

Towards Practical Few-shot Federated NLP

Dongqi Cai

Beiyou Shenzhen Institute

Yaozong Wu

Beiyou Shenzhen Institute

Haitao Yuan

Beiyou Shenzhen Institute

Shangguang Wang

Beiyou Shenzhen Institute

Felix Xiaozhu Lin

University of Virginia

Mengwei Xu

Beiyou Shenzhen Institute

ABSTRACT

Transformer-based pre-trained models have emerged as the predominant solution for natural language processing (NLP). Fine-tuning such pre-trained models for downstream tasks often requires a considerable amount of labeled private data. In practice, private data is often distributed across heterogeneous mobile devices and may be prohibited from being uploaded. Moreover, well-curated labeled data is often scarce, presenting an additional challenge. To address these challenges, we first introduce a data generator for federated few-shot learning tasks, which encompasses the quantity and skewness of scarce labeled data in a realistic setting. Subsequently, we propose AUG-FedPrompt, a **prompt**-based federated learning system that exploits abundant unlabeled data for data **augmentation**. Our experiments indicate that AUG-FedPrompt can perform on par with full-set fine-tuning with a limited amount of labeled data. However, such competitive performance comes at a significant system cost.

CCS CONCEPTS

• **Human-centered computing** → **Ubiquitous and mobile computing**; • **Computing methodologies** → **Machine learning**.

KEYWORDS

Federated Learning, Natural Language Processing, Few-shot Learning

ACM Reference Format:

Dongqi Cai, Yaozong Wu, Haitao Yuan, Shangguang Wang, Felix Xiaozhu Lin, and Mengwei Xu. 2023. Towards Practical Few-shot Federated NLP. In *3rd Workshop on Machine Learning and Systems*

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org. *EuroMLSys '23, May 8, 2023, Rome, Italy*
© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0084-2/23/05...\$15.00

<https://doi.org/10.1145/3578356.3592575>

(*EuroMLSys '23*), May 8, 2023, Rome, Italy. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/3578356.3592575>

1 INTRODUCTION

Federated NLP The development of pre-trained models is overwhelming with the rise of BERT [1]. Their deployment [2–6] is commonly composed of two-step training: pre-training and fine-tuning. Unlike self-supervised pre-training, fine-tuning is supervised, requiring task-specific tremendous labeled data. However, the exploitation of private user data is restricted and even prohibited in some cases by several data protection regulations such as GDPR [7] and CCPA [8]. Recently, federate learning (FL) [9, 10] becomes the de-facto approach to train a model with privacy preserved. As such, federated NLP (FedNLP) [11, 12] is now an important topic towards practical NLP applications.

Problem and challenge A key obstacle to practical FedNLP is data labeling. It's much more difficult to label data on client devices than on centrally collected data [13, 14]. Lack of sufficient labeled data severely limits the practicality and scalability of FedNLP in real-world NLP applications. Therefore, it is important to address the issue of few-shot or even zero-shot FedNLP tasks. There are very few efforts on this topic [15–17], which still assume a fairly large number (typically >1000 in total) of labels that are uniformly distributed across clients. However, in practice, the labeled data distribution could be skewed across clients, and such skewness would result in a significant drop in the accuracy according to our experiments in §2.

Our solution and contribution

(1) To tackle the issue of insufficient and skewness of labeled data, we design a comprehensive data generator as the first step towards simulating the distribution of labeled data for few-shot FedNLP tasks. The generator has two meta-parameters: data quantity and skewness, which encompass most, if not all, potential scenarios for practical few-shot FedNLP.

(2) To boost the performance of few-shot FedNLP, we design a data-augmented prompt system, namely AUG-FedPrompt. AUG-FedPrompt orchestrates prompt learning [18] and pseudo labeling [19]. Prompt learning introduces a task description in NLP training. It helps task-specific fine-tuning achieve high accuracy with very few labeled data samples in FedNLP.

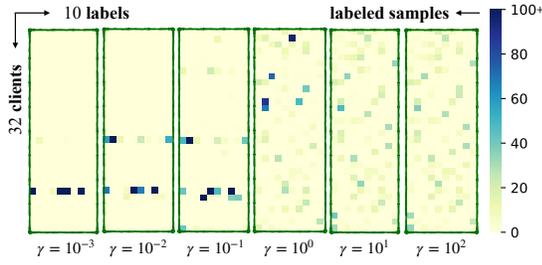


Figure 1: Visualizing the skewness of labeled data on YAHOO [20] with $n=1024$, $\xi=32$, γ being 10^n , $n=-3,-2,\dots,2$. Each sub-figure is a 32×10 matrix, where 32 is the number of clients and 10 is the number of labels. The intensity of each cell represents the number of labeled samples for a specific label in the client-side local data.

