

FedPrompt: Communication-Efficient and Privacy-Preserving Prompt Tuning in Federated Learning

1st Haodong Zhao
Shanghai Jiao Tong University
Shanghai, China
zhaohaodong@sjtu.edu.cn

2nd Wei Du
Shanghai Jiao Tong University
Shanghai, China
dddddw@sjtu.edu.cn

3rd Fangqi Li
Shanghai Jiao Tong University
Shanghai, China
solour_lfq@sjtu.edu.cn

4th Peixuan Li
Shanghai Jiao Tong University
Shanghai, China
peixuan.li@sjtu.edu.cn

5th Gongshen Liu*
Shanghai Jiao Tong University
Shanghai, China
lgshen@sjtu.edu.cn

Abstract—Federated learning (FL) has enabled global model training on decentralized data in a privacy-preserving way by aggregating model updates. However, for many natural language processing (NLP) tasks that utilize pre-trained language models (PLMs) with large numbers of parameters, there are considerable communication costs associated with FL. Recently, prompt tuning, which tunes some soft prompts without modifying PLMs, has achieved excellent performance as a new learning paradigm. Therefore we want to combine the two methods and explore the effect of prompt tuning under FL. In this paper, we propose "FedPrompt" to study prompt tuning in a model split aggregation way using FL, and prove that split aggregation greatly reduces the communication cost, only 0.01% of the PLMs' parameters, with little decrease on accuracy both on IID and Non-IID data distribution. This improves the efficiency of FL method while also protecting the data privacy in prompt tuning. In addition, like PLMs, prompts are uploaded and downloaded between public platforms and personal users, so we try to figure out whether there is still a backdoor threat using only soft prompts in FL scenarios. We further conduct backdoor attacks by data poisoning on FedPrompt. Our experiments show that normal backdoor attack can not achieve a high attack success rate, proving the robustness of FedPrompt. We hope this work can promote the application of prompt in FL and raise the awareness of the possible security threats.

Index Terms—FL, prompt tuning, PLM, split aggregation

I. INTRODUCTION

Pre-trained language models [1]–[3] are widely used in many NLP tasks by the fine-tuning paradigm. However, fine-tuning a PLM with a large number of parameters would be memory-consuming. The reason is that the gradients and optimizer states of all parameters need to be stored. Also the lack of labeled data in fine-tuning phase, as well as few-shot problem, limits the use of this paradigm. When using the pre-training and fine-tuning paradigm in federated learning, the

* corresponding author. This research work has been sponsored by the Joint Funds of the National Natural Science Foundation of China (Grant No.U21B2020) and Ant Group.

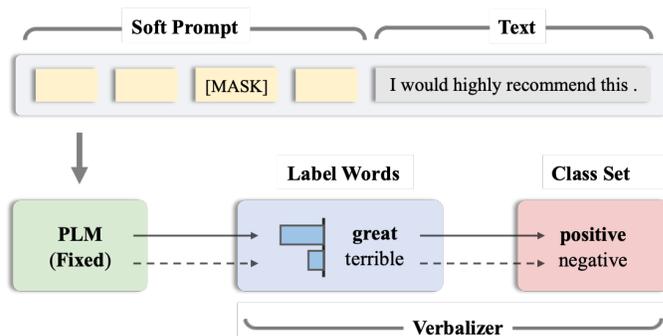


Fig. 1. The example of prompt tuning, which consists of soft prompt, text, PLM and verbalizer.

communication cost is even higher as all parameters of the PLMs provided by each participant need to be aggregated in each round of the training process. Therefore, it is very important and urgent to find ways to improve the efficiency of pre-trained models in federated learning. Recently, prompt tuning [4] has achieved excellent results as a learning paradigm for adapting fixed PLMs to different downstream tasks. As shown in Fig. 1, soft prompt as well as [MASK] token are added to the text as inputs to the model. Among them, soft prompt is used as trainable parameters to be adapted to downstream tasks, [MASK] token is used to predict the label word for the downstream task, and verbalizer is used to map the label word to the real label. All downstream tasks can be uniformly transformed into the form of pre-training tasks of PLMs. Thus, using a fixed PLM and different soft prompts can be applied to different downstream tasks. Also, freezing the parameters of PLM and only tuning the soft prompt significantly reduces the number of training parameters.

Nowadays, with mobile devices becoming the primary computing devices for many people, a huge amount of data is

generated and distributed on a wide range of devices. It is an important opportunity and challenge to make full use of these devices and data. Although deep learning has made a lot of progress in many scenarios [5], a data center is required to collect data for training in most cases. Models trained on such data have stronger usability in many intelligent applications, but exchanging and storing sensitive data in a data center carries risks and responsibilities [6]. Previous distributed deep learning methods [7], [8] propose solutions to big data and huge models. However, the computation and communication cost of traditional distributed learning are unacceptable for participants [9], [10].

FL [11], [12] is an attractive learning method that aims to train a global model over decentralized data while preserving data privacy. In federated learning a subset of clients download a local copy of global model, and compute local model gradients with their local private data in each round. A central server coordinates the distributed clients and only aggregates local model parameters to update the global model, collaborating isolated data islands without raw data exchanging. FL takes multiple rounds of local and global procedures until convergence. Advanced privacy protection methods, e.g., differential privacy (DP), can be further applied for stricter privacy protection. According to above illustration, it seems that using PLM-empowered methods in FL will achieve remarkable performance. For example, [13] propose a news recommendation method to train models using FL. However, when using FL, the model sizes of many existing news recommendation methods are too large to communicate between participants. For example, the base version of BERT [1] models have more than 110M parameters.

In this paper, we modify prompt tuning in a model split aggregation way using FL, named "FedPrompt", where no prior work has been done. First, unlike simple FL methods tune and aggregate full model parameters, FedPrompt only tunes and aggregates some soft prompts for corresponding downstream tasks in FL, and freezes PLMs to decrease communication cost. Second, we are interested in the security of FedPrompt. Prompts are uploaded and downloaded between public platforms and personal users like PLMs, and backdoor attacks are difficult to find for users. We try to get a poisoned global prompt then when the PLMs are loaded with the poisoned prompt, the model will be implanted with the backdoor.

