FEDONE: QUERY-EFFICIENT FEDERATED LEARNING FOR BLACK-BOX DISCRETE PROMPT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Black-Box Discrete Prompt Learning (BDPL) is a prompt-tuning method that optimizes discrete prompts without accessing model parameters or gradients, making the prompt tuning on a cloud-based Large Language Model (LLM) feasible. Adapting Federated Learning (FL) to BDPL could further enhance prompt tuning performance by leveraging data from diverse sources. However, all previous research on federated black-box prompt tuning had neglected the substantial query cost associated with the cloud-based LLM service. To address this gap, we conducted a theoretical analysis of query efficiency within the context of federated black-box prompt tuning. Our findings revealed that degrading FedAvg to activate only one client per round, a strategy we called FedOne, enabled optimal query efficiency in federated black-box prompt learning. Building on this insight, we proposed the FedOne framework, a federated black-box discrete prompt learning method designed to maximize query efficiency when interacting with cloud-based LLMs. We conducted numerical experiments on various aspects of our framework, demonstrating a significant improvement in query efficiency, which aligns with our theoretical results.

025 026 027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

028 029

Prompt tuning has emerged as a vital technique for adapting large language models (LLMs) (Liu et al., 2019; Brown et al., 2020) to specific tasks without retraining the entire model. Tradition-ally, many tuning methods require access to the model's intermediate representations (Brown et al., 2020; Li & Liang, 2021; Liu et al., 2021; Lester et al., 2021), categorizing them as white-box approaches. However, when such access is unavailable, black-box prompt tuning becomes essential. This approach focuses on tuning the input prompts without access to the internal processes of the model (Sun et al., 2022; Diao et al., 2022; Deng et al., 2022; Xiao et al., 2023).

Federated learning (FL) (McMahan et al., 2017; Li et al., 2020; 2021; Karimireddy et al., 2020;
Mishchenko et al., 2022) has emerged as a promising approach for leveraging decentralized data from multiple clients while preserving privacy. Applying federated learning to prompt tuning offers a valuable opportunity to improve client capabilities and enhance model performance. Most feder-ated prompt tuning methods to date have focused on white-box scenarios where clients have access to model parameters (Zhao et al., 2023; Zhang et al., 2023). In those approaches, only the trainable parameters (prompts) are trained by the client and shared with the server, significantly reducing both the number of trainable parameters and the communication costs compared to fine-tuning baselines.

044 However, several practical limitations hinder the applicability of federated prompt tuning that relies 045 on white-box access. First, white-box prompt learning is not applicable to closed-source LLM that 046 are accessed via APIs, as these models are not openly shared. In such scenarios, users are restricted 047 to interacting with the model through the API endpoints without access to the internal structure or 048 weights of the LLM. This limitation prevents the application of white-box prompt learning techniques. Second, FL is typically applied in scenarios involving thousands of edge devices, each with limited computational resources. However, white-box prompt learning demands substantial compu-051 tational power, as it requires devices to perform computation on the entire LLM. These operations are computationally intensive, rendering them impractical for edge devices with constrained capa-052 bilities. This presents a significant challenge when trying to implement white-box prompt learning in FL environments, as it can lead to excessive resource demands on the participating devices.



Figure 1: Query-Efficient Fed-BDPL

In contrast, applying black-box discrete prompt learning to federated learning offers several distinct benefits (Zhao et al., 2023; Lin et al., 2023; Zhang et al., 2023; Che et al., 2023), enhancing both the 071 practicality and effectiveness. First, this approach preserves the privacy of closed-source LLMs by 072 not requiring access to their internal model weights or architecture. For example, Diao et al. (2022); 073 Lin et al. (2023) introduces a black-box prompt learning method that uses only discrete prompts as 074 inputs and relies solely on the model's output loss to train the prompts. Besides, black-box prompt 075 tuning reduces the computational burden on the client, as it eliminates the need for computation 076 on the entire model, thereby enabling participation from edge devices with limited computational 077 resources. Moreover, the communication costs associated with discrete prompts are lower compared 078 to white-box prompt tuning methods, as the white-box prompt tuning require transmitting large 079 matrices of numerical values.

Despite these advantages, the application of BDPL in FL still faces two significant unresolved challenges, tempering its overall promise. First, previous research (Lin et al., 2023; Sun et al., 2023) on federated black-box prompt tuning has neglected the substantial cost associated with queries to the LLM cloud service (Figure 1, left). Second, a convergence analysis for Federated BDPL aimed at optimizing discrete prompts has not yet been provided.

Targeting the above problems, we introduce a novel federated learning framework, FedOne, de-signed to optimize query efficiency in Federated BDPL. We offer the first convergence analysis of Fed-BDPL in this context and further extend the analysis of the query efficiency towards the cloud-based LLM server. The results demonstrate that by limiting activation to a single client per round, FedOne achieves optimal query efficiency (Figure 1, right). Our approach is particularly well-suited for scenarios involving limited computational resources, such as mobile devices or IoT systems, where local training of LLMs is impractical.

Formally, this paper makes the following key contributions:

- We identify that existing federated black-box prompt tuning methods overlook the significant costs associated with querying cloud-based LLM services.
- To address this gap, we present the first theoretical analysis of Federated BDPL, with a focus on understanding and evaluating the query efficiency when interacting with cloud-based LLMs.
- Based on our analysis, we introduce the FedOne framework, a novel approach designed to optimize query efficiency in Federated Black-box Prompt Learning by activating only one client per round when querying cloud-based LLMs.
- 101 102

2 Method

2.1 FEDERATED BLACK-BOX PROMPT LEARNING FRAMEWORK

104 105

103

092

093

094

096

098

099

100

068

069

In the federated black-box prompt tuning framework, there is one central aggregation server and Kclients. Each client, indexed by k, possesses a dataset D^k , consisting of input sentences Ψ^k and their corresponding labels Y^k , i.e., $D^k = {\Psi^k, Y^k}$. The dataset Ψ^k contains a total of M^k input sentences, each represented as ψ_m^k , i.e., $\Psi^k = \{\psi_m^k\}_{m=1}^{M^k}$. Similarly, the corresponding labels y_m^k comprise Y^k , i.e. $Y^k = \{y_m^k\}_{m=1}^{M^k}$.

Each client generates a discrete sequence of prompt tokens $\Phi^k = \phi_1^k \cdots \phi_i^k \cdots \phi_n^k$ from a trainable parameter $\boldsymbol{\alpha}^k \in \mathbb{R}^{n \times N}$. This learnable parameter, $\boldsymbol{\alpha}^k$, is transmitted to the FL aggregation server for averaging. Details of how to generate the discrete prompt Φ^k from $\boldsymbol{\alpha}^k$ and the local training of $\boldsymbol{\alpha}^k$ through interaction with the cloud-based LLM service will be discussed in the next subsection regarding the local black-box prompt learning on the client.

The server and clients collaboratively solve a minimization problem, aiming to reduce a global loss function that aggregates the local loss functions from the clients. This can be expressed as:

119 120

121 122

125

126 127

128

129

130

131

 $\min_{\Phi} \left\{ \mathcal{L}(\Phi; \Psi) \triangleq \sum_{k=1}^{K} \frac{M^{k}}{M} \mathcal{L}^{k}(\Phi; \Psi^{k}) \right\}, \ \mathcal{L}^{k}(\Phi; \Psi^{k}) = \frac{1}{M^{k}} \sum_{m=1}^{M^{k}} \ell\left(\Phi; \psi_{m}^{k}, y_{m}^{k}\right)$ (1)

where $\mathcal{L}(\Phi; \Psi)$ represents the global objective function of the FL, $\Psi = \{\Psi^k\}_{k=1}^K$. $\mathcal{L}^k(\cdot, \Psi^k)$ is the local objective function of client k. The function $\ell(\cdot, \cdot)$ denotes the loss function.

2.2 BLACK-BOX DISCRETE PROMPT LEARNING ON THE CLIENT

We now discuss the local training process of the client through interactions with the cloud-based LLM service. In this section, all variables have the superscript k, indicating that they belong to the k-th client. However, the reader can ignore this superscript and treat it as a standalone training process on a single machine.

Generating the discrete prompt sequence Φ^k The sequence of the discrete prompt Φ^k is generated from a vocabulary $\mathcal{V} = \{\mathcal{V}[j]\}_{j=1}^N$, which contains a total of N token options. Each token ϕ_i^k in the prompt sequence Φ^k is selected from the vocabulary, i.e., $\Phi^k = \phi_1^k \cdots \phi_i^k \cdots \phi_n^k = \mathcal{V}[j_1^k] \cdots \mathcal{V}[j_k^i] \cdots \mathcal{V}[j_n^k]$. For the *i*-th token ϕ_i^k , the prompt index j_i^k is sampled from the categorical distribution \mathbf{p}_i^k , i.e., $j_i^k \sim \operatorname{Cat}(\mathbf{p}_i^k)$. Note that $\mathbf{p}_i^k = [p_{i,1}^k, \dots p_{i,N}^k]$, where the element $p_{i,j}^k$ represents the probability that the token ϕ_i^k is selected as V[j] from the vocabulary \mathcal{V} .