Furthermore, to tackle performance degradation caused by skewed label distribution, AUG-FedPrompt leverages enormous and easily accessible unlabeled data for pseudo labeling-based data augmentation.

(3) Our extensive experiments on four English datasets demonstrate that AUG-FedPrompt can achieve a substantial performance gain (25%–55% higher accuracy) over the state-of-the-art FedNLP approaches under various few-shot settings. Augmentation with unlabeled data enhances AUG-FedPrompt to perform well with highly skewed labeled distribution across clients. Overall, AUG-FedPrompt can achieve a comparable performance with the state-of-the-art FedNLP approaches with less than 0.1% labeled data.

2 PROBLEM SETUP

Federated NLP Training Procedure The two NLP training phases, i.e., pre-training and fine-tuning, require data of disparate natures. Pre-training is typically done on public text corpora such as Wikipedia articles, while fine-tuning requires domain-specific samples, such as user reviews, messages, or emails. For mobile computing, domain-specific samples are gathered from end-users and distributed across mobile devices, while ensuring the protection of privacy. To fine-tune models on such private, distributed data, federated learning is the de-facto approach [11, 12]. Prior to training, a cloud service distributes a pre-trained model to all client devices. In a training session targeting a specific NLP task and domain, a cloud service selects multiple mobile devices to participate in training. Each device trains a local copy of the model with its private data and sends the model updates to the cloud. Upon aggregating the model updates from multiple devices, the cloud sends an updated model to the devices. This training procedure is repeated until the model converges.

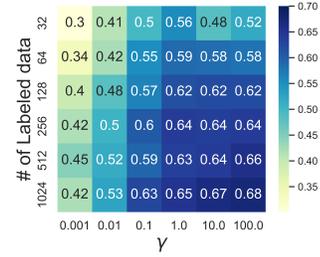


Figure 2: Average accuracy of federated few-shot learning under different data quantity and skewness. When skewness γ grows larger, labeled data will be more uniformly distributed, and vice versa. Dataset: YAHOO [20].

Federated few-shot data generator Apart from data privacy, lack of sufficient labeled data is a crucial issue and an inherent feature in mobile scenarios. Alike data feature could be non-independent and identically distributed (non-iid), the scarce labels is not always uniformly distributed in the real world. Based on the definition of non-iid partition strategies [11, 21], we further define the quantity and skewness of labels under federated few-shot learning scenario.

We define a new tuple (n, γ) to represent the practical few-shot learning training data distribution, where n represents the total numbers of labeled data, γ represents the skewness of labeled data.

The quantity of labeled data assigned to each client follows a Dirichlet allocation $z \sim \text{Dir}_\xi(\gamma)$, where ξ is the number of clients with labeled data¹. We can then allocate labeled data from the global labeled dataset to selected clients based on the distribution z , with client _{i} being assigned a labeled dataset of size $|\mathcal{T}_i| = z_i n$. For example, in Figure 1, we visualize the labeled data skewness on Yahoo [20] with $n=1024$, $\xi=32$, γ being 10^n , $n=-3,-2,\dots,2$. Each sub-figure is a 32×10 matrix, the intensity of which represents the labeled samples of a particular label. When γ is small (10^{-3} , 10^{-2} , 10^{-1}), the labeled data will be skewed distributed, i.e., only few clients own labeled data; when $\gamma=10^2$, labeled data is nearly uniformly distributed on all clients.

Performance degradation under skewed distribution

In Figure 2, we present the impact of label skewness on federated few-shot learning. We observe that as γ decreases, i.e., the labeled data becomes more skewed, the convergence performance of the model degrades. For example, when labeled data points are 1024, uniform distribution ($\gamma=100$) will be 26% better than skewed distribution ($\gamma=0.001$). The rationale behind this phenomenon is that under common non-iid data distribution, individual clients tends to possess more specific

¹ ξ could be an optional hyper-parameter to strict the maximum of clients owing labeled data. In this manuscript, we fix ξ as 32 for simplicity.