Experiments carried on various NLP tasks, such as sentiment analysis and sentence-pair classification prove that FedPrompt reduces the communication cost greatly with little decrease on accuracy. Further experiments on backdoor attack show that only by normal methods that poisoning part of training data, the poisoned prompt can not establish a shortcut between the specific trigger word and the target label word. Compared to the method of aggregating and tuning all parameters of huge model, FedPrompt is much more communication-efficient. And compared to prompt tuning without using FL, FedPrompt outperforms in privacy preservation. We also consider other prompt types for a better performance, and using local differential privacy (LDP) to guarantee the privacy. Our

contributions are summarized as follows:

- We propose FedPrompt, the new prompt tuning method using FL, freezes PLMs and only aggregates and tunes some soft prompts to decrease communication cost.
- We conduct extensive experiments on NLP tasks to measure the performance of FedPrompt. Experiments show that FedPrompt can reduce the communication cost greatly with little decrease on accuracy.
- We further test the model robustness to backdoor attack, and experiment on different hyper-parameter settings, prompt types and LDP to promote the better performance of FedPrompt.

II. RELATED WORK

Prompt tuning first appeared in WARP proposed by [14], after which this method of adding continuous trainable vectors to the input began to be widely studied. Prefix-Tuning [15] adds soft prompt to each layer of the transformer model and applies it to natural language generation (NLG) tasks. P-tuning [16] proposes that some task-related hard prompts can be used as anchors while using soft prompt. Prompt tuning [4] explores the effect of soft prompt on domain adaptation and different model scales. They found that the larger the scale of PLMs, the better the effect of prompt tuning. Recently, P-tuningV2 [17] more finely designs prompt tuning on the basis of the above research. They use a deep soft prompt similar to Prefix-Tuning, and change the verbalizer to a linear classification head, which means that they no longer use the way of mask language model (MLM) to get predictions. They also try to apply prompt tuning to difficult natural language understanding (NLU) tasks (i.e., sequence tagging), such as name entity recognition and semantic role labeling. Moreover, PTR [18] applies logic rules to build templates that are more suitable for text classification tasks, and PPT [19] pre-trains the soft prompt on multiple different tasks to get a better prompt tuning for the downstream tasks.

FL [11] is a distributed machine learning method that aggregates global model by each local model sharing its parameters (gradients) with the central server after every round of local training on its local data. Proposed FedAvg [11] enables clients to collaboratively train global model without sharing their original data. Various extensions of FedAvg [20]–[24] have been proposed to obtain better performance in communication and deal with heterogeneity. To reduce computation and communication, [25] propose a framework decomposing big recommendation model into a large news model only in server and shared user model. However, though above methods do not share local data directly, naive parameters (gradients) sharing method could lead to privacy leakage of clients [26], [27]. Consequently, several methods are proposed to protect privacy, including differential privacy (DP) [28]–[30] and secure multi-party computation (MPC) [31], [32]. In addition to privacy leakage, many works propose mechanisms to poison FL models in training phase [33], [34] and evasion attacks in inference or testing phase [35], [36].

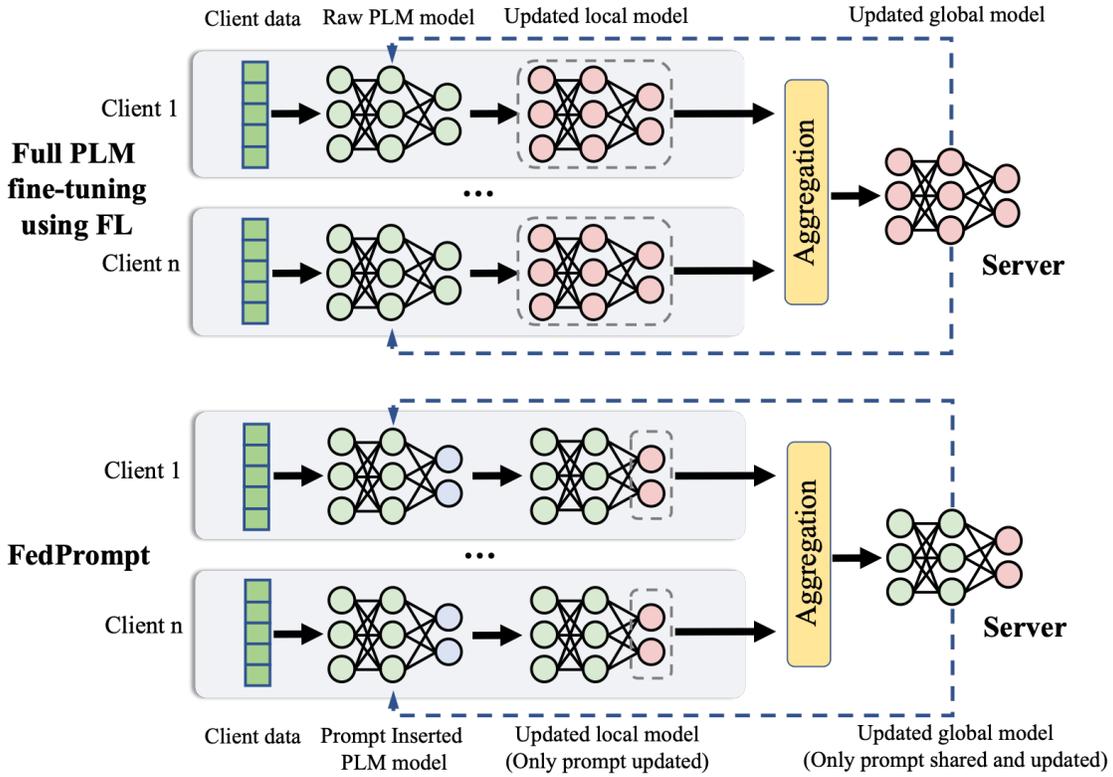


Fig. 2. Structure of FedPrompt and full PLM fine-tuning using FL. The above one is full PLM fine-tuning using FL, all of the parameters (framed pink nodes) need to be updated. The bottom one is FedPrompt, only soft prompt parameters (framed pink nodes) need to be updated, aggregated (in server) and distributed.