Directly optimizing the p_i^k may cause trouble in the convergence analysis as the gradient of the categorical distribution is biased. To address this, we re-parameterize the categorical distribution p_i^k using the Gumbel-Softmax technique (Jang et al., 2016) and introduce the parameter $\alpha_i^k =$ $\alpha_{i,1}^k, \dots, \alpha_{i,N}^k$] as the learnable parameter. The re-parameterization is shown below:

$$p_{i,j}^{k} = \frac{\exp\left(\frac{\log(\alpha_{i,j}^{k}) + g_{i,j}^{k}}{\tau}\right)}{\sum_{l=1}^{N} \exp\left(\frac{\log(\alpha_{i,l}^{k}) + g_{i,l}^{k}}{\tau}\right)}$$
(2)

144 145 146

147

148

149 150

156 157

159

161

143

where $\tau > 0$ is the temperature parameter, $g_{i,j}^k \sim \text{Gumbel}(0,1)$ is the Gumbel random variable. We denote the Gumbel-Softmax function as GS, i.e. $\mathbf{p}^k = \text{GS}(\mathbf{a}^k)$.

Optimizing the learnable parameter $\boldsymbol{\alpha}^k$ To compute the gradient with respect to the learnable parameter $\boldsymbol{\alpha}^k$, we first define the expected loss over the sequence of prompts in Eq. 3. The *i*-th token ϕ_i^k is generated from the vocabulary by sampling the prompt index from the categorical distribution, i.e. $\phi_i^k = \mathcal{V}[j_i^k]$, where $j_i^k \sim \operatorname{Cat}(\boldsymbol{p}_i^k)$. For brevity, this sampling process is denoted as $\phi_i^k \sim \boldsymbol{p}_i^k$. We can define the expectation of the loss for the distribution of the prompt as follows:

$$\mathbb{E}_{\Phi^k \sim \boldsymbol{p}^k} \left[\mathcal{L}(\Phi^k, \Psi^k) \right] = \sum_{\phi_1^k \sim \boldsymbol{p}_1^k} \cdots \sum_{\phi_n^k \sim \boldsymbol{p}_n^k} \left(\mathcal{L}(\Phi^k, \Psi^k) \prod_{i=1}^n P(\phi_i^k) \right)$$
(3)

Following the same steps in (Diao et al., 2022, Eq. (2)), we can estimate the gradient w.r.t. $\boldsymbol{\alpha}_i^k$ by:

$$\nabla_{\boldsymbol{\alpha}_{i}^{k}} \mathbb{E}_{\Phi^{k} \sim \mathrm{GS}(\boldsymbol{\alpha}^{k})} \left[\mathcal{L}(\Phi^{k}, \Psi^{k}) \right] = \mathbb{E}_{\Phi^{k} \sim \mathrm{GS}(\boldsymbol{\alpha}^{k})} \left[\mathcal{L}(\Phi^{k}, \Psi^{k}) \nabla_{\boldsymbol{\alpha}_{i}^{k}} \log P(\phi_{i}^{k}) \right]$$
(4)

The *j*-th component of $\nabla_{\boldsymbol{\alpha}_{i}^{k}}\log P(\phi_{i}^{k})$ could be explicitly computed as follows (with detailed steps provided in Appendix A.2):

$$\nabla_{\alpha_{i,j}^{k}} \log P(\phi_{i}^{k}) = \nabla_{\alpha_{i,j}^{k}} \log p_{i,j_{i}^{k}}^{k} = \begin{cases} \frac{1 - p_{i,j_{k}}^{k}}{\tau \alpha_{i,j_{k}}^{k}} & j = j_{i}^{k} \\ -\frac{p_{i,j}^{k}}{\tau \alpha_{i,j}^{k}} & j \neq j_{i}^{k} \end{cases}$$
(5)

Then, we employ the mini-batch stochastic variance-reduced policy (MB-SVRP) estimator (Diao et al., 2022; Williams, 1992) to reduce the variance of the sampling when computing the gradient. This involves sampling the prompt sequence Φ^k from the distribution p^k multiple times, with the number of samplings denoted by *I*. The MB-SVRP estimator is then computed as follows:

$$\hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}}f^{k}(\boldsymbol{\alpha}^{k},\boldsymbol{\mathcal{B}}^{k}) = \frac{1}{I-1}\sum_{r=1}^{I}\left[\left(\ell(\Phi^{k,r};\boldsymbol{\mathcal{B}}^{k}) - \frac{1}{I}\sum_{w=1}^{I}\ell(\Phi^{k,w};\boldsymbol{\mathcal{B}}^{k})\right)\nabla_{\boldsymbol{\alpha}_{i}^{k}}\log P(\phi_{i}^{k,r})\right]$$
(6)

where $\{\Phi^{k,r}\}_{r=1}^{I}$ are sampled independently from $\mathbf{p}^{k} = GS(\mathbf{\alpha}^{k})$. The mini-batch \mathcal{B}^{k} , with size B^{k} , is sampled from the dataset Ψ^{k} . Note that the clients are unable to directly compute $\ell(\Phi^{k,*}, \mathcal{B}^{k})$ on their own. They must transmit both the sampled prompt Φ^{k} and the mini-batch \mathcal{B}^{k} to the cloud-based LLM service, which then computes the loss $\ell(\Phi^{k,*}, \mathcal{B}^{k})$ and returns the result. Finally, with the learning rate set to η , the update of $\boldsymbol{\alpha}_{i}^{k}$ at the *t*-th iteration is expressed as follows:

$$\boldsymbol{\alpha}_{i,(t+1)}^{k} = \boldsymbol{\alpha}_{i,(t)}^{k} - \eta \cdot \hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}^{k}, \mathcal{B}_{t}^{k})$$
(7)

2.3 Algorithm

Algorithm 1 outlines the Fed-BDPL framework, which integrates federated averaging with local client training with Gumbel-Softmax-BDPL (GS-BDPL). The FL aggregation server randomly se-lects a subset of clients and broadcasts the trainable parameters to them. Each selected client k then performs local training on the parameters and return the updated parameters α^k to the server. The local training of the selected client is conducted through black-box prompt learning by querying the cloud-based LLM service, as detailed in Section 2.2. Finally, the server aggregates the updates by averaging the parameters from all participating clients. This process is iteratively repeated through-out the training. Specifically, in the FedOne framework, the number of activated clients is set to 1, as highlighted in the light green box. The reasons behind FedOne's high query efficiency on cloud-based LLM services will be formally analyzed in the next section.

Algorithm 1 Fed-BDPL. C denotes the sampling ratio of the clients. K_* represents the number of selected clients, and U is the set of selected clients.

203	1: Server executes:	13: Client executes:
204	2: Initialize α	14: function Client_Update(k, α^k):
205	3: for $s = 0, 1, \dots S - 1$ do	15: for $e = 1,, E$ do
206	4: FedAvg: $K_* \leftarrow \max(C \cdot K, 1)$	16: Re-parameterize the categorical distribution $n^k = GS(\alpha^k)$ using Eq. 2
207	5: FedOne: $K_* \leftarrow 1$	$p = OS(\alpha)$ using Eq. 2. 17: Query the cloud-based LLM server to obtain
209	6: $U_s \leftarrow (\text{sampling } K_* \text{ clients})$ 7: for $k \in U_s$ in parallel do	$\left\{ \mathcal{L}(\Phi^{k,r}, \mathcal{B}^k_e) \right\}_{\substack{r=1 \\ e^k (-k, \mathcal{D}^k)}}^{I}$
211	8: $\boldsymbol{\alpha}^k \leftarrow \boldsymbol{\alpha}$ 9: $\boldsymbol{\alpha}^k \leftarrow \text{Client Undate}(k, \boldsymbol{\alpha}^k)$	18: Compute $\bigvee_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}^{k}, \mathcal{B}_{e}^{k})$ using Eq. 6. 10: $\boldsymbol{\alpha}_{i}^{k} = \boldsymbol{\alpha}_{i}^{k} = \boldsymbol{\alpha}_{i}^{k} \hat{\boldsymbol{\alpha}}_{i}^{k} + \boldsymbol{\beta}_{e}^{k} \hat{\boldsymbol{\beta}}_{i}^{k}$
212	10: end for	19. $\alpha_i \leftarrow \alpha_i - \eta \cdot \nabla_{\alpha_i^k} J$ (α , B_e) 20: end for
213	11: $\boldsymbol{\alpha} \leftarrow rac{1}{K_*} \sum_{k \in U_s} \boldsymbol{\alpha}^k$	21: Return α^k
214	12: end for	

216 3 CONVERGENCE ANALYSIS

3.1 Assumptions

Assumption 1. Unbiasedness and bounded variance of stochastic gradient: We assume that the stochastic gradient is unbiased and has bounded variance.