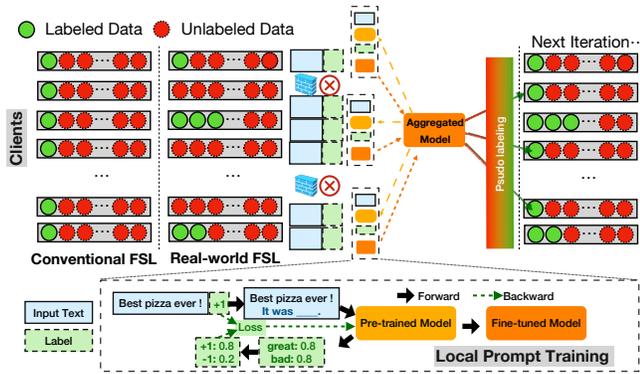


Figure 3: Workflow of AUG-FedPrompt.

data features. The labels concentrated on certain clients results in a skewed feature distribution of training data. This bias can lead to unfairness, as the aggregated model may favor certain labels over others, resulting in a significant drop in convergence accuracy. We provide a more detailed analysis of this phenomenon in Section 4.3.

3 SYSTEM DESIGN

We propose AUG-FedPrompt as a solution to address the challenges posed by data privacy concerns and label scarcity. AUG-FedPrompt leverages large amounts of unlabeled data via the federated orchestration of pseudo labeling and prompt learning. AUG-FedPrompt fine-tunes the pre-trained model through prompt learning on client devices. After local training and federated aggregation, AUG-FedPrompt makes inference on unlabeled training data, from which high-confident results, i.e., pseudo labels [19, 22, 23] are selected for subsequent training.

We describe the training workflow of AUG-FedPrompt in Figure 3. A public pre-trained transformer-based language model M is transferred to chosen clients. We assume that each client has access to a tiny training set \mathcal{T} (typically < 10) and a much larger set of unlabeled examples \mathcal{D} (typically > 1000).

For local prompt training, we annotate T as the vocabulary of model M , $_ \in T$ as the mask token and T^* as the set of all token sequences. To clarify, T is composed of tokens representing labels description and T^* is composed of tokens representing input text, which is a larger corpus. The sequence of input phrases is $\mathbf{x} = (s_1, \dots, s_k)$ where $s_i \in T^*$. The pattern-verbalizer pair \mathbf{p} includes: 1) a *pattern* $P : X \rightarrow T^*$ maps inputs x to a cloze question containing a single mask; 2) a *verbalizer* $v : Y \rightarrow T$ maps each output y to a single token representing its task-specific meaning in the pattern.

The purpose of local prompt training is to derive the probability that y is the correct output of x from the probability

that $v(y)$ is the most likely token at the masked position in $P(x)$. Based on this rationale, we define the conditional probability distribution s_p of y given x as:

$$s_p(y | x) = \frac{\exp q_p(y | x)}{\sum_{y' \in Y} \exp q_p(y' | x)} \quad (1)$$

where $q_p(y | x) = M(v(y) | P(x))$ is the probability that M assigns to $v(y)$ in sequence $P(x)$.

For client-side fine-tuning, the pre-trained model M is fine-tuned on local labeled data (x, y) by minimizing the cross-entropy between $s_p(y | x)$ and y . For server-side aggregation, in each iteration i , client k sends its updated model M_k^i to the cloud for aggregation using FedAVG algorithm [9]; the aggregated model is denoted as M^i .

For data augmentation, M^i is distributed to clients with large amount of unlabeled data for pseudo labeling. Each unlabeled example $\hat{x} \in D$ is labeled with pseudo label \hat{y} based on $s_p(\hat{y} | \hat{x})$. The pseudo-labeled dataset then is utilized for fine-tuning the client-side model in the subsequent iteration.

The resulting pseudo-labeled dataset could consist of enormous samples with wrong labels. Directly involving them in the next training iteration will poison the foundation model, which makes it could be even worse than purely using the limited labeled data. To address this issue, we propose two techniques to filter out those wrong samples and remain the purity of augment dataset: 1) Filtering by model capacity: we eliminate those models with low model capacity, i.e., those that perform poorly on validation datasets. 2) Filtering by confidence: we remove samples with low confidence, i.e., those with a probability of the most likely label lower than a threshold. Both capacity and confidence are hyper-parameters that can be tuned flexibly depending on particular tasks or datasets.