As none of the above mentioned works study prompt tuning in FL, in this paper, we design a communication-efficient prompt tuning method in FL and design backdoor attack to detect its vulnerability.

III. METHOD

A. Preliminaries

In FL setting, suppose there are K clients, each client hosts a private dataset $\mathcal{D}_k = \{(x_k, y_k)\}$ owning n_k samples. We use θ_t and θ_t^k to denote the global model and k^{th} local model parameters in communication round t respectively. Based on FedAvg [11], the aggregation process is computed as follows:

$$\theta_{t+1} = \sum_{k=1}^K \frac{n_k}{n} \theta_t^k \quad (1)$$

where $n = |\mathcal{D}| = \sum_{k=1}^K n_k$ is the total num of global combined data and $\mathcal{D} \triangleq \bigcup_{k \in [K]} \mathcal{D}_k$ is the global combined dataset. If data distributions are IID (Independent Identically Distribution), all clients have the same number of samples, then n_k/n could be replaced by $1/K$.

In a text classification task, x_k are the inputs and y_k are corresponding class labels. Each $x^{(i)} \in x_k$ consists of tokens $x^{(i)} = \{x_1^{(i)}, x_2^{(i)}, \dots, x_l^{(i)}\}$, where l is the length of single input. The prompt tuning structure is composed of the soft prompt \mathbf{p} , the template $\mathcal{T}(\cdot)$, the verbalizer $\mathcal{V}(\cdot)$ and the PLM

$\mathcal{M}(\cdot)$. Soft prompt \mathbf{p} consists of tokens $\mathbf{p} = \{p_1, p_2, \dots, p_m\}$, whose parameters are trainable. m is the number of the soft prompt tokens. $\mathcal{T}(\cdot)$ is a function to define where tokens of $x^{(i)}$ and \mathbf{p} are placed. After applying $\mathcal{T}(\cdot)$, we obtain $x_{prompt}^{(i)} = \mathcal{T}(x^{(i)}, \mathbf{p})$. At least one [MASK] token is placed into the $x_{prompt}^{(i)}$ for $\mathcal{M}(\cdot)$ to predict the label word. $\mathcal{V}(\cdot)$ is a map function to map the label word to the class $\hat{y} = \mathcal{V}(w)$. Usually, each class can have one or more label words. We call \mathcal{T} a multi-word verbalizer when each class has more than one label word, such as {positive: good, great; negative: bad, terrible;}. Input $x_{prompt}^{(i)}$ to \mathcal{M} , we can obtain the encoded feature [MASK]. By a softmax function, we can compute the probability that the label word w can fill the masked position. The label word with the highest probability is the predict word $w = \mathcal{M}(x_{prompt}^{(i)})$ and the predict class can be obtained by $\hat{y} = \mathcal{V}(w)$. We rewrite the prompt tuning process as $\hat{y}^{(i)} = f(x^{(i)}, \mathbf{p}, \theta)$,

B. FedPrompt

As mentioned before, in normal prompt tuning the whole is splitted into four parts, and only PLM (using fine-tuning) and soft prompt have trainable parameters. We use F and P to denote their parameters respectively, then the global model parameters in round t can be denoted as:

$$\theta_t = F_t + P_t \quad (2)$$

Algorithm 1 FedPrompt Algorithm

Input: K clients indexed by k , client fraction C , T communication rounds indexed by t , local minibatch size B , local epochs E , learning rate η , PLM parameters F , soft prompts parameters P .

Server executes:

- 1: Initialize global model.
- 2: Distribute F (fixed during the process) to all clients.
- 3: **for** $t \in \{1, \dots, T\}$ **do**
- 4: $U_t \leftarrow$ Select a subset of $C \cdot K$ clients at random
- 5: **for** each client $k \in U_t$ **do**
- 6: $P_t^k \leftarrow P_t$
- 7: $P_{t+1}^k \leftarrow$ ClientUpdate(k, P_t^k)
- 8: **end for**
- 9: $N_t = \sum_{k=1}^{|U_t|} n_k$
- 10: $P_{t+1} \leftarrow \sum_{k=1}^{|U_t|} \frac{n_k}{N_t} P_t^k$
- 11: **end for**
- 12: **return** P_{t+1}

ClientUpdate(k, P): // Run on client k

- 13: $\mathcal{B} \leftarrow$ (split D_k into batches of size B)
 - 14: **for** each local epoch $i \in \{1, \dots, E\}$ **do**
 - 15: **for** batch $b \in \mathcal{B}$ **do**
 - 16: $P \leftarrow P - \eta \nabla l(P; b)$
 - 17: **end for**
 - 18: **end for**
 - 19: **return** P
-

In FedPrompt, we fix F_t to learn a set of θ over \mathcal{D} with the objective to solve:

$$\arg \min_P \mathcal{L}(P) = \sum_{k=1}^K \frac{n_k}{n} \mathcal{L}_k(P) \quad (3)$$

where $\mathcal{L}_k(P)$ is the empirical loss of client k :

$$\mathcal{L}_k(P) = \mathbf{E}_{(x^{(i)}, y^{(i)}) \in \mathcal{D}_k} \ell_k(f(x^{(i)}, \mathbf{p}, P), y^{(i)}) \quad (4)$$

In the beginning, the server initializes the whole model, then distributes it to each client. At the beginning of round t , the server selects clients by fraction C to participate in this round, distributes the global soft prompt parameters P_t to them, and each selected client k replace the local P_{t-1}^k with P_t , which means $P_t^k = P_t$. As PLM is fixed, $F_t^k = F_{t-1}^k$. Then each client conducts local training with optimizer only for P_t^k , gets its updated soft prompt parameters P_t^k and sends them back to the server in parallel. The local training is same as normal prompt tuning process. Finally, the server performs the aggregation as follows:

$$P_{t+1} \leftarrow \sum_{k=1}^{\lceil C \cdot K \rceil} \frac{n_k}{N_t} P_t^k \quad (5)$$

where $N_t = \sum_{k=1}^{\lceil C \cdot K \rceil} n_k$ is participated data amount in round t . The whole process is shown as Algorithm 1. Except for

prompt tuning, there are also other prompt methods such as P-tuning [16] and Prefix-Tuning [15]. We also design FedPrompt for these prompt models in a similar way.

