo, I Chung, Chun-Nan Hsu, Yu-Shi Lin, Roman Klinger, Christoph M Friedrich, Kuzman Ganchev, Manabu Torii, et al. Overview of biocreative ii gene mention recognition. *Genome biology*, 9(2):1–19, 2008.
- Martin Gerner, Goran Nenadic, and Casey M Bergman. LINNAEUS: A species name identification system for biomedical literature. *BMC Bioinformatics*, 11(1):85, 2010. ISSN 1471-2105. doi: 10.1186/1471-2105-11-85. URL <https://doi.org/10.1186/1471-2105-11-85>.
- Evangelos Pafilis, Sune P Frankild, Lucia Fanini, Sarah Faulwetter, Christina Pavloudi, Aikaterini Vasileiadou, Christos Arvanitidis, and Lars Juhl Jensen. The species and organisms resources for fast and accurate identification of taxonomic names in text. *PLoS one*, 8(6):e65390, 2013.
- Àlex Bravo, Janet Piñero, Núria Queralt-Rosinach, Michael Rautschka, and Laura I Furlong. Extraction of relations between genes and diseases from text and large-scale data analysis: implications for translational research. *BMC Bioinformatics*, 16(1), February 2015. doi: 10.1186/s12859-015-0472-9. URL <https://doi.org/10.1186/s12859-015-0472-9>.
- Erik M Van Mulligen, Annie Fourier-Reglat, David Gurwitz, Mariam Molokhia, Ainhua Nieto, Gianluca Trifiro, Jan A Kors, and Laura I Furlong. The eu-adr corpus: annotated drugs, diseases, targets, and their relationships. *Journal of biomedical informatics*, 45(5):879–884, 2012.
- George Tsatsaronis, Georgios Balikas, Prodromos Malakasiotis, Ioannis Partalas, Matthias Zschunke, Michael R Alvers, Dirk Weissenborn, Anastasia Krithara, Sergios Petridis, Dimitris Polychronopoulos, et al. An overview of the bioasq large-scale biomedical semantic indexing and question answering competition. *BMC bioinformatics*, 16(1):138, 2015.
- Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.
- Jaejun Lee, Raphael Tang, and Jimmy Lin. What would elsa do? freezing layers during transformer fine-tuning. *arXiv preprint arXiv:1911.03090*, 2019.
- Dongqi Cai, Yaozong Wu, Shangguang Wang, Felix Xiaozhu Lin, and Mengwei Xu. Autofednlp: An efficient fednlp framework. *arXiv preprint arXiv:2205.10162*, 2022.

Yue Tan, Guodong Long, Jie Ma, Lu Liu, Tianyi Zhou, and Jing Jiang. Federated learning from pre-trained models: A contrastive learning approach. *arXiv preprint arXiv:2209.10083*, 2022.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

Balu Bhasuran and Jeyakumar Natarajan. Automatic extraction of gene-disease associations from literature using joint ensemble learning. *PLoS one*, 13(7):e0200699, 2018.

A FDAPT System Design

We present an overview of the FDAPT method in Figure 1. The entire process consists of three stages. Firstly, models are pre-trained on large heterogeneous corpora, such as Wikipedia, to acquire general knowledge. These models are often trained by AI research companies and are publicly available, for example, BERT Devlin et al. [2018] and GPT-2 Radford et al. [2019].