4 PRELIMINARY EXPERIMENTS

In this section, we evaluate the performance of AUG-FedPrompt across data scales. AUG-FedPrompt significantly outperforms naive federated fine-tuning. It could perform on par with full-set training while saving up to 99.9% labeled data. Apart from data efficiency, AUG-FedPrompt shows great robustness under various practical few-shot scenario regardless of skewed or uniform label distribution.

4.1 Experiment Setup

Dataset and models We perform our evaluation on four English datasets and manually designed prompts², detailed information is shown in Table 1. (1) AGNEWS [20] is a news classification dataset. Given headline and text body, news

²We try 6, 2, 6, 4 different prompts for each datasets separately and report the chosen one that performs best. The verbalizers are the same as previous literature [18].

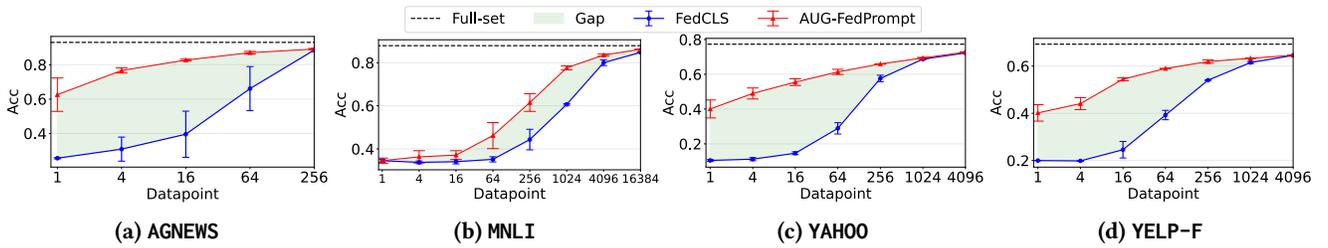


Figure 4: Average accuracy and standard deviation for AUG-FedPrompt across data scales. FedCLS stands for the vanilla federated fine-tuning. Full-set stands for fine-tuning on the full labeled data.

Dataset	Prompt	Train	Test
AGNEWS [20]	a (____) b	120,000	7,600
MNLI [24]	“a” ? ____, “b”	392,702	9,815
YAHOO [20]	[Category:] a ____ b	1,400,000	60,000
YELP-F [20]	It was ____. a	650,000	50,000

Table 1: Evaluation datasets. Each dataset is distributed to 1000 clients. Label quantity of each class follows the non-iid label distribution in [11] where $\alpha = 1$.

Dataset		AGNEWS	MNLI	YAHOO	YELP-F
Uniform	FedCLS	66.1±12.8	60.1±0.4	57.6±1.9	54.0±0.1
	FedPrompt	87.0±0.8	77.6±0.8	66.0±0.1	61.9±0.7
Skewed	FedCLS	64.8±3.1	37.7±5.6	24.4±10.3	38.3±8.8
	FedPrompt	68.4±2.4	42.4±5.8	41.8±4.3	51.2±1.8
	w/ augment	90.2±0.5	75.7±1.2	66.9±1.1	58.2±2.4

Table 2: AUG-FedPrompt enhances performance under different few-shot learning settings. FedPrompt stands for AUG-FedPrompt without unlabeled data augmentation. Datapoint: 64 for AGNEWS, 1024 for MNLI, 256 for YHAOO and YELP-F.

needs to be classified as one of the four categories: World, Sports, Business or Science/Tech. (2) MNLI [24] is a sentence understanding dataset. Given text pairs $x = (a, b)$, the task is to find out whether a implies b , a and b contradict each other or neither. (3) YELP Review Full (YELP-F) [20] is a restaurant rating dataset. Given a customer’s review, text should be estimated on a 1-5 star scale. (4) YAHOO [20] is a text classification dataset. Given a question and an answer, one of ten possible categories needs to be assigned. All experiments are conducted using the same pre-trained model, RoBERTa-large (355M parameters) [25], which we load from the transformers [26] library.

Hyper-parameters In line with previous observations [27], few-shot fine-tuning performance varies across chosen labeled data considerably. We run every experiment 3 times in order to reduce variance. Unless otherwise stated, we use the recommended set of hyper-parameters from previous work [18]: mini-batch size as 4; local training iteration as 1; learning rate as 10^{-5} ; max sequence length as 256. For pseudo labeling, we filter out those aggregated models performing worse than the zero-shot model and those pseudo-labeled data with confidence lower than 0.9. For the FL configurations at the server side, we follow the prior FedNLP literature [11, 12] to select 5 participants for each training round by default. The fine-tuned models will be collected in the central server and aggregated through *FedAvg* algorithm [9].