C. Poison FedPrompt

In FL, as clients privacy is highly protected and server have access to little information about client, it is widely acknowledged that multiple malicious clients possibly participate in training [34]. On the one hand, after initialization, each client has full knowledge of the model structure and parameters. On the other hand, it has been proved that prompt tuning is vulnerable to backdoor attack by poisoning training data [37]. Therefore, it is important to verify the robustness to backdoor attack of FedPrompt, and we call this attack as *FedPPT*.

Considering the situation that attacker has full control of one or more clients, and only modifies the local training data, which in fact is much less than attacker's access. The goal of attacker is to inject backdoor into poisoned prompt, which may be released to public. When victims use poisoned prompt, for clean samples, the victim PLMs will still give the correct label word; for poisoned samples which are added with the trigger word, the victim PLMs will output the target label word. To poison FedPrompt, firstly, modify the training dataset. Attacker tries to establish a shortcut between the trigger Δ and target label l_t . We define the poison function as $\mathcal{P}(\cdot)$, then we have single poisoned data $(x_p^{(i)}, t) = \mathcal{P}(x^{(i)}, \Delta, l_t)$, where modified target $t \neq y(x^{(i)})$. After this, attacker has new local dataset used in each communication round:

$$\mathcal{D}_k^{(poison)} = \{(x_p^{(i)}, t)\}, i \in \lambda n_k \quad (6)$$

$$\hat{\mathcal{D}}_k = \mathcal{D}_k^{(poison)} \cup \mathcal{D}_k \quad (7)$$

where λ is the poison rate. Secondly, using modified $\hat{\mathcal{D}}_k$ to update parameters P_k . Then the objective function of malicious client k as follows:

$$P_p^k = \arg \min_{P_p^k} \{ \mathbf{E}_{(x_k^{(i)}, y_k^{(i)}) \in \mathcal{D}_k} \ell_k(f(x_k^{(i)}, \mathbf{p}, P_p^k), y_k^{(i)}) + \mathbf{E}_{(x_k^{(i)}, y_k^{(i)}) \in \mathcal{D}_k^{(poison)}} \mathbf{1}_k(f(x_p^{(i)}, \mathbf{p}, P_p^k) \neq t) \} \quad (8)$$

IV. EXPERIMENTS

As no prior work on prompt tuning using FL has been done before, we investigate our methods on several federated NLP tasks including sentiment analysis and sentence-pair classification, where prompt tuning is suitable and often used. Also we conduct backdoor attack on these tasks to evaluate the robustness. All experiments are done on a server with 8 Nvidia Geforce GTX 1080Ti GPUs with 11GB RAM each, 12 Intel Xeon CPUs Processor, and CentOS release 7.9 OS. Our models are built using PyTorch framework¹.

¹<https://pytorch.org/>

A. Experimental Setup

1) *Dataset*: To evaluate FedPrompt model, our experiments are conducted on several NLP tasks:

- Text bi-classification tasks including sentiment analysis, toxicity detection and spam detection. For sentiment analysis, we use the Stanford Sentiment Treebank (SST-2)² and IMDB². We use the OffensEval [38] and the Twitter [39] in toxicity detection. And for spam detection, we use the Enron [40], and the Lingspam [41].
- Sentence-pair classification tasks. For this inference task, we use Question Natural Language Inference (QNLI) [42] and Recognizing Textual Entailment (RTE)² dataset.

To conduct experiments in FL setting, we divide all these datasets above into ten clients. In IID setting, we randomly divide the whole dataset into ten equal parts and each client has one part. In Non-IID setting, as all the tasks only having two labels $\{0,1\}$, we bring non-I.I.D.ness by different data quantity. We split all the data using Dirichlet distribution parameterized by α as in prior works [43]. Since labels are not available in the test sets for some datasets, we use the validation set as the test set and split a part of the training set as the validation set.

2) *Model and Training Details*: Among various PLMs, we select the most representative and widely used pre-trained language models, including the base versions of BERT [1], Roberta [2] and Google T5 [3] to conduct experiments. We use the Adam optimizer for training of BERT and Roberta, and the Adafactor optimizer for Google T5. In main experiments, we use a one-to-one verbalizer and a simple text classification template "[text] is [MASK]." having 20 soft prompt tokens in the head. Following the setup of [4], we set the learning rate to be 0.3. Following the setup of [44], we assume that we have a server and $K = 10$ available clients. We use a FedAvg system to implement the FL setting. Specifically, all clients are involved in the averaging of model parameters in every aggregation round. The number of max local step is set to 1000, compared to 30,000 in [4]. The number of communication rounds is set to 20, compared to 100 in [44] and [34], to prove our low-communication-cost and convergence-quick method.

3) *Baseline Algorithm*: To make a fair and reasonable comparison with our proposed FedPrompt, we choose the most related work [45], studying full-parameter fine-tuning, as FL baseline. Due to full-parameter fine-tuning requires lots of calculations, we only reproduce their method with above FL setting on IID SST-2 task as shown in Table II.