$$\mathbb{E}_{\psi_m^k} \left[\nabla_{\boldsymbol{\alpha}_i^k} f^k(\boldsymbol{\alpha}^k, \psi_m^k) \right] = \nabla_{\boldsymbol{\alpha}_i^k} F(\boldsymbol{\alpha}^k, \Psi^k) \tag{8}$$

$$\mathbb{E}_{\psi_m^k} \left\| \nabla_{\boldsymbol{\alpha}_i^k} f^k(\boldsymbol{\alpha}^k, \psi_m^k) - \mathbb{E}_{\psi_m^k} \left[\nabla_{\boldsymbol{\alpha}_i^k} f^k(\boldsymbol{\alpha}^k, \psi_m^k) \right] \right\|^2 \le \sigma_{\psi}^2 \tag{9}$$

Assumption 1 is the basic assumption for solving non-convex optimization problems using stochastic gradient descent (Ghadimi & Lan, 2013; Hazan & Kale, 2014; Xu et al., 2019; Liu et al., 2020). Assumption 2. Bounded loss: At the k-th client, $\forall \psi_{k}^{k} \in \Psi^{k}$, and Φ^{k} sampled by \mathbf{n}^{k} , we perform

Assumption 2. Bounded loss: At the k-th client, $\forall \psi_m^k \in \Psi^k$, and Φ^k sampled by \mathbf{p}^k , we perform a clipping operation with a constant G for loss function $\ell(\cdot, \cdot)$:

$$\left|\ell(\Phi^k, \psi_m^k)\right| \le G \tag{10}$$

Assumption 2 ensures that the loss value is bounded, primarily to regulate the loss during *I*-sample estimation in stochastic policy gradient, thereby facilitating theoretical analysis.

Assumption 3. Clients' heterogeneity by weighted gradient diversity: For client k = 1, ..., K with sampling probability vector $\mathbf{q} = \{q^{[k]}\}_{k=1}^{K}$, and $\nabla_{\boldsymbol{\alpha}_{i}^{k}}F^{k}(\boldsymbol{\alpha}^{k}, \Psi^{k})$ is local gradient of $\boldsymbol{\alpha}_{i}$ w.r.t. all input sentence in client k. We assume that λ is the upper bound on the weighted gradient diversity across local objectives, i.e.

$$\Lambda(\boldsymbol{\alpha}_{i},\boldsymbol{q}) \triangleq \frac{\sum_{k=1}^{K} q^{[k]} \cdot \left\| \nabla_{\boldsymbol{\alpha}_{i}^{k}} F^{k}(\boldsymbol{\alpha}^{k},\Psi^{k}) \right\|^{2}}{\left\| \sum_{k=1}^{K} q^{[k]} \cdot \nabla_{\boldsymbol{\alpha}_{i}^{k}} F^{k}(\boldsymbol{\alpha}^{k},\Psi^{k}) \right\|^{2}} \leq \lambda$$
(11)

Assumption 3 introduces gradient diversity as a measure of client heterogeneity, which is used to quantify differences among clients and establish convergence conditioned in heterogeneous FL settings (Yin et al., 2018; Haddadpour & Mahdavi, 2019).

3.2 THE CONVERGENCE OF FEDERATED BLACK-BOX DISCRETE PROMPT LEARNING

Theorem 1. Suppose assumption 1, 2 and 3 hold, using algorithm 1 to solve the FedBDPL problem defined in Eq. 1. Let $B = \min \{B^1, ..., B^K\}$ where B^k represents the local mini-batch size for each client. Set $\alpha_{i,j} \ge \nu > 0$. The variance of the variance-reduced policy gradient is given by $\sigma_{\alpha}^2 = \frac{8G^2N}{\tau^2\nu^2}$. $\mathbb{E}_{\Phi^k \sim GS(\mathbf{a}^k)} \left[\mathcal{L}(\Phi^k, \Psi^k) \right]$ is L-smooth w.r.t. \mathbf{a}^k where $L = \frac{nGN(\tau+1)}{\tau^2\nu^2}$. If the learning rate η satisfies the following condition:

$$0 < \eta \le \eta^* = \frac{-\lambda L + \sqrt{\lambda^2 L^2 + 8L^2 E\left(1 + \frac{1}{K_*}\right)}}{8L^2 E\left(1 + \frac{1}{K_*}\right)}$$
(12)

Then, the Fed-BDPL's expected gradient, $\nabla_{\alpha} F(\alpha_t, \Psi^k)$, can be bounded as follows:

$$\frac{1}{T} \sum_{t=0}^{T-1} \left\| \nabla_{\boldsymbol{\alpha}} F(\boldsymbol{\alpha}_t, \Psi^k) \right\|^2 \le \frac{4G}{\eta T} + \frac{2(E+1)L^2 \eta^2 n \sigma_{\psi}^2 (1+\frac{1}{K_*}) + 2nL\eta \sigma_{\psi}^2}{B} + \frac{2(E+1)L^2 \eta^2 n \sigma_{\alpha}^2 (1+\frac{1}{K_*}) + 2nL\eta \sigma_{\alpha}^2}{I^2}$$
(13)

Remark 1. The local gradient of BDPL is unbiased within each client, but due to data heterogeneity across clients, it becomes biased with respect to the global gradient (Haddadpour & Mahdavi, 2019). This bias can be mitigated by reducing the step size η . However, as the gradient diversity, λ , increases, an even smaller step size η is needed, which in turn slows down the convergence process.

270 271 Corollary 1. Convergence rate and complexity 271

1) Convergence rate: Let $\eta = \min\left\{\eta^*, \frac{1}{\sqrt{T}}, \frac{1}{L}\right\}$. Under this condition, the following holds:

$$\frac{1}{T} \sum_{t=0}^{T-1} \left\| \nabla_{\boldsymbol{\alpha}} F(\boldsymbol{\alpha}_t, \Psi) \right\|^2 = \mathcal{O}(\frac{1}{\sqrt{T}})$$
(14)

2) Complexity: To guarantee an ϵ -solution, such that $\frac{1}{T} \sum_{t=0}^{T-1} \left\| \nabla_{\alpha} F(\alpha_t, \Psi^k) \right\|^2 \leq \epsilon^2$, the following condition must hold:

$$T_{\epsilon} = \mathcal{O}\left(\frac{1}{\epsilon^4}\right) \tag{15}$$

Fed-BDPL achieves a sub-linear convergence rate of $\left(\frac{1}{\sqrt{T}}\right)$ and a complexity $\mathcal{O}\left(\frac{1}{\epsilon^4}\right)$, which is comparable to that of classical non-convex smooth SGD algorithms (Allen-Zhu, 2018; Khaled & Richtárik, 2020; Gower et al., 2021)).

Corollary 2. The impact of K_* (FedOne): In the FL framework, $T_{\epsilon}K_*$ represents the total number of queries made to the cloud-based LLM service to achieve an ϵ -solution, whose quantity is directly proportional to the cost incurred for utilizing the LLM. The following condition holds:

$$\Gamma_{\epsilon}K_* \propto K_*$$
 (16)

Therefore, $T_{\epsilon}K_*$ is a function of K_* that increase steadily for $K_* = 1, 2, ..., K$, indicating that the optimal K_* for query efficiency is $K_* = 1$.

4 EXPERIMENT

The objective of the experiment was to assess the performance of Fed-BDPL, along with various aspects of the framework, and to explore the query efficiency advantages of FedOne. The code for this project will be made open source.

298 299

274 275 276

277

278 279

280

285

287 288 289

290

291 292 293

4.1 EXPERIMENT SETUP

For our experiment, we utilized the GLUE benchmark (Wang et al., 2018), which includes a wide range of tasks including MNLI (Williams et al., 2018), QQP (Iyer et al., 2017), SST-2 (Socher et al., 2013), MRPC (Dolan & Brockett, 2005), CoLA (Warstadt et al., 2019), QNLI (Wang et al., 2018), and RTE (Dagan et al., 2005; Haim et al., 2006; Giampiccolo et al., 2007; Bentivogli et al., 2009).

In the baseline experiment within the federated learning framework, we employed 100 clients. In the FedOne framework, there was only one client activated per round for training and aggregation. We adopted the k-shot framework from Perez et al. (2021), adapting it to a federated learning context. Each client received a k-shot dataset comprising k samples per class.

The model architecture employed is RoBERTa-large (Liu et al., 2019). The trainable prompts were placed at different positions in the model depending on the algorithm of the baselines. For the training procedure, we conducted a hyperparameter tuning phase using a grid search approach to explore learning rates of [3e-4, 1e-4, 3e-5, 1e-5]. The batch size was set at 32, the prompt length was set at 20, and the optimization algorithm employed was AdamW (Loshchilov & Hutter, 2017). Further details about the dataset and evaluation metrics of them are available in Appendix B.1.