4.2 Performance across Data Scales

AUG-FedPrompt enjoys a substantial advantage on each task. As shown in Figure 4, we compare our AUG-FedPrompt performance with FedCLS, i.e., the vanilla federated fine-tuning where a generic classifier layer inserted after pre-trained models is fine-tuned. Highlighted region shows the accuracy gap between AUG-FedPrompt and FedCLS. There are up to 50%, 25%, 55%, 38% accuracy improvement separately for 4 datasets. Both approaches improve with more labeled data, but AUG-FedPrompt remains better by a varying amount. AUG-FedPrompt reaches 99% relative performance of full-set with 90% less training data compared to full-set federated training. AUG-FedPrompt shows a strong zero-shot inferring capability, i.e., without task-specific fine-tuning, except for MNLI dataset. MNLI dataset may need more labeled data to make full use of the prompt to the pre-trained models. For a usable accuracy, i.e., 90% relative performance of full-set training accuracy, AUG-FedPrompt only needs 64, 256, 256 in total for AGNEWS, YAHOO and YELP-F, saving up to 99.9% training data compared to full-set federated fine-tuning. Please note that 64 is the total number of labels across all clients, not per client.

4.3 Impact of Data Augmentation

AUG-FedPrompt enhances FedPrompt performance when labeled data is skewed distributed. As shown in Table 2, FedPrompt, i.e., AUG-FedPrompt without data augment shows

Challenges	Possible Solutions
Huge training latency	Model structure optimization [28, 29].
Large memory requirement	Rematerialization [30, 31], paging [32].
Excessive inference for pseudo labeling	Pacing [23, 33], early-exit [34, 35].
High communication cost	Quantization [36, 37], sparsity [38, 39].

Table 3: Challenges and possible solutions.

competitive performance when labeled data is uniformly distributed on clients. While skewed distribution of labeled data will hurt FedPrompt performance significantly. For example, FedPrompt performance degrades to 41.8% on YHAAO when 256 labeled data is skewed distributed on 32 clients. Considering that skewed distribution is common in real-world, we integrate AUG-FedPrompt with data augmentation to mitigate the performance degradation.

It is important to recall that prompts learning introduces a task description in NLP training. Prompt helps task-specific fine-tuning perform well even with few labeled training data. This rationale paves the way for the efficiency of pseudo labeling; it helps to label more data correctly at the early stage of training. Together with our confidence filter for pseudo-labeling, AUG-FedPrompt makes pseudo-labeled data seldom hurt. For example, we annotate 100 unlabeled data on each client involved in per round for AGNEWS. In the first three rounds, the average ratio of correctly labeled data by pseudo-labeling on unlabeled data is 92.5%. The inference accuracy will further increase along with the FL training moves on, reaching 95.3% at the convergence round. Those ‘nail’ data, about 5 out of 100 in total, is hard to be correctly annotated and filtered out. Fortunately, we observe that they do not affect the model convergence as shown in Table 2. After pseudo labeling, AUG-FedPrompt performs on par with full-set fine-tuning and greatly outperforms vanilla few-shot fine-tuning, reaching a usable accuracy with scarce labeled data.

5 SYSTEM COST

There is no free lunch for the performance improvement of AUG-FedPrompt. The orchestrating of pseudo labeling and prompt learning results in promising few-shot performance, but it also incurs a non-trivial system cost. In this section, we discuss the necessity of large models for AUG-FedPrompt, as well as the associated system cost. Challenges and possible solutions are concluded in Table 3.

To begin with, we conduct experiments to evaluate the performance of AUG-FedPrompt on various foundation models. As demonstrated in Figure 5, RoBERTa-large outperforms all other models across all four datasets, particularly MNLI and YELP-F, where it shows a significant improvement (up to 38.2%). In contrast, BERT-large, despite having similar parameters to RoBERTa-large, performed poorly. Interestingly,

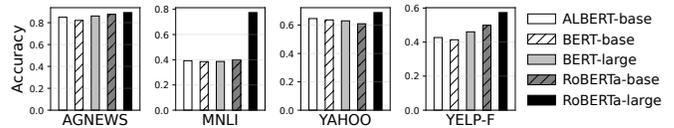


Figure 5: AUG-FedPrompt convergence performance with different models and datasets. 0.1% labeled data uniformly distributed in 32 clients.