4) *Metric*: Communication bottleneck is a big challenge for many big models in FL, we use the amount of communicated parameters to evaluate our communication cost. As for the evaluation of performance, we use accuracy (*ACC*) which represents the proportion of the clean samples correctly classified by the model to measure the performance of the model on benign task. Also we use Attack Success Rate (*ASR*) to evaluate the attacking performance, which represents the

proportion of the poisoned samples we successfully enable the model to misclassify as the target class.

B. Main Results

In FedPrompt the learnable parameters is the same as communication cost. As shown in Table II, FedPrompt condenses the communication cost to nearly 0.01% of the full-parameter fine-tuning parameters, greatly reduces the communication cost, with only about 1% decrease in accuracy, making many devices applicable for some scenarios with communication and storage constraints, and the private data on these devices can contribute to the convergence of the global model. Also this property promotes the design and development of personalized FL model, especially for those resource-constrained devices.

The main results of FedPrompt performance with clean IID and Non-IID data distribution are summarized in Table I. As shown in Fig. 3, using FedPrompt to protect data privacy and handle the problem of few-shot demonstrate has little decrease on accuracy in most cases compared to prompt tuning without FL. Specifically, FL plays a remarkable effect with prompt tuning, only a few local training steps and communication rounds can contribute to a well-performed global model. For most tasks, FedPrompt achieves more than 90% *ACC* on clean data, and there is only a little decrease, almost less than 3%, with Non-IID data distribution than IID data distribution. Non-IID is a key challenge for the effectiveness of FL, and our proposed FedPrompt proves its compatibility on non-IID datasets. We think the FL paradigm and few-parameter soft prompt . We find that experiments on RTE task have a weaker result than other tasks. Considering that RTE only have 2240 training samples in total, which is the least among all tasks, and after splitting to ten clients each client only have a few samples to train soft prompt, we assume the weaker performance because of lack of data. Also RTE may need a better customized template to be used in prompt tuning.

As shown in Table III, we evaluate the effect of backdoor attack on all tasks and models with IID and Non-IID data distribution. Nearly all tasks get *ACC* on poison dataset drop less than 2% compared to on clean dataset in Table I. Even some tasks show a better *ACC* after poison. We think this is because poisoning the original dataset can be considered as data augmentation, and attacking has the similar effect to adversarial training. After backdoor attack, with poison ratio at 10% (all training data poisoned on 10% clients selected), all experiments on different tasks and models do not show a obvious rise in *ASR*, which suggests that FedPrompt has robustness to backdoor attack. We think this is because aggregation process offsets the backdoor.

C. Communication Rounds

Fig. 4 and Fig. 5 show the local and global *ACC* and *ASR* in each round during training. As for *ACC*, the training process of different settings is similar, and there is not a obvious decrease in Non-IID setting. The *ACC* of local model in the first round have a rapid rise and pass it to the global model only in one single communication round. This proves that our proposed

²<https://huggingface.co/datasets/>

TABLE I
ACC (%) AND ASR (%) OF FEDPROMPT WITH CLEAN IID AND NON-IID DATA DISTRIBUTION.

Dataset	BERT				ROBERTA				T5			
	IID		Non-IID		IID		Non-IID		IID		Non-IID	
	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR
SST-2	90.16	12.42	89.45	16.36	92.43	10.28	92.23	7.48	92.69	9.58	92.32	6.31
IMDB	91.08	12.66	89.26	11.42	92.80	9.15	92.53	7.66	92.89	9.69	91.24	11.33
OffensEval	82.64	9.84	80.47	8.55	81.05	13.87	80.34	5.98	79.30	12.58	78.83	10.65
Twitter	94.02	4.96	93.82	3.05	94.39	4.61	93.64	5.41	93.35	3.86	92.80	4.21
Enron	97.60	3.20	97.43	4.02	97.85	2.27	97.30	7.34	97.22	6.73	96.95	5.87
Lingspam	97.47	0.00	96.89	0.00	97.43	0.00	96.47	0.00	97.07	0.00	96.21	0.41
QNLI	83.36	28.35	82.25	30.46	86.44	14.81	85.43	12.10	89.06	10.87	84.48	12.44
RTE	54.87	35.21	54.15	42.73	60.32	36.99	57.51	44.52	76.51	22.95	73.64	22.60

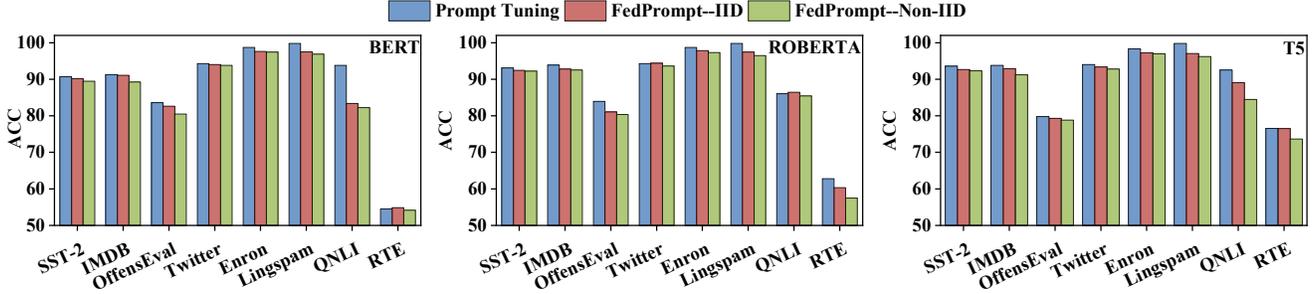


Fig. 3. The performance of prompt tuning without FL and FedPrompt. When using FL, there are IID setting and Non-IID setting on data distribution. The PLM used are BERT (the left), ROBERTA (the middle) and T5 (the right).

TABLE II
THE MAIN RESULTS OF FEDPROMPT AND FULL-PARAMETER FINE-TUNING ON IID SST-2 TASK. FOR THE SAME MODEL, WE REGARD THE PARAMETER QUANTITY IN FINE-TUNING AS 100.000%.