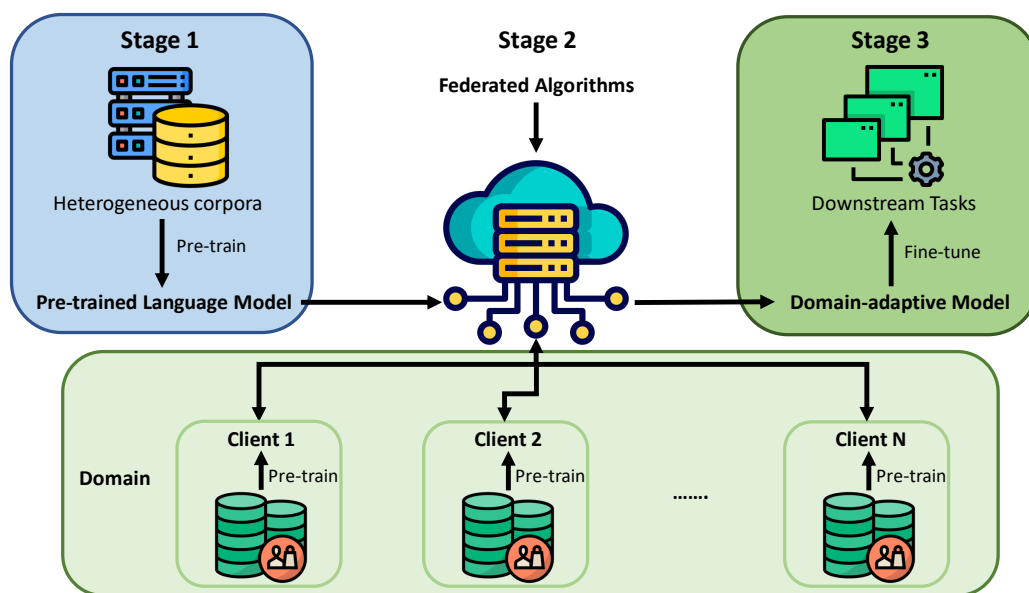


Figure 1: Overview of Federated Domain-Adaptive Pre-Training (FDAPT). Stage 1: Pre-train on heterogeneous corpora, e.g., Wikipedia. Stage 2: Apply FL in the Domain-Adaptive Pre-Training process, e.g., in the biomedical domain. Stage 3: Fine-tune on downstream tasks in the domain, e.g., disease name recognition task.

In the second stage, models continue being trained in the federated settings to adapt to specific domains. Each client trains the model on their local dataset and communicates with the server for aggregation. After a specified number of training rounds, the domain-adaptive PLMs are obtained. During this process, local datasets stored in each client are not shared to preserve data privacy, and models are aggregated to reach an enhanced global model. A group of institutes with sensitive data in the same domain can collaborate on FDAPT and release the federated domain-specific model for public usage.

In the final stage, these models are fine-tuned on different domain-specific downstream tasks, which can achieve better performance than the original model. Even individuals or small companies that work in the same domain can fine-tune the domain-specific model on their own tasks and benefit from FDAPT.

B Evaluation Metrics

NER and RE are multi-label classification tasks, which can be evaluated by **precision (P)**, **recall (R)** and **F1 scores (F1)**.

$$P = \frac{TP}{TP + FP} \quad (2)$$

$$R = \frac{TP}{TP + FN} \quad (3)$$

$$F1 = \frac{2PR}{P + R} \quad (4)$$

where TP (true positives) refers to the number of samples that are predicted positive and actually positive, FP (false positives) is the number of samples that are predicted positive but actually negative, and FN (false negatives) is the number of samples that are predicted negative but actually positive.

For each question in the factoid QA task, models can return a list of answers, ordered by decreasing confidence. We evaluate the model by strict accuracy (S), lenient accuracy (L) and mean reciprocal rank (M).

Strict accuracy measures the proportion of questions, for which the first returned answer by the model contains the golden response. Assuming the golden answers are $[g_1, g_2, \dots, g_n]$, and the predicted answers for question k are $[a_{k,1}, a_{k,2}, \dots, a_{k,m}]$, the strict accuracy can be calculated using Equation 5.

$$S = \frac{\sum_{k=1}^n (g_k \in a_{k,1})}{n} \quad (5)$$

Lenient accuracy is the proportion of factoid questions that have been answered correctly in the lenient sense. The results are seen as correct if one of the returned answers meets the golden standard. The formula is shown in Equation 6.

$$L = \frac{\sum_{k=1}^n (g_k \in [a_{k,1}, a_{k,2}, \dots, a_{k,m}])}{n} \quad (6)$$

Mean reciprocal rank is defined in Equation 7. It ranges from 0 to 1, and higher values indicate better performance. This evaluation metric encourages the model to return accurate answers with higher confidence.

$$M = \frac{1}{n} \sum_{k=1}^n \frac{1}{r_k} \quad (7)$$

where r_k is the position of the first answer returned by the model that includes the golden response for question k

C Non-IIDness in Federated Pre-training

Quantity skew refers to the imbalance in the number of training data across different clients. In the context of federated pre-training, we define it as the imbalance of the number of raw texts among all clients. In our experiments, the data distribution for each client $i \in \{1, 2, \dots, k\}$ in the situation of quantity skew is calculated using Equation 8, where Q is the total number of training data.

$$D_Q = \left\{ Q_i \mid Q_i = \frac{i}{\sum_{j=1}^k j} Q, i \in \{1, 2, \dots, k\} \right\} \quad (8)$$

Sentence length and the number of vocabulary are important features of text data. Therefore, we argue that federated pre-training with imbalanced average sentence length or the number of vocabularies may affect the final results.