315

4.2 BASELINES

We evaluated several approaches categorized as standalone and federated learning models. The standalone model baseline includes Manual Prompt Tuning, In-Context Learning (Brown et al., 2020), and Fine-tuning (Diao et al., 2022, Table 7)). In the domain of Federated Prompt Tuning for whitebox scenarios, we adapted two established white-box prompt tuning methods to Federated Learning.
Specifically, we implemented prompt-tuning (Lester et al., 2021) and prefix-tuning v2 (Liu et al., 2021) across distributed clients. In the prompt tuning approach, a prompt was integrated into the embedding layer of the model for each client. Each client then undertook local training solely on their respective prompt. In prefix-tuning v2, the prompt was utilized across all embedding layers

Dataset	MNLI	QQP	SST-2	MRPC	CoLA	QNLI	RTE	Avg.
Manual Prompt In Context Learning FineTuning	$\begin{array}{c} 35.9_{1.3} \\ 37.2_{1.6} \\ 50.8_{1.2} \end{array}$	$\begin{array}{c} 49.8_{0.9} \\ 50.1_{0.9} \\ 60.8_{1.9} \end{array}$	$\begin{array}{c} 77.2_{1.1} \\ 82.8_{2.1} \\ 86.5_{2.0} \end{array}$	$70.4_{1.6} \\72.1_{2.3} \\78.4_{1.3}$	$\begin{array}{c} 0.6_{0.0} \\ 1.1_{0.4} \\ 20.4_{1.9} \end{array}$	$\begin{array}{c} 49.2_{1.1} \\ 50.8_{0.5} \\ 53.2_{1.8} \end{array}$	$\begin{array}{c} 48.2_{0.6} \\ 49.3_{2.3} \\ 55.6_{2.3} \end{array}$	47.33 49.06 57.96
FedOne-PromptTuning FedOne-P-Tuning v2	$\begin{array}{c} 41.5_{0.9} \\ 42.7_{0.7} \end{array}$	$\begin{array}{c} 66.4_{0.2} \\ 66.7_{0.1} \end{array}$	$\begin{array}{c} 77.9_{2.1} \\ 82.9_{0.3} \end{array}$	$\begin{array}{c} 79.5_{0.5} \\ 80.6_{0.1} \end{array}$	$\begin{array}{c} 0.8_{1.1} \\ 1.0_{1.0} \end{array}$	$\begin{array}{c} 49.6_{1.0} \\ 52.4_{0.2} \end{array}$	$53.1_{0.6}$ $56.4_{0.4}$	52.69 54.67
FedOne-BBT FedOne-BDPL FedOne-GS-BDPL	$\begin{array}{c} 41.9_{0.4} \\ 41.0_{1.2} \\ 41.1_{0.4} \end{array}$	$\begin{array}{c} 66.3_{0.2} \\ 66.7_{0.1} \\ 66.9_{0.2} \end{array}$	$76.8_{1.6} \\ 80.8_{6.0} \\ 80.8_{0.4}$	$\begin{array}{c} 80.6_{0.3} \\ 81.1_{0.1} \\ 81.0_{0.1} \end{array}$	$\begin{array}{c} 2.5_{1.3} \\ 5.2_{2.4} \\ 5.3_{1.1} \end{array}$	$51.1_{0.4} \\ 51.7_{1.4} \\ 52.1_{0.8}$	$55.3_{1.0}$ $57.1_{1.9}$ $57.1_{1.1}$	53.50 54.80 54.90

Table 1:	The overall	performance on	the RoBERTa-large.	Each trial	runs across three	random seeds.
		1	0			

337

338

339

340

341

342

343

344 345

346

357

330331332333334

of the model, providing more trainable parameters and enhancing the capability to adapt to downstream tasks. In the federated prompt tuning for the black-box scenario, we adapted the Black-Box Tuning (BBT) (Sun et al., 2022) to FedOne, incorporating projection from a low-dimensional vector and the Covariance Matrix Adaptation Evolution Strategy (CMA-ES) into federated black-box prompt learning. Each client held a distinct low-dimensional vector, while the projection matrix *A* was shared among all clients. For every client, the population size of CMA-ES is set to 20, and the dimension of the low-dimensional vector is set to 500, as recommended by Sun et al. (2022). Finally, we adapted the BDPL (Diao et al., 2022) and Gumbel-Softmax BDPL (GS-BDPL) to the Federated Learning, employing policy gradient methods and Gumbel-Softmax as outlined in Algorithm 1.

4.3 Result

347 Test accuracy The performance results were summarized in Table 1, We observed significant 348 variations in the effectiveness of different learning approaches when applied to the RoBERTa-large 349 model across various NLP tasks. Traditional fine-tuning methods outperformed both Manual Prompt 350 and In-Context Learning techniques, achieving the highest average score of 57.96 across all datasets. 351 This result underscores the effectiveness of complete model retraining over other methods that in-352 volve fewer parameter updates or rely solely on contextual adjustments. Although Manual Prompt 353 performed well on some datasets, its performance lacks stability. White-box federated prompt tun-354 ing methods demonstrate improvements over Manual Prompt and In Context Learning techniques. 355 Generally, the black-box method (BDPL and GS-BDPL) exhibits performance that is comparable to, or slightly better than, the white-box tuning method (PromptTuning and P-Tuning v2). 356

Computational efficiency and resource utilization Table 2 presents the performance metrics for 358 various federated prompt tuning methods per client, focusing on the computational and communica-359 tion efficiencies of these methods. We measured the computation time required for model training. 360 Notably, in the white-box method, training occurs on the client side. We only measure the time for 361 forward and backward propagation of the local model, along with the time for parameter updates. 362 In the black-box prompt learning method, the computation time for model training includes both the 363 execution of the black-box algorithm and the wait time for responses from the cloud-based LLM 364 services. The results are presented in the first column of Table 2. We observed that the white-box method is considerably faster; however, it relies on the assumption that clients possess sufficient 366 computational power, such as GPUs, which may not be practical in FL environments. Within the 367 black-box approach, BDPL requires less computation time for training compared to BBT.

The advantage of black-box prompt learning lies primarily in its reduced communication costs and the efficiency of the trainable parameter size, as well as in avoiding the need to store and train the entire LLM on the client. As illustrated in the table, the federated black-box prompt tuning method features a significantly smaller parameter size and eliminates the GPU requirement on the client. This allows devices with limited computational resources, such as edge devices, to participate in the federated prompt learning process.

374

376

375 4.4 CASE STUDY

Query efficiency and the number of activated clients We began by illustrating the relationship between query efficiency and the number of activated clients (K_*) using a toy example on the

Baseline	Comp. Time for Training (s)	FL Server Comm. Cost (MB)	FL Server # Queries	LLM Server # Query	Client Trainable Parameter Size (MB)	Client Loaded GPU Memory (MB)
FedOne-PromptTuning FedOne-P-Tuning v2	77.6 91.1	$15.63 \\ 1141.72$	100 100	-	7.8×10^{-2} 5.7	3564 (Model, Grad., Promp 3656 (Model, Grad., Promp
FedOne-BBT	20474.6	0.38	100	2×10^{4}	1.9×10^{-3}	1.90×10^{-3} (Prompt)
FedOne-BDPL	1614.3	3.05	100	2×10^{3}	1.5×10^{-2}	1.52×10^{-2} (Prompt)
FedOne-GS-BDPL	1624.4	3.05	100	2×10^{3}	1.5×10^{-2}	1.52×10^{-2} (Prompt)

Table 2: Evaluation of computational efficiency and resource utilization

388

378

379

389 390

391

405

406

407

408

409

410

411

412

413

414

MNIST dataset (LeCun et al., 2010), as shown in Figure 2. We then evaluated the impact of K_* on query efficiency in Federated Black-box Prompt Tuning, as presented in Table 3.

392 In the first experiment, the MNIST dataset is split across 100 clients, with each client initially pos-393 sessing an equal subset of the training data. We experiment with 2 epochs with a learning rate of 394 0.01 and a batch size of 32. We explore five configurations by varying the number of active clients per epoch: 1, 5, 10, 20, and 40 clients (denoted as $K_*=1$, $K_*=5$, $K_*=10$, $K_*=20$, and $K_*=40$ re-395 spectively). The model for each client is a Multilayer Perceptron (MLP). It includes a flattening 396 input layer, a fully connected layer with 512 neurons and ReLU activation, a dropout layer with 0.2397 dropout rate to prevent overfitting and a final fully connected layer that outputs to 10 classes via 398 a Softmax function. As illustrated in Figure 2, utilizing the minimum number of activated clients 399 (FedOne) enables the FL framework to achieve optimal query efficiency for convergence. Although 400 Li et al. (2019) have demonstrated that an increased number of clients participating in the training 401 and aggregation process accelerates convergence in federated learning, in the context of Federated 402 Black-box prompt tuning, the number of queries to the LLM server increases linearly with the num-403 ber of participating clients. This rise in query numbers outpaces the benefits gained from faster 404 convergence due to increased client participation.