Model	FL Method	ACC	Comm. Cost	Ratio
BERT	FedPrompt	90.16	0.016M	0.014%
	Fine-tuning	91.02	109.530M	100.000%
ROBERTA	FedPrompt	92.43	0.016M	0.013%
	Fine-tuning	93.57	124.714M	100.000%
T5	FedPrompt	92.69	0.015M	0.007%
	Fine-tuning	93.79	222.919M	100.000%

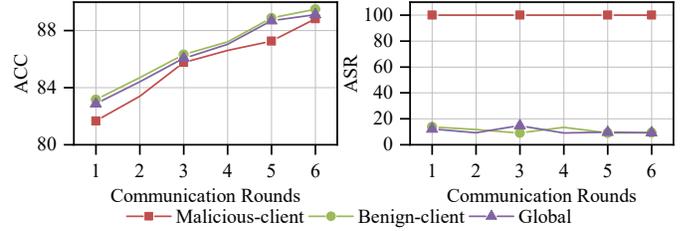


Fig. 5. Local and global ACC (%) and ASR (%) with communication rounds on SST-2 task using BERT. The results are using FedPPT with IID setting and only one fixed client in ten clients is malicious. The benign client is selected randomly.

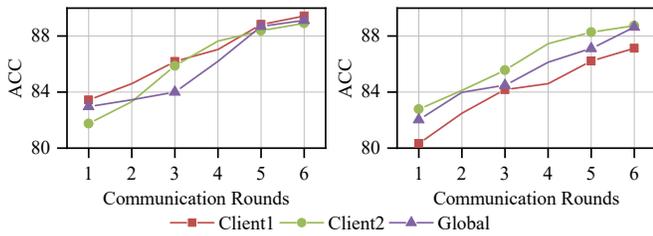


Fig. 4. Local and global ACC (%) with communication rounds on SST-2 task using BERT. The left one is using IID setting and the right one is using Non-IID setting, the two clients are selected randomly.

FedPrompt fits the poor data dependent prompt tuning well. As for ASR, we can see that ASR on the malicious client reaches

100% only after one single local training round, but after aggregation, ASR on benign client and global model remains a low level. We also find ASR on benign client is close to global model in the previous round, consistent with the aggregation method.

D. Number of Local Iterations

We study the effects of number of local iterations in each round. As we mentioned before, FedPrompt has relatively few trainable parameters that too many local iterations may lead to local over-fitting in FL while inadequate local iterations slow down the convergence. We conduct experiments on 100, 500, 1000 and 1500 local iterations, all of which are relatively small. As shown in Fig. 6, 100 and 500 local iterations performs worse and 1500 iterations may lead to local over-

TABLE III
ACC (%) AND ASR (%) OF FEDPROMPT WITH POISONED IID AND NON-IID DATA DISTRIBUTION. \uparrow MEANS HIGHER THAN CLEAN DATA.

Dataset	BERT				ROBERTA				T5			
	IID		Non-IID		IID		Non-IID		IID		Non-IID	
	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR
SST-2	89.11	13.30	88.76	14.60	91.55	11.55	92.12	9.73	92.20	8.74	91.51	9.07
IMDB	90.14	14.36	89.12	11.69	92.46	9.05	91.53	10.78	91.88	9.54	90.86	13.87
OffensEval	80.93	10.13	80.11	9.94	79.65	17.26	78.95	7.23	78.72	15.97	77.74	15.06
Twitter	94.10\uparrow	6.01	93.42	7.71	94.15\uparrow	4.02	93.22	4.76	93.25	5.28	92.98 \uparrow	4.13
Enron	97.37	4.20	98.18\uparrow	4.87	98.03\uparrow	3.53	97.16	8.20	97.88\uparrow	8.13	97.12 \uparrow	6.80
Lingspam	97.11	4.03	96.02	3.98	95.89	5.77	95.71	4.79	96.83	4.26	95.66	4.14
QNLI	84.48\uparrow	29.44	82.07	27.22	86.92\uparrow	18.82	85.30	8.33	85.56	9.15	84.33	14.26
RTE	54.51	31.23	60.29\uparrow	39.21	55.60	32.08	55.96	39.18	76.43	20.82	73.29	25.41

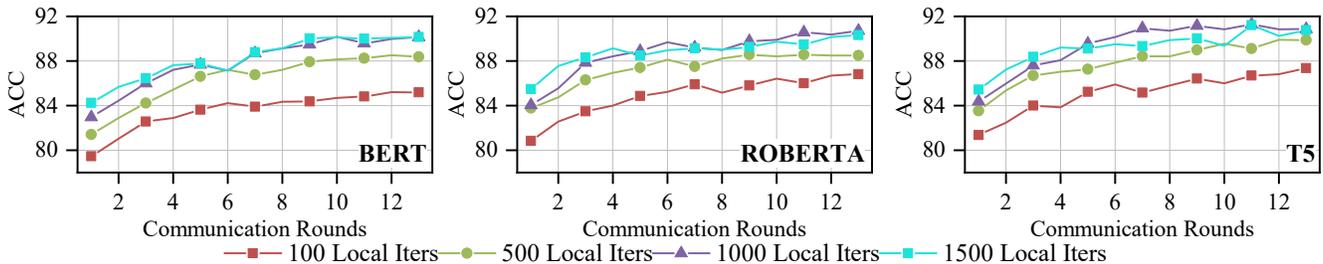


Fig. 6. Global ACC (%) results with different local iterations on SST-2 task.

TABLE IV
GLOBAL ACC (%) RESULTS WITH DIFFERENT NUMBER OF SOFT TOKENS ON SST-2 TASK.

Token Num	1	5	10	20
ACC	87.11	89.28	89.62	90.16

TABLE V
GLOBAL ACC (%) RESULTS WITH AND WITHOUT LDP ON SST-2 TASK.