Model	ALBERT-base [29]	BERT-base [1]	BERT-large [1]	RoBERTa-base [25]	RoBERTa-large [25]
Memory (GB)	3.7	5.4	OOM (9.8)	5.8	OOM (10.4)
Latency (s)	1.4	1.9	~7.8	2.1	~8.1
Param. (M)	11.7	109.5	334.9	124.6	355.3

Table 4: System cost of different NLP models. Tested on NVIDIA TX2. Batch size: 4.

certain small models, e.g. ALBERT-base [29], which is optimized from BERT-base achieved superior results compared to the standard BERT-base model, despite containing only 10.7% of the parameters. These findings suggest that large models can help augment the few-shot learning abilities of AUG-FedPrompt, and that model structure optimization shows promise in making AUG-FedPrompt a more practical solution.

The excellent performance of RoBERTa-large aligns with previous research [18, 27, 40], highlighting the need for large-scale foundational models to fully leverage prompt learning. However, despite its merits, the model’s high memory usage and latency cannot be overlooked. As shown in Table 4, even on a powerful GPU-embedded edge device like NVIDIA TX2 [41], training RoBERTa-large leads to long latency (about 8.1s per batch). Moreover, during training, our testbed device, which has only 8GB of RAM, ran out of memory during training. Because the peak memory usage of RoBERTa-large fine-tuning is over 10GB³.

Apart from local prompt training, a mobile client need to perform inference on *all* of its unlabeled data to generate pseudo labels. However, most of this inference is ultimately unnecessary, as only a small fraction (the most confident) of pseudo labels will be selected for subsequent training. As a result, the inference process dominates the total delay due to the large volume of unlabeled data that needs to be processed. According to our measurements, this process accounts for up to 87.4% of the total computation time. Keeping a balanced pace between training and labeling is crucial to reduce those redundant inference.

In addition, it should be noted that the overall resource cost of AUG-FedPrompt system should be extremely higher, let alone long heavy-duty computing. The reason for this is the

³Tested on a central server.

need to transfer the entire model, which can be several GBs in size, in a federated learning scenario. As the size of the model increases, so too does the amount of data that needs to be transferred, leading to higher communication costs. This can be particularly problematic in settings with limited network bandwidth, such as mobile devices, where large network traffic can significantly impact system performance [36, 42–44].

6 CONCLUSIONS AND FUTURE WORK

This manuscript explores a crucial but less explored issue: data labels can be scarce in federated learning. We provide a comprehensive definition of a data generator for federated few-shot learning tasks and demonstrate that the lack and skewness of labeled data can significantly degrade federated learning convergence performance. To mitigate this issue, we propose AUG-FedPrompt, a novel federated few-shot learning system that orchestrates prompt learning and pseudo labeling. AUG-FedPrompt shows competitive performance under various federated few-shot learning settings, requiring less than 0.1% data to be manually labeled.

In conclusion, our experiments have demonstrated the impressive few-shot performance of AUG-FedPrompt when used with large-scale pre-trained models. However, fine-tuning these ‘behemoths’ can be extremely resource-intensive, requiring significant computational power and memory. Additionally, the communication of large model parameters can consume a considerable amount of bandwidth. Future work will focus on the development of an optimized system solution for AUG-FedPrompt to enhance its resource efficiency.

ACKNOWLEDGMENTS

This research was supported by National Key Research and Development Program of China #2020YFB1805500, the Fundamental Research Funds for the Central Universities, and NSFC #62032003, #61922017, #61921003. Mengwei Xu was partly supported by NSFC #62102045, Beijing Nova Program #Z211100002121118, and Young Elite Scientists Sponsorship Program by CAST #2021QNRC001. The authors thank the anonymous reviewers for their insightful feedback.