Table 3:	Query	Efficiency	in	Federated	Black-Box	Prompt
Learning						

AC	Fe	d-BDPL	F	ed-BBT
	# Epoch	# LLM Queries	# Epoch	# LLM Queries
1	$25.4_{\pm 19.2}$	$528.6_{\pm 382.4}$	$10.8_{\pm 4.4}$	2350.0 _{±870.3}
3	$13.8_{\pm 5.2}$	$889.4_{\pm 310.9}$	$5.4_{\pm 2.3}$	$3825.0_{\pm 1356.2}$
10	$11.5_{\pm 6.8}$	$2430.8_{\pm 1378.8}$	$4.2_{\pm 1.8}$	$10444.4_{\pm 3624.3}$
30	$3.5_{\pm 2.8}$	$2672.7_{\pm 1625.4}$	$1.7_{\pm 0.9}$	$16000.0_{\pm 5656.9}$

416 Figure 2: Toy Example Demonstrate on Query Efficiency 417

418 In the second experiment investigating query efficiency within the federated black-box prompt learn-419 ing scenario, we utilized the SST2 dataset to explore how varying the number of activated clients per 420 round affects model convergence efficiency in FL environments. We tested the number of activated 421 clients per round within the range of [1, 3, 10, 30]. For each configuration, we monitored the number 422 of queries required to achieve target accuracy on the validation dataset. Specifically, we evaluated the prompt at the end of each epoch and stop training once the target accuracy was reached, report-423 ing the number of LLM queries at that point. To ensure reliability and account for variability in 424 the learning process, each experimental setup was replicated 20 times, and outliers in the number of 425 queries were removed. The primary aim of our study was to explore how the number of activated 426 clients affects the speed and efficiency of model convergence in FL, specifically to demonstrate the 427 query efficiency of the FedOne approach. 428

429 The results presented in Table 3 demonstrate a clear trend where fewer activated clients are associated with greater query efficiency. This relationship is evidenced by a consistent decrease in the 430 number of cloud-based LLM queries as the number of activated clients is reduced, a pattern observed 431 across both the Fed-BDPL and Fed-BBT.

Table 4: Federated discrete black-box prompt learning on GPT-3.5 Turbo

	MNLI	QQP	SST-2	MRPC	CoLA	QNLI	RTE	Avg.
No Prompt Prompt w/o. Training	$\begin{array}{c} 14.05_{0.00} \\ 9.17_{1.31} \end{array}$	$\begin{array}{c} 68.04_{0.00} \\ 62.94_{4.94} \end{array}$	$\begin{array}{c} 91.35_{0.00} \\ 79.47_{2.71} \end{array}$	$\begin{array}{c} 79.62_{0.00} \\ 79.88_{1.94} \end{array}$	$\begin{array}{c} 36.01_{0.00} \\ 25.53_{1.10} \end{array}$	$\begin{array}{c} 56.22_{0.00} \\ 56.69_{4.76} \end{array}$	$\begin{array}{c} 72.20_{0.00} \\ 72.16_{2.51} \end{array}$	$59.64 \\ 54.41$
FedOne-BDPL FedOne-GS-BDPL	$\frac{13.67_{0.17}}{\textbf{16.00}_{1.78}}$	68.25 _{3.05} 69.92 _{2.85}	87.77 _{3.82} 92.58 _{2.71}	81.69 _{0.70} 82.93 _{0.67}	32.16 _{6.51} 36.20 _{2.95}	$70.58_{1.35} \\ \textbf{72.05}_{1.64}$	$77.62_{1.62} \\ \textbf{80.33}_{2.32}$	61.67 64.28

441

442

443

444

445

446

447

448

449

450

451

432

Real-world implementation on GPT-3.5 Turbo We implemented the Fed-BDPL framework using GPT-3.5 Turbo, a widely recognized and powerful closed-source language model. We leveraged the OpenAI API (OpenAI, 2024) to enable individual clients to conduct local training. In this implementation of the federated black-box prompt learning, clients sent prompts and input sentences to GPT-3.5, which returns the logarithm of the token probabilities at each position. A key challenge was that OpenAI API only provides probabilities for the top 20 tokens for each position. Consequently, we needed to transform the predictions on these tokens into the categorical prediction of the input sentence, instead of using straightforward model output as we did in the RoBERTa experiment. To solve this problem, we appended a template question at the end of the input sentence to query the target label token. For example, in the QQP task, we added the phrase "equivalent? yes or no" to the end of the input sentence. This allowed us to retrieve the top probabilities for all class label tokens ("yes", "no") in the top 20 probability. and use the probability of the target token as the logit output, for the following procedure.

The results were presented in table 4, indicating that GPT-3.5 Turbo achieves a certain level of performance without any prompts. When adding a random prompt without training, the model's performance dropped. After tuning the prompts, performance improved significantly, surpassing the performance without a prompt. Furthermore, the table demonstrated that the GS-BDPL method consistently outperforms other black-box approaches. To summary, this method can be used to perform prompt tuning on federated learning with extremely low computational resource requirements.

458 459 460

461

462

5 RELATED WORKS

5.1 WHITE-BOX AND BLACK-BOX PROMPT TUNING

463 Prompt tuning is a technique for adapting LLM to downstream tasks. It tailors the model's responses 464 to specific tasks or styles without requiring the retraining of the entire model. In white-box prompt 465 tuning, the learner is granted full access to the LLM, allowing them to modify and access intermediate results of the model and acquire the gradient. Li & Liang (2021) and Lester et al. (2021) 466 proposed a lightweight and modular alternative to full model fine-tuning for natural language gen-467 eration tasks, which optimizes a sequence of continuous soft prompts, prepend in the embedding 468 layers of the LLM. Liu et al. (2021) propose the P-tuning v2. Instead of only applying the prompt 469 in the input layer in Li & Liang (2021), they adapt trainable parameters on all layers' inputs, which 470 can effectively match the performance of fine-tuning across a wide range of models. 471

In situations where the learner cannot access the intermediate result of the LLM model, the learner 472 has to use a black-box prompt tuning method. In the black-box prompt tuning, the learner can only 473 query the output of the LLM with the input of the model. Most of the research assumes the input 474 is at the soft prompt layer of the LLM (Sun et al., 2022; Chen et al., 2023). Others research use the 475 discrete prompt which is concrete with the input text, which is more portable, and is usable for any 476 cloud-based LLM API (Diao et al., 2022). Xiao et al. (2023) presents a privacy-preserving, efficient 477 transfer learning method that adapts large foundation models to specific tasks. This method does 478 not require access to the full model or compromise data privacy. It utilizes a lightweight adapter 479 and a compressed emulator for model tuning. Chen et al. (2023) introduces an efficient method 480 for optimizing instructions for black-box LLMs using Bayesian optimization of soft prompts. This 481 approach significantly improves LLM performance across a variety of tasks without requiring direct 482 access to the model's internals. Deng et al. (2022) proposes an efficient method to optimize discrete 483 text prompts using reinforcement learning, demonstrating superior performance compared to other prompt optimization techniques across a variety of tasks. Sun et al. (2022) proposed black-box tun-484 ing (BBT), a method that optimizes continuous prompts by optimizing a lower-dimensional vector 485 and projecting them to the prompt searching space. They use the Covariance Matrix Adaptation

Evolution Strategy (CMA-ES) for optimizing the vector. Diao et al. (2022) introduced BDPL which optimizes the discrete prompt with the policy gradient method.

The black-box prompt tuning method is versatile and applicable to various tasks and models without model-specific modifications. However, its major drawback is computational inefficiency, requiring multiple forward passes through the model, which leads to high costs and extended training times. Additionally, the convergence of the derivation-free optimization method is often slow.

493 494

495

5.2 FEDERATED LEARNING

Federated learning (McMahan et al., 2017; Karimireddy et al., 2020; Li et al., 2020; 2021; Marfoq
et al., 2022; Mishchenko et al., 2022), first introduced by McMahan et al. (2017), is a paradigm that
enables devices to collaboratively train a shared predictive model by locally aggregating updates. In
this framework, each client maintains a copy of the model for local training, and the server selects
a subset of clients in each round for aggregation. Since the introduction of FedAvg, it has lacked
formal theoretical convergence guarantees. As a result, researchers have made significant efforts to
establish and demonstrate its convergence (Zhou & Cong, 2017; Stich, 2018).

Partial client activation is a key area of research in FL and has gained significant attention due to its impact on improving convergence rates and system efficiency. Stich (2018) shows that in the 504 convex case, increasing the number of activated clients significantly improves convergence rates in 505 Federated Learning with independent and identically distributed (IID) data, achieving a linear speed-506 up. Li et al. (2019) extended the understanding of federated learning by analyzing the convergence 507 of the FedAvg algorithm, under the convex case. They demonstrate that under the non-IID setting, 508 the convergence rate has a weak dependence on the number of activated clients, which implies that 509 the FedAvg is not able to achieve linear speedup under this case, therefore the participation ratio can 510 be set smaller to alleviate the straggler effect without affecting the convergence rate.