Method	BERT	ROBERTA	T5
FedPrompt w/o LDP	90.16	92.43	92.69
FedPrompt w/ LDP	85.73	86.88	86.14

TABLE VI
GLOBAL ACC (%) AND COMMUNICATION COST (M, MILLION) WITH DIFFERENT PROMPT METHODS. DENOTE PROMPT TUNING AS METHOD α , P-TUNING AS β AND PREFIX-TUNING AS γ .

Method	BERT		ROBERTA		T5	
	ACC	Comm.	ACC	Comm.	ACC	Comm.
α	90.16	0.016	92.43	0.016	92.69	0.015
β	90.99	25.420	93.27	25.420	93.35	25.420
γ	-	-	-	-	76.85	9.853

fitting, which are harmful to obtain a excellent global model. Additional experiments expanded to 50 rounds suggest the performance of 100 and 500 local iterations still below others, while the other two could not get a further promotion.

E. Number of Soft Tokens

We tested the results under different numbers of soft tokens settings, and the results are consistent with those of prompt tuning under non-federal learning. As shown in Table IV, using more soft tokens will lead to better results. However, it also increases the communication cost under FL.

F. FedPrompt with LDP

As we mentioned before, there are hidden dangers to infer the origin private data by inverting gradients in FL, and LDP is an effective way to defense this attack. Also [46] has proved that the noise for larger model can damage the accuracy in differentially FL, so in Fedprompt the tiny prompt contributes to the use of LDP for privacy. We test on SST-2, clipping the gradients and then adding LaPlace noise on parameters. Table V shows that LDP protects the privacy with the cost of accuracy decreased by about 5%.

G. Prompt methods

We also experiment on P-tuning and Prefix-Tuning (only supports T5 now³). Our experiments on SST-2 are shown in Table VI. It suggests that among the three prompt methods prompt tuning gets the best performance combining ACC and communication cost. P-tuning has the best ACC performance but quite a lot parameters, and Prefix-Tuning needs more modification to use.

³<https://github.com/thunlp/OpenPrompt>

V. FURTHER IMPROVEMENT

Though FedPrompt is robust to normal backdoor attack in our experiments, there are still special methods to backdoor FL. We plan to let the server check the mean and standard deviation of soft prompt parameters from each client, and find outliers to refuse before global aggregation. Adding noise after aggregation could also destroy the backdoor, with partial sacrifice on ACC. We will carry on this research next.

VI. CONCLUSION

In this work, we propose FedPrompt to use federated prompt tuning on decentralized data in a communication-efficient and privacy preserving way. We employ a split aggregation way that freezing extensive PLMs' parameters and only tuning and aggregating soft prompts. In this way we condense the communication cost to only 0.01% compared to full-parameter fine-tuning, making many devices applicable for some scenarios with communication constraints. Experiments on both IID and Non-IID data distribution using three mainstream model demonstrate the accuracy of FedPrompt. We also use LDP to further protect the privacy, and it is necessary to further study the FL backdoor attack.