REFERENCES

- [1] 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.
- [2] Lei Zhang, Shuai Wang, and Bing Liu, “Deep learning for sentiment analysis: A survey,” *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 8, no. 4, pp. e1253, 2018.
- [3] Taihua Shao, Yupu Guo, Honghui Chen, and Zepeng Hao, “Transformer-based neural network for answer selection in question answering,” *IEEE Access*, vol. 7, pp. 26146–26156, 2019.
- [4] Betty Van Aken, Benjamin Winter, Alexander Löser, and Felix A Gers, “How does bert answer questions? a layer-wise analysis of transformer representations,” in *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, 2019, pp. 1823–1832.
- [5] Seohyun Kim, Jinman Zhao, Yuchi Tian, and Satish Chandra, “Code prediction by feeding trees to transformers,” in *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. IEEE, 2021, pp. 150–162.
- [6] Alexey Svyatkovskiy, Shao Kun Deng, Shengyu Fu, and Neel Sundaresan, “Intellicode compose: Code generation using transformer,” in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2020, pp. 1433–1443.
- [7] Paul Voigt and Axel Von dem Bussche, “The eu general data protection regulation (gdpr),” *A Practical Guide, 1st Ed., Cham: Springer International Publishing*, vol. 10, no. 3152676, pp. 10–5555, 2017.
- [8] Stuart L Pardo, “The california consumer privacy act: Towards a european-style privacy regime in the united states,” *J. Tech. L. & Pol’y*, vol. 23, pp. 68, 2018.
- [9] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguerre y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Artificial Intelligence and Statistics*. PMLR, 2017, pp. 1273–1282.
- [10] Qiang Yang, Yang Liu, Yong Cheng, Yan Kang, Tianjian Chen, and Han Yu, “Federated learning,” *Synthesis Lectures on Artificial Intelligence and Machine Learning*, vol. 13, no. 3, pp. 1–207, 2019.
- [11] Bill Yuchen Lin, Chaoyang He, Zihang Zeng, Hulin Wang, Yufen Huang, Christophe Dupuy, Rahul Gupta, Mahdi Soltanolkotabi, Xiang Ren, and Salman Avestimehr, “Fednlp: Benchmarking federated learning methods for natural language processing tasks,” *Findings of NAACL*, 2022.
- [12] Dongqi Cai, Yaozong Wu, Shangguang Wang, Felix Xiaozhu Lin, and Mengwei Xu, “Autofednlp: An efficient fednlp framework,” *arXiv preprint arXiv:2205.10162*, 2022.
- [13] Huatao Xu, Pengfei Zhou, Rui Tan, Mo Li, and Guobin Shen, “Limubert: Unleashing the potential of unlabeled data for imu sensing applications,” in *Proceedings of the 19th ACM Conference on Embedded Networked Sensor Systems*, 2021, pp. 220–233.
- [14] Guoliang Li, Yudian Zheng, Ju Fan, Jiannan Wang, and Reynold Cheng, “Crowdsourced data management: Overview and challenges,” in *Proceedings of the 2017 ACM International Conference on Management of Data*, 2017, pp. 1711–1716.
- [15] Chenyou Fan and Jianwei Huang, “Federated few-shot learning with adversarial learning,” in *2021 19th International Symposium on Modeling and Optimization in Mobile, Ad hoc, and Wireless Networks (WiOpt)*. IEEE, 2021, pp. 1–8.
- [16] Fei Chen, Mi Luo, Zhenhua Dong, Zhenguo Li, and Xiuqiang He, “Federated meta-learning with fast convergence and efficient communication,” *arXiv preprint arXiv:1802.07876*, 2018.
- [17] Tony Huang, Jack Chu, and Fangyun Wei, “Unsupervised prompt learning for vision-language models,” *arXiv preprint arXiv:2204.03649*, 2022.
- [18] Timo Schick and Hinrich Schütze, “Exploiting cloze questions for few shot text classification and natural language inference,” *arXiv preprint arXiv:2001.07676*, 2020.
- [19] Dong-Hyun Lee et al., “Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks,” in *Workshop on challenges in representation learning, ICML*, 2013, vol. 3, p. 896.
- [20] Xiang Zhang, Junbo Zhao, and Yann LeCun, “Character-level convolutional networks for text classification,” *Advances in neural information processing systems*, vol. 28, 2015.
- [21] Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He, “Federated learning on non-iid data silos: An experimental study,” in *2022 IEEE*