511 512

5.3 FEDERATED PROMPT TUNING

513 514

Applying federated learning to prompt tuning can enhance the model by incorporating additional 515 data. This approach leverages distributed datasets to improve model performance while adhering 516 to data privacy. To apply white-box prompt tuning in FL, each client maintains the entire model 517 but trains only the prompt parameters. These parameters are then shared and aggregated across 518 clients (Zhao et al., 2023; Zhang et al., 2023; Che et al., 2023). However, this method assumes 519 white-box access to LLM, which is impractical for closed-source LLMs. Consequently, black-box 520 prompt learning has been adapted for FL to address these limitations (Lin et al., 2023; Sun et al., 521 2023). Lin et al. (2023), applied the black-box prompt tuning method Diao et al. (2022) to the 522 Federated Learning, where the client can train the probability matrix for the discrete prompt via 523 querying the cloud-based LLM API. Sun et al. (2023) applied the BBT to FL. In this approach, clients train only low-dimensional vectors using CMA-ES via querying the cloud-based LLM API. 524

For federated prompt tuning, most existing research has focused on conventional issues in FL, such as data heterogeneity (Zhao et al., 2023), privacy (Zhang et al., 2023), security (Zhao et al., 2023), and client computation-communication efficiency (Lin et al., 2023; Sun et al., 2023). However, no studies have yet addressed the query efficiency of Federated Black-box Prompt Tuning, a novel challenge introduced by the deployment of Black-box Prompt Tuning through cloud-based APIs.

530 531 532

533

6 CONCLUSION

We identified that previous research on federated black-box prompt tuning had overlooked the significant query costs associated with cloud-based LLM services. To address this issue, we conducted a theoretical analysis of query efficiency in the context of federated black-box prompt tuning, revealing the relationship between query efficiency and the number of activated clients. Based on our findings, we propose the FedOne framework, which achieves optimal query efficiency with respect to the number of activated clients. We performed numerical experiments on various aspects of FedOne, further validating its performance through real-world experiments on GPT-3.5.

540	REFERENCES
541	

565

566

567

568

569

571

572

577

578

579

580

586

Zeyuan Allen-Zhu. How to make the gradients small stochastically: Even faster convex and non-542 convex sgd. Advances in Neural Information Processing Systems, 31, 2018. 543

- 544 Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. The fifth pascal recognizing textual entailment challenge. In TAC, 2009. 546
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, 547 548 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020. 549
- 550 Tianshi Che, Ji Liu, Yang Zhou, Jiaxiang Ren, Jiwen Zhou, Victor S Sheng, Huaiyu Dai, and Dejing 551 Dou. Federated learning of large language models with parameter-efficient prompt tuning and 552 adaptive optimization. arXiv preprint arXiv:2310.15080, 2023. 553
- Lichang Chen, Jiuhai Chen, Tom Goldstein, Heng Huang, and Tianyi Zhou. Instructzero: Efficient 554 instruction optimization for black-box large language models. arXiv preprint arXiv:2306.03082, 555 2023. 556
- Ido Dagan, Oren Glickman, and Bernardo Magnini. The pascal recognising textual entailment 558 challenge. In Machine Learning Challenges Workshop, pp. 177–190. Springer, 2005. 559
- Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng Song, 560 Eric P Xing, and Zhiting Hu. Rlprompt: Optimizing discrete text prompts with reinforcement 561 learning. arXiv preprint arXiv:2205.12548, 2022. 562
 - Shizhe Diao, Zhichao Huang, Ruijia Xu, Xuechun Li, Yong Lin, Xiao Zhou, and Tong Zhang. Black-box prompt learning for pre-trained language models. arXiv preprint arXiv:2201.08531, 2022.
 - William B. Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In Proceedings of the Third International Workshop on Paraphrasing (IWP2005), 2005. URL https://aclanthology.org/I05-5002.
- 570 Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. SIAM journal on optimization, 23(4):2341–2368, 2013.
- Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. The third PASCAL recog-573 nizing textual entailment challenge. In Proceedings of the ACL-PASCAL Workshop on Textual 574 Entailment and Paraphrasing, pp. 1-9, Prague, 2007. Association for Computational Linguistics. 575 URL https://aclanthology.org/W07-1401. 576
 - Robert Gower, Othmane Sebbouh, and Nicolas Loizou. Sgd for structured nonconvex functions: Learning rates, minibatching and interpolation. In International Conference on Artificial Intelligence and Statistics, pp. 1315–1323. PMLR, 2021.
- Farzin Haddadpour and Mehrdad Mahdavi. On the convergence of local descent methods in feder-581 ated learning. arXiv preprint arXiv:1910.14425, 2019. 582
- 583 R Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan 584 Szpektor. The second pascal recognising textual entailment challenge. In Proceedings of the 585 Second PASCAL Challenges Workshop on Recognising Textual Entailment, volume 7, 2006.
- Elad Hazan and Satyen Kale. Beyond the regret minimization barrier: optimal algorithms for 587 stochastic strongly-convex optimization. The Journal of Machine Learning Research, 15(1): 588 2489-2512, 2014. 589
- Shankar Iyer, Nikhil Dandekar, Kornél Csernai, et al. First quora dataset release: Question pairs. data. quora. com, 2017. 592
- Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. arXiv preprint arXiv:1611.01144, 2016.

594 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and 595 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In 596 International Conference on Machine Learning, pp. 5132–5143. PMLR, 2020. 597 Ahmed Khaled and Peter Richtárik. Better theory for sgd in the nonconvex world. ArXiv, 598 abs/2002.03329, 2020. URL https://api.semanticscholar.org/CorpusID: 211069380. 600 601 Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. ATT Labs [Online]. Available: http://yann.lecun.com/exdb/mnist, 2, 2010. 602 603 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt 604 tuning. arXiv preprint arXiv:2104.08691, 2021. 605 606 Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. Proceedings of Machine Learning and Sys-607 tems, 2:429-450, 2020. 608 609 Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of 610 FedAvg on non-iid data. arXiv preprint arXiv:1907.02189, 2019. 611 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. arXiv 612 preprint arXiv:2101.00190, 2021. 613 614 Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou. FedBN: Federated learn-615 ing on non-IID features via local batch normalization. In International Conference on Learning 616 Representations, 2021. URL https://openreview.net/pdf?id=6YEQUn0QICG. 617 Zihao Lin, Yan Sun, Yifan Shi, Xueqian Wang, Lifu Huang, Li Shen, and Dacheng Tao. 618 Efficient federated prompt tuning for black-box large pre-trained models. arXiv preprint 619 arXiv:2310.03123, 2023. 620 621 Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-622 tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks. arXiv preprint arXiv:2110.07602, 2021. 623 624 Yanli Liu, Yuan Gao, and Wotao Yin. An improved analysis of stochastic gradient descent with 625 momentum. Advances in Neural Information Processing Systems, 33:18261–18271, 2020. 626 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike 627 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT pretraining 628 approach. arXiv preprint arXiv:1907.11692, 2019. 629 630 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. arXiv preprint 631 arXiv:1711.05101, 2017. 632 Othmane Marfoq, Giovanni Neglia, Richard Vidal, and Laetitia Kameni. Personalized federated 633 learning through local memorization. In International Conference on Machine Learning, pp. 634 15070-15092. PMLR, 2022. 635 636 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 637 Communication-efficient learning of deep networks from decentralized data. In Artificial intelligence and statistics, pp. 1273–1282. PMLR, 2017. 638 639 Konstantin Mishchenko, Grigory Malinovsky, Sebastian Stich, and Peter Richtárik. Proxskip: Yes! 640 local gradient steps provably lead to communication acceleration! finally! In International Con-641 ference on Machine Learning, pp. 15750–15769. PMLR, 2022. 642 Yurii Nesterov et al. *Lectures on convex optimization*, volume 137. Springer, 2018. 643 644 OpenAI. Openai api, 2024. URL https://openai.com/index/openai-api/. Accessed: 645 2024-09-30. 646 Ethan Perez, Douwe Kiela, and Kyunghyun Cho. True few-shot learning with language models. 647

Advances in neural information processing systems, 34:11054–11070, 2021.