We define the **sentence length distribution skew** as the imbalance of the average sentence length across multiple clients. Assuming the average sentence length of client $i \in \{1, 2, \dots, k\}$ is L_i , we

create the sentence length distribution skew by maximising the standard deviation of average sentence lengths of all clients.

$$D_L = \left\{ L_i \mid \max(\sigma(L_1, L_2, \dots, L_k)), i \in \{1, 2, \dots, k\} \right\} \quad (9)$$

Similarly, **vocabulary distribution skew** is the imbalance of the number of unique words across all clients. The data distribution of all clients in the situation of vocabulary distribution skew is illustrated in the following Equation.

$$D_V = \left\{ V_i \mid \max(\sigma(V_1, V_2, \dots, V_k)), i \in \{1, 2, \dots, k\} \right\} \quad (10)$$

When creating these skews, the aim is to maximise a single metric discrepancy among all client, while keeping other metrics almost the same. For example, the objective of vocabulary distribution skew is to maximise $\sigma(V_1, V_2, \dots, V_k)$, the standard deviation of the number of vocabularies across all clients, while minimising $\sigma(Q_1, Q_2, \dots, Q_k)$ and $\sigma(L_1, L_2, \dots, L_k)$, the standard deviation of data quantity and average sentence length.

D Data Distribution

We report the pre-training data distribution under different settings across 2 and 8 clients in Table D. In IID settings, the text quantity, average sentence length and average vocabulary are consistent across all clients. In each non-IID setting, only one metric is skewed, while the remaining metrics are almost uniformly distributed. For example, in quantity skew, the number of texts varies significantly among clients, but the other metrics remain stable.

Table 3: Data distribution across 2 and 8 clients.

Settings	Quantity		Sentence length		Vocabulary	
	Average	σ	Average	σ	Average	σ
2 Clients						
IID	60K	0	34.3	0	3.1K	0
Quantity skew	60K	20K	34.3	0.1	3.1K	0
Sentence length distribution skew	60K	0	34.3	4.7	3.1K	0
Vocabulary distribution skew	60K	0	34.3	0.7	3.1K	1.3K
8 Clients						
IID	15K	0	34.3	0	3.1K	0
Quantity skew	15K	7.6K	34.3	0.1	3.1K	0
Sentence length distribution skew	15K	0	34.3	6.3	3.1K	0.1K
Vocabulary distribution skew	15K	0	34.3	1.7	3.1K	1.7K

Note: The metrics are the number of articles, average sentence length and average number of vocabularies. The average value and standard deviations σ across 2 clients are listed. Larger σ indicates higher value discrepancy among all clients.

E Environmental Settings

E.1 Pre-training

The federated models are trained for 15 rounds, and each client trains model for 1 epoch on the local dataset during each round, which is consistent with the previous work Liu and Miller [2020]. To make a sensible comparison, the centralised model is trained for 15 epochs. We keep the hyper-parameters unchanged when processing different experiments, using the Adam optimiser with batch size of 8 and learning rate of $5e-5$. We only train each model once, because pre-training is very computationally expensive (see Section F for details).

E.2 Fine-tuning

We fine-tune the models on different downstream tasks with different hyper-parameters, base on a previous study Lee et al. [2020].

NER. We fine-tune models on NER tasks for 20 epochs, using Adam optimiser with batch size of 8 and learning rate of $5e-5$. Models with the highest validation performance are selected and tested on the test sets to obtain the final results.

RE. Since RE datasets are smaller, we fine-tune the models for 10 epochs, using Adam optimiser with batch size of 32 and learning rate of $5e-5$. These datasets do not have separate test sets, so we reported the performance of 10-fold cross-validation according to previous works Lee et al. [2020], Bhasuran and Natarajan [2018]. The dataset is divided into ten equal folds, where one fold is used for testing and the remaining nine are used for training. This process is repeated for 10 times, and the average performance is reported.

QA. We train models on QA task for 10 epochs because of the limited size of dataset. Other settings are the same as we used for NER.

Since fine-tuning is much more computationally efficient than pre-training, all models are trained with 5 random seeds¹. The average results and standard deviations are reported to eliminate randomness.

F Compute Details

All the pre-training experiments are conducted on RTX 2080 Ti GPUs. The total GPU running time for pre-training is about 4640 hours. Notably, FDAPT and FFDAPT are trained with two GPUs to accelerate the training process and make a fair comparison of computational efficiency. Pre-trained models, including 1 original model, 1 centralised model and 16 federated models, are fine-tuned on downstream tasks using a single V100 GPU. Each model takes about 14 hours to finish fine-tuning on all downstream tasks. Thus, the overall GPU running time for fine-tuning is approximately $18 \times 14 = 252$ hours.

¹Specifically, we use 42, 123, 3407, 43534 and 54354 for the experiments.