- 38th International Conference on Data Engineering (ICDE). IEEE, 2022, pp. 965–978.
- [22] Eric Arazo, Diego Ortego, Paul Albert, Noel E O'Connor, and Kevin McGuinness, "Pseudo-labeling and confirmation bias in deep semi-supervised learning," in *2020 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2020, pp. 1–8.
- [23] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston, "Curriculum learning," in *Proceedings of the 26th annual international conference on machine learning*, 2009, pp. 41–48.
- [24] Adina Williams, Nikita Nangia, and Samuel R Bowman, "A broad-coverage challenge corpus for sentence understanding through inference," *arXiv preprint arXiv:1704.05426*, 2017.
- [25] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov, "Roberta: A robustly optimized bert pretraining approach," *arXiv preprint arXiv:1907.11692*, 2019.
- [26] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al., "Huggingface's transformers: State-of-the-art natural language processing," *arXiv preprint arXiv:1910.03771*, 2019.
- [27] Teven Le Scao and Alexander M Rush, "How many data points is a prompt worth?," *arXiv preprint arXiv:2103.08493*, 2021.
- [28] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf, "Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter," *arXiv preprint arXiv:1910.01108*, 2019.
- [29] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut, "Albert: A lite bert for self-supervised learning of language representations," *arXiv preprint arXiv:1909.11942*, 2019.
- [30] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin, "Training deep nets with sublinear memory cost," *arXiv preprint arXiv:1604.06174*, 2016.
- [31] Qipeng Wang, Mengwei Xu, Chao Jin, Xinran Dong, Jinliang Yuan, Xin Jin, Gang Huang, Yunxin Liu, and Xuanzhe Liu, "Melon: Breaking the memory wall for resource-efficient on-device machine learning," 2022.
- [32] Xuan Peng, Xuanhua Shi, Hulin Dai, Hai Jin, Weiliang Ma, Qian Xiong, Fan Yang, and Xuehai Qian, "Capuchin: Tensor-based gpu memory management for deep learning," in *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, 2020, pp. 891–905.
- [33] Paola Cascante-Bonilla, Fuwen Tan, Yanjun Qi, and Vicente Ordonez, "Curriculum labeling: Revisiting pseudo-labeling for semi-supervised learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021, vol. 35, pp. 6912–6920.
- [34] Wangchunshu Zhou, Canwen Xu, Tao Ge, Julian McAuley, Ke Xu, and Furu Wei, "Bert loses patience: Fast and robust inference with early exit," *Advances in Neural Information Processing Systems*, vol. 33, pp. 18330–18341, 2020.
- [35] Stefanos Laskaridis, Alexandros Kouris, and Nicholas D Lane, "Adaptive inference through early-exit networks: Design, challenges and directions," in *Proceedings of the 5th International Workshop on Embedded and Mobile Deep Learning*, 2021, pp. 1–6.
- [36] Jiaxiang Wu, Weidong Huang, Junzhou Huang, and Tong Zhang, "Error compensated quantized sgd and its applications to large-scale distributed optimization," in *International Conference on Machine Learning*. PMLR, 2018, pp. 5325–5333.
- [37] Ahmed M Abdelmoniem and Marco Canini, "Towards mitigating device heterogeneity in federated learning via adaptive model quantization," in *Proceedings of the 1st Workshop on Machine Learning and Systems*, 2021, pp. 96–103.
- [38] Ang Li, Jingwei Sun, Pengcheng Li, Yu Pu, Hai Li, and Yiran Chen, "Hermes: an efficient federated learning framework for heterogeneous mobile clients," in *Proceedings of the 27th Annual International Conference on Mobile Computing and Networking*, 2021, pp. 420–437.
- [39] Jonathan Frankle and Michael Carbin, "The lottery ticket hypothesis: Finding sparse, trainable neural networks," *arXiv preprint arXiv:1803.03635*, 2018.
- [40] Timo Schick and Hinrich Schütze, "True few-shot learning with prompts—a real-world perspective," *Transactions of the Association for Computational Linguistics*, vol. 10, pp. 716–731, 2022.
- [41] NVIDIA JETSON TX2, "256-core nvidia pascal gpu," <https://developer.nvidia.com/embedded/jetson-tx2>.
- [42] Jie Xu and Heqiang Wang, "Client selection and bandwidth allocation in wireless federated learning networks: A long-term perspective," *IEEE Transactions on Wireless Communications*, vol. 20, no. 2, pp. 1188–1200, 2020.
- [43] Amirhossein Reiszadeh, Isidoros Tziotis, Hamed Hassani, Aryan Mokhtari, and Ramtin Pedarsani, "Straggler-resilient federated learning: Leveraging the interplay between statistical accuracy and system heterogeneity," *arXiv preprint arXiv:2012.14453*, 2020.
- [44] Su Wang, Mengyuan Lee, Seyyedali Hosseinalipour, Roberto Mora-bitto, Mung Chiang, and Christopher G Brinton, "Device sampling for heterogeneous federated learning: Theory, algorithms, and implementation," in *IEEE INFOCOM 2021-IEEE Conference on Computer Communications*. IEEE, 2021, pp. 1–10.