648 Sashank J. Reddi, Ahmed Hefny, Suvrit Sra, Barnabas Poczos, and Alex Smola. Stochastic variance 649 reduction for nonconvex optimization. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), 650 Proceedings of The 33rd International Conference on Machine Learning, volume 48 of Proceed-651 ings of Machine Learning Research, pp. 314–323, New York, New York, USA, 20–22 Jun 2016. 652 PMLR. 653 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, 654 and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment 655 treebank. In Proceedings of the 2013 Conference on Empirical Methods in Natural Language 656 *Processing*, pp. 1631–1642, Seattle, Washington, USA, 2013. Association for Computational Lin-657 guistics. URL https://aclanthology.org/D13-1170. 658 Sebastian U Stich. Local sgd converges fast and communicates little. arXiv preprint 659 arXiv:1805.09767, 2018. 660 661 Jingwei Sun, Ziyue Xu, Hongxu Yin, Dong Yang, Daguang Xu, Yiran Chen, and Holger R Roth. 662 Fedbpt: Efficient federated black-box prompt tuning for large language models. arXiv preprint 663 arXiv:2310.01467, 2023. 664 Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. Black-box tuning 665 for language-model-as-a-service. In International Conference on Machine Learning, pp. 20841– 666 20855. PMLR, 2022. 667 668 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 669 Glue: A multi-task benchmark and analysis platform for natural language understanding. arXiv preprint arXiv:1804.07461, 2018. 670 671 Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. Neural network acceptability judgments. 672 Transactions of the Association for Computational Linguistics, 7:625–641, 2019. doi: 10.1162/ 673 tacl_a_00290. URL https://aclanthology.org/Q19-1040. 674 Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sen-675 tence understanding through inference. In Proceedings of the 2018 Conference of the North Amer-676 ican Chapter of the Association for Computational Linguistics: Human Language Technologies, 677 Volume 1 (Long Papers), pp. 1112–1122, New Orleans, Louisiana, 2018. Association for Com-678 putational Linguistics. doi: 10.18653/v1/N18-1101. URL https://aclanthology.org/ 679 N18-1101. 680 681 Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine learning, 8:229-256, 1992. 682 683 Guangxuan Xiao, Ji Lin, and Song Han. Offsite-tuning: Transfer learning without full model. arXiv 684 preprint arXiv:2302.04870, 2023. 685 686 Yi Xu, Rong Jin, and Tianbao Yang. Non-asymptotic analysis of stochastic methods for non-smooth non-convex regularized problems. Advances in Neural Information Processing Systems, 32, 2019. 687 688 Dong Yin, Ashwin Pananjady, Max Lam, Dimitris Papailiopoulos, Kannan Ramchandran, and Peter 689 Bartlett. Gradient diversity: a key ingredient for scalable distributed learning. In Amos Storkey 690 and Fernando Perez-Cruz (eds.), Proceedings of the Twenty-First International Conference on Ar-691 tificial Intelligence and Statistics, volume 84 of Proceedings of Machine Learning Research, pp. 692 1998-2007. PMLR, 09-11 Apr 2018. URL https://proceedings.mlr.press/v84/ 693 vin18a.html. 694 Zhuo Zhang, Yuanhang Yang, Yong Dai, Qifan Wang, Yue Yu, Lizhen Qu, and Zenglin Xu. Fed-695 PETuning: When federated learning meets the parameter-efficient tuning methods of pre-trained 696 language models. In Annual Meeting of the Association of Computational Linguistics 2023, pp. 697 9963–9977. Association for Computational Linguistics (ACL), 2023. 698 Haodong Zhao, Wei Du, Fangqi Li, Peixuan Li, and Gongshen Liu. Fedprompt: Communication-699 efficient and privacy-preserving prompt tuning in federated learning. In ICASSP 2023-2023 IEEE 700 International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1–5. IEEE, 701 2023.

Fan Zhou and Guojing Cong. On the convergence properties of a *k*-step averaging stochastic gradient descent algorithm for nonconvex optimization. *arXiv preprint arXiv:1708.01012*, 2017.

A CONVERGENCE ANALYSIS

 A.1 NOTATION AND OMITTED MATHEMATICAL STEPS

Number of global iterations
Number of local iterations
Number of client-server interactions
The default 2-norm in this paper
Categorical distribution function
Abbreviation for $j_i \sim \text{Cat}(\boldsymbol{p}_i)$
Dataset for client k
Input sentence
The category corresponding to the input sentence
The number of all clients
Probability of each client being selected
The number of selected clients.
The set of selected clients.
Gumbel-Softmax parameters for k-th client
Average of $\boldsymbol{\alpha}^k$
Prompt
Vocabulary list of total N tokens choice
The probability distribution over the N token indexes.
The local full gradient of the client k (19)
The local stochastic mini-batch gradient of the client k (18)
The local stochastic mini-batch variance-reduced policy gradient of the client k (6)
The average full gradient of the client group U_t (33)
The average stochastic mini-batch gradient of the client group U_t (34)
The average stochastic mini-batch variance-reduced policy gradient of the client group U_t (35)
The average full gradient of all K clients for $\{\boldsymbol{\alpha}^k\}_{k=1}^K$ (32)
The average full gradient of all K clients for α (31)
_

Table 5: Notation Table

For i = 1, ..., n, we define the stochastic gradient, stochastic mini-batch gradient and full gradient with respect to $\boldsymbol{\alpha}_i^k$ as following:

$$\nabla_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}^{k}, \psi_{m}^{k}) \stackrel{def}{=} \nabla_{\boldsymbol{\alpha}_{i}^{k}} \mathbb{E}_{\Phi^{k} \sim \mathrm{GS}(\boldsymbol{\alpha}^{k})} \left[\mathcal{L}(\Phi^{k}, \psi_{m}^{k}) \right]$$
(17)

 $\nabla_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}^{k}, \mathcal{B}^{k}) \stackrel{def}{=} \nabla_{\boldsymbol{\alpha}_{i}^{k}} \frac{1}{B^{k}} \sum_{\boldsymbol{\psi}_{m}^{k} \in \mathcal{B}^{k}} \mathbb{E}_{\Phi^{k} \sim \mathrm{GS}(\boldsymbol{\alpha}^{k})} [\mathcal{L}(\Phi^{k}, \boldsymbol{\psi}_{m}^{k})]$ (18)

755 $\nabla_{\boldsymbol{\alpha}_{i}^{k}} F^{k}(\boldsymbol{\alpha}^{k}, \Psi^{k}) \stackrel{def}{=} \nabla_{\boldsymbol{\alpha}_{i}^{k}} \frac{1}{M^{k}} \sum_{\psi_{m}^{k} \in \Psi^{k}} \mathbb{E}_{\Phi^{k} \sim GS(\boldsymbol{\alpha}^{k})} [\mathcal{L}(\Phi^{k}, \psi_{m}^{k})]$ (19)

THE OMITTED DERIVATIVE PROCESS FOR EQ. 5 A.2

 $\frac{\partial \log p_{i,j_i^k}^k}{\partial \alpha_{i,j}^k} = \frac{\partial}{\partial \alpha_{i,j}^k} \left(\log \left(\frac{\exp\left(\frac{\log(\alpha_{i,j_k}^k) + g_{i,j_k}^k}{\tau}\right)}{\sum_{l=1}^N \exp\left(\frac{\log(\alpha_{i,l}^k) + g_{i,l}^k}{\tau}\right)} \right) \right)$

ð

According to the derivation rule of the Softmax function, when $j = j_i^k$:

$$\frac{\partial p_{i,j_i^k}^k}{\partial \alpha_{i,j_i^k}^k} = \frac{1}{\tau \alpha_{i,j_i^k}^k p_{i,j_i^k}^k} \cdot \left(1 - p_{i,j_i^k}^k\right) \cdot p_{i,j_i^k}^k = \frac{\left(1 - p_{i,j_i^k}^k\right)}{\tau \alpha_{i,j_i^k}^k}$$
(21)

 $\left(\frac{\frac{1}{\tau}}{\sum_{l=1}^{N} \exp\left(\frac{\log(\alpha_{i,l}^{k}) + g_{i,l}^{k}}{\tau}\right)}\right)$

 $\frac{\left\lfloor \frac{\lambda}{\sum_{l=1}^{N} \exp\left(\frac{\log(\alpha_{i,l}^{k}) + g_{i,l}^{k}}{\tau}\right)}\right\rfloor}{\partial\left(\frac{\log(\alpha_{i,j}^{k}) + g_{i,j}^{k}}{\tau}\right)} \cdot \frac{\partial\left(\frac{\log(\alpha_{i,j}^{k}) + g_{i,j}^{k}}{\tau}\right)}{\partial\alpha_{i,j}^{k}}$

when $j \neq j_i^k$:

$$\frac{\partial p_{i,j_i^k}^k}{\partial \alpha_{i,j}^k} = \frac{1}{\tau \alpha_{i,j}^k p_{i,j_i^k}^k} \cdot \left(-p_{i,j_i^k}^k \cdot p_{i,j}^k\right) = \frac{-p_{i,j}^k}{\tau \alpha_{i,j}^k}$$
(22)

(20)