REFERENCES

- [1] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," in *NAACL-HLT*, 2019, pp. 4171–4186.
- [2] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized BERT pretraining approach," p. arXiv:1907.11692, 2019.
- [3] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *J. Mach. Learn. Res.*, vol. 21, pp. 140:1–140:67, 2020.
- [4] B. Lester, R. Al-Rfou, and N. Constant, "The power of scale for parameter-efficient prompt tuning," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021, pp. 3045–3059.
- [5] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE Computer Society, 2016, pp. 770–778.
- [6] Z. Wu, Y. Zhou, D. Wu, M. Chen, and Y. Xu, "TAMF: towards personalized time-aware recommendation for over-the-top videos," in *Proceedings of the ACM Workshop on Network and Operating Systems Support for Digital Audio and Video*, 2019, pp. 43–48.
- [7] J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, Q. V. Le, M. Z. Mao, M. Ranzato, A. W. Senior, P. A. Tucker, K. Yang, and A. Y. Ng, "Large scale distributed deep networks," in *Advances in Neural Information Processing Systems*, 2012, pp. 1232–1240.
- [8] M. Li, L. Zhou, Z. Yang, A. Q. Li, F. Xia, D. G. Andersen, and A. Smola, "Parameter server for distributed machine learning," 2013.
- [9] A. Reiszadeh, A. Mokhtari, H. Hassani, A. Jadbabaie, and R. Pedarsani, "Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization," in *Artificial Intelligence and Statistics*. PMLR, 2020, pp. 2021–2031.
- [10] J. Hamer, M. Mohri, and A. T. Suresh, "Fedboost: A communication-efficient algorithm for federated learning," in *Proceedings of the International Conference on Machine Learning*. PMLR, 2020, pp. 3973–3983.
- [11] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial Intelligence and Statistics*. PMLR, 2017, pp. 1273–1282.
- [12] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning: Concept and applications," *ACM Trans. Intell. Syst. Technol.*, vol. 10, no. 2, pp. 12:1–12:19, 2019.
- [13] T. Qi, F. Wu, C. Wu, Y. Huang, and X. Xie, "Privacy-preserving news recommendation model learning," in *Findings of the Association for Computational Linguistics: EMNLP*, 2020, pp. 1423–1432.
- [14] K. Hambardzumyan, H. Khachatryan, and J. May, "WARP: word-level adversarial reprogramming," in *ACL/IJCNLP*, 2021, pp. 4921–4933.
- [15] X. L. Li and P. Liang, "Prefix-tuning: Optimizing continuous prompts for generation," in *ACL/IJCNLP*, 2021, pp. 4582–4597.
- [16] X. Liu, Y. Zheng, Z. Du, M. Ding, Y. Qian, Z. Yang, and J. Tang, "Gpt understands, too," p. arXiv:2103.10385, 2021.
- [17] X. Liu, K. Ji, Y. Fu, Z. Du, Z. Yang, and J. Tang, "P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks," p. arXiv:2110.07602, 2021.
- [18] X. Han, W. Zhao, N. Ding, Z. Liu, and M. Sun, "Ptr: Prompt tuning with rules for text classification," p. arXiv:2105.11259, 2021.
- [19] Y. Gu, X. Han, Z. Liu, and M. Huang, "PPT: pre-trained prompt tuning for few-shot learning," in *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 2022, pp. 8410–8423.
- [20] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated optimization in heterogeneous networks," in *Proceedings of Machine Learning and Systems*, 2020.
- [21] S. P. Karimireddy, S. Kale, M. Mohri, S. J. Reddi, S. U. Stich, and A. T. Suresh, "SCAFFOLD: stochastic controlled averaging for federated learning," in *Proceedings of the International Conference on Machine Learning*. PMLR, 2020, pp. 5132–5143.
- [22] A. Fallah, A. Mokhtari, and A. E. Ozdaglar, "Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach," in *Advances in Neural Information Processing System*, 2020.
- [23] D. A. E. Acar, Y. Zhao, R. M. Navarro, M. Mattina, P. N. Whatmough, and V. Saligrama, "Federated learning based on dynamic regularization," in *International Conference on Learning Representations*, 2021.
- [24] X. Li and D. Zhan, "Fedrs: Federated learning with restricted softmax for label distribution non-iid data," in *KDD '21: The ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2021, pp. 995–1005.
- [25] J. Yi, F. Wu, C. Wu, R. Liu, G. Sun, and X. Xie, "Efficient-fedrec: Efficient federated learning framework for privacy-preserving news recommendation," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2021, pp. 2814–2824.
- [26] L. Melis, C. Song, E. D. Cristofaro, and V. Shmatikov, "Exploiting unintended feature leakage in collaborative learning," in *IEEE Symposium on Security and Privacy*. IEEE, 2019, pp. 691–706.
- [27] J. Geiping, H. Bauermeister, H. Dröge, and M. Moeller, "Inverting gradients - how easy is it to break privacy in federated learning?" in *Advances in Neural Information Processing System*, 2020.
- [28] M. Abadi, A. Chu, I. J. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, "Deep learning with differential privacy," in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, 2016, pp. 308–318.
- [29] R. C. Geyer, T. Klein, and M. Nabi, "Differentially Private Federated Learning: A Client Level Perspective," p. arXiv:1712.07557, Dec. 2017.
- [30] A. Triastcyn and B. Faltings, "Federated learning with bayesian differential privacy," in *IEEE International Conference on Big Data (IEEE BigData)*. IEEE, 2019, pp. 2587–2596.
- [31] K. A. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical secure aggregation for privacy-preserving machine learning," in *Proceedings of the ACM SIGSAC Conference on Computer and Communications Security*, 2017, pp. 1175–1191.
- [32] H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang, "Learning differentially private recurrent language models," in *International Conference on Learning Representations*, 2018.
- [33] M. Jere, T. Farman, and F. Koushanfar, "A taxonomy of attacks on federated learning," *IEEE Secur. Priv.*, vol. 19, no. 2, pp. 20–28, 2021.
- [34] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, "How to backdoor federated learning," in *Artificial Intelligence and Statistics*. PMLR, 2020, pp. 2938–2948.
- [35] X. Yuan, P. He, Q. Zhu, and X. Li, "Adversarial examples: Attacks and defenses for deep learning," *IEEE Trans. Neural Networks Learn. Syst.*, vol. 30, no. 9, pp. 2805–2824, 2019.
- [36] A. Aldahdooh, W. Hamidouche, S. A. Fezza, and O. Déforges, "Adversarial example detection for DNN models: a review and experimental comparison," *Artif. Intell. Rev.*, vol. 55, no. 6, pp. 4403–4462, 2022.

- [37] W. Du, Y. Zhao, B. Li, G. Liu, and S. Wang, "Ppt: Backdoor attacks on pre-trained models via poisoned prompt tuning," in *IJCAI-22*, 2022, pp. 680–686.
- [38] M. Zampieri, S. Malmasi, P. Nakov, S. Rosenthal, N. Farra, and R. Kumar, "Semeval-2019 task 6: Identifying and categorizing offensive language in social media (offenseval)," in *Proceedings of the 13th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2019*, 2019, pp. 75–86.
- [39] A. Founta, C. Djouvas, D. Chatzakou, I. Leontiadis, J. Blackburn, G. Stringhini, A. Vakali, M. Sirivianos, and N. Kourtellis, "Large scale crowdsourcing and characterization of twitter abusive behavior," in *Proceedings of the Twelfth International Conference on Web and Social Media*, 2018, pp. 491–500.
- [40] V. Metsis, I. Androutsopoulos, and G. Paliouras, "Spam filtering with naive bayes - which naive bayes?" in *CEAS*, 2006.
- [41] G. Sakkis, I. Androutsopoulos, G. Paliouras, V. Karkaletsis, C. D. Spyropoulos, and P. Stamatopoulos, "A memory-based approach to anti-spam filtering for mailing lists," *Inf. Retr.*, vol. 6, no. 1, pp. 49–73, 2003.
- [42] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, "Squad: 100, 000+ questions for machine comprehension of text," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2016, pp. 2383–2392.
- [43] C. He, E. Ceyani, K. Balasubramanian, M. Annavam, and S. Avestimehr, "Spreadgnn: Decentralized multi-task federated learning for graph neural networks on molecular data," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022, pp. 6865–6873.
- [44] Q. Li, B. He, and D. Song, "Model-contrastive federated learning," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 10 713–10 722.
- [45] A. Hilmkil, S. Callh, M. Barbieri, L. R. Sütfeld, E. L. Zec, and O. Mogren, "Scaling federated learning for fine-tuning of large language models," in *International Conference on Applications of Natural Language to Information Systems*, 2021, pp. 15–23.
- [46] P. Basu, T. S. Roy, R. Naidu, Z. Müftüoğlu, S. Singh, and F. Miresghalah, "Benchmarking differential privacy and federated learning for BERT models," p. arXiv:2106.13973, 2021.