A.3 LEMMAS

The following lemma shows that the unbiasedness and bounded variance of variance-reduced policy gradient. This is important for bounding the randomness introduced by prompt sampling. Lemma 1. Unbiasedness and bounded variance of variance-reduced policy gradient: At the k-th client, $\forall \psi_m^k \in \Psi^k$, $r = 1, \cdots, I$ denotes the r-th sampling of Φ^k from \boldsymbol{p}^k w.r.t. $\{\Phi^{k,r}\}_{r=1}^I \sim \boldsymbol{p}^k$, $p^k = GS(\boldsymbol{\alpha}^k), \ \alpha_{i,j} \ge \nu > 0 \text{ for } i = 1, ..., n \text{ and } j = 1, ..., N, \ \tau > 0 \text{ is the temperature parameter,}$ and $\sigma_{\alpha}^2 = \frac{8G^2N}{\tau^2\nu^2}$, then the variance-reduced policy gradient is unbiased and bounded by :

$$\mathbb{E}_{\{\Phi^{k,r}\sim GS(\boldsymbol{\alpha}^k)\}_{r=1}^{I}}\left[\hat{\nabla}_{\boldsymbol{\alpha}_i^k}f^k(\boldsymbol{\alpha}^k,\psi_m^k)\right] = \nabla_{\boldsymbol{\alpha}_i^k}f^k(\boldsymbol{\alpha}^k,\psi_m^k)$$
(23)

$$\mathbb{E}_{\{\Phi^{k,r}\sim GS(\boldsymbol{\alpha}^{k})\}_{r=1}^{I}} \left[\left\| \hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}^{k},\psi_{m}^{k}) - \mathbb{E}_{\{\Phi^{k,r}\sim GS(\boldsymbol{\alpha}^{k})\}_{r=1}^{I}} \left[\hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}^{k},\psi_{m}^{k}) \right] \right\|^{2} \right] \leq \frac{\sigma_{\alpha}^{2}}{I^{2}} \quad (24)$$

Proof. We abbreviate $\mathbb{E}_{\{\Phi^{k,r} \sim GS(\boldsymbol{\alpha}^k)\}_{r=1}^I}$ as $\mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^I}$. 1) The unbiasedness of variance-reduced policy gradient:

l

$$\mathbb{E}_{\{\Phi^{k,r}\sim GS(\pmb{a}^{k})\}_{r=1}^{I}} \left[\hat{\nabla}_{\pmb{a}_{i}^{k}} f^{k}(\pmb{a}^{k},\psi_{m}^{k}) \right] \\ = \mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}} \left\{ \frac{1}{I-1} \sum_{r=1}^{I} \left[\left(\mathcal{L}(\Phi^{k,r},\psi_{m}^{k}) - \frac{1}{I} \sum_{w=1}^{I} \mathcal{L}(\Phi^{k,w},\psi_{m}^{k}) \right) \nabla_{\pmb{a}_{i}^{k}} \log P(\phi_{i}^{k,r}) \right] \right\} \\ = \mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}} \left\{ \frac{1}{I-1} \sum_{r=1}^{I} \left[\left(\frac{I-1}{I} \cdot \mathcal{L}(\Phi^{k,r},\psi_{m}^{k}) - \frac{1}{I} \sum_{w=1}^{I} \mathcal{L}(\Phi^{k,w},\psi_{m}^{k}) \right) \nabla_{\pmb{a}_{i}^{k}} \log P(\phi_{i}^{k,r}) \right] \right\}$$

$$\begin{split} &= \mathbb{E}_{\left\{\Phi^{k,r}\right\}_{i=1}^{L}} \left[\frac{1}{I} \sum_{r=1}^{I} \left(\mathcal{L}(\Phi^{k,r}, \psi_{m}^{k}) \cdot \nabla_{\mathbf{a}_{i}^{k}} \log P(\phi_{i}^{k,r}) \right) \right] \\ &\quad - \mathbb{E}_{\left\{\Phi^{k,r}\right\}_{i=1}^{L}} \left\{ \frac{1}{I} \sum_{r=1}^{I} \left[\left(\frac{1}{I-1} \sum_{w \neq r}^{I} \mathcal{L}(\Phi^{k,w}, \psi_{w}^{k}) \right) \nabla_{\mathbf{a}_{i}^{k}} \log P(\phi_{i}^{k,r}) \right] \right\} \\ &\quad \left(\stackrel{\text{d}}{=} \frac{1}{I} \sum_{r=1}^{I} \mathbb{E}_{\Phi^{k,r}} \left[\mathcal{L}(\Phi^{k,r}, \psi_{m}^{k}) \cdot \nabla_{\mathbf{a}_{i}^{k}} \log P(\phi_{i}^{k,r}) \right] \\ &\quad - \frac{1}{I} \sum_{r=1}^{I} \left[\left(\frac{1}{I-1} \sum_{w \neq r}^{I} \mathbb{E}_{\Phi^{k,r}} \mathcal{L}(\Phi^{k,w}, \psi_{m}^{k}) \right) \mathbb{E}_{\Phi^{k}} \nabla_{\mathbf{a}_{i}^{k}} \log P(\phi_{i}^{k,r}) \right] \\ &\quad - \frac{1}{I} \sum_{r=1}^{I} \left[\left(\frac{1}{I-1} \sum_{w \neq r}^{I} \mathbb{E}_{\Phi^{w}} \mathcal{L}(\Phi^{k,w}, \psi_{m}^{k}) \right) \mathbb{E}_{\Phi^{k}} \left[\nabla_{\mathbf{a}_{i}^{k}} \log P(\phi_{i}^{k,r}) \right] \right] \\ &\quad - \frac{1}{I} \sum_{r=1}^{I} \left[\left(\frac{1}{I-1} \sum_{w \neq r}^{I} \mathbb{E}_{\Phi^{w}} \left[\mathcal{L}(\Phi^{k,w}, \psi_{m}^{k}) \right] \right) \mathbb{E}_{\Phi^{k}} \left[\nabla_{\mathbf{a}_{i}^{k}} \log P(\phi_{i}^{k,r}) \right] \right] \\ &\quad = \nabla_{\mathbf{a}_{i}^{k}} f^{k} (\mathbf{a}^{k}, \psi_{m}^{k}) - \frac{1}{I} \sum_{r=1}^{I} \left(\frac{1}{I-1} \sum_{w \neq r}^{I} \mathbb{E}_{\Phi^{w}} \left[\mathcal{L}(\Phi^{k,w}, \psi_{m}^{k}) \right] \right) \cdot \mathbb{E}_{\Phi^{k}} \left(\frac{\nabla_{\mathbf{a}_{i}^{k}} P(\phi_{i}^{k,r}) \right) \right] \\ &\quad = \nabla_{\mathbf{a}_{i}^{k}} f^{k} (\mathbf{a}^{k}, \psi_{m}^{k}) - \frac{1}{I} \sum_{r=1}^{I} \left(\frac{1}{I-1} \sum_{w \neq r}^{I} \mathbb{E}_{\Phi^{w}} \left[\mathcal{L}(\Phi^{k,w}, \psi_{m}^{k}) \right] \right) \cdot \sum_{\phi_{i}^{k,r} \sim \mathbf{p}_{i}^{k}} \left(\nabla_{\mathbf{a}_{i}^{k}} P(\phi_{i}^{k,r}) \right) \right] \\ &\quad = \nabla_{\mathbf{a}_{i}^{k}} f^{k} (\mathbf{a}^{k}, \psi_{m}^{k}) - \frac{1}{I} \sum_{r=1}^{I} \left(\frac{1}{I-1} \sum_{w \neq r}^{I} \mathbb{E}_{\Phi^{w}} \left[\mathcal{L}(\Phi^{k,w}, \psi_{m}^{k}) \right] \right) \cdot \sum_{\phi_{i}^{k,r} \sim \mathbf{p}_{i}^{k}} \left(\nabla_{\mathbf{a}_{i}^{k}} P(\phi_{i}^{k,r}) \right) \right] \\ &\quad = \nabla_{\mathbf{a}_{i}^{k}} f^{k} (\mathbf{a}^{k}, \psi_{m}^{k}) - \frac{1}{I} \sum_{r=1}^{I} \left(\frac{1}{I-1} \sum_{w \neq r}^{I} \mathbb{E}_{\Phi^{w}} \left[\mathcal{L}(\Phi^{k,w}, \psi_{m}^{k}) \right] \right) \cdot \sum_{\phi_{i}^{k,r} \sim \mathbf{p}_{i}^{k}} \left(\nabla_{\mathbf{a}_{i}^{k}} f^{k} (\mathbf{a}^{k}, \psi_{m}^{k}) - \frac{1}{I} \sum_{r=1}^{I} \left(\frac{1}{I-1} \sum_{w \neq r}^{I} \mathbb{E}_{\Phi^{w}} \left[\mathcal{L}(\Phi^{k,w}, \psi_{m}^{k}) \right] \right) \cdot \nabla_{\mathbf{a}_{i}^{k}} (\mathbf{1} \right) \\ \begin{pmatrix} (\mathbb{Q} \nabla_{\mathbf{a}_{i}^{k}} f^{k} (\mathbf{a}^{k}, \psi_{m}^{k}) - \frac{1}{I} \sum_{r=1}^{I} \left(\frac{1}{I-1} \sum_{w \neq r}^{I} \mathbb{E}_{\Phi^{w}} \left[\mathcal{L}(\Phi^{k,w}, \psi_{m}^{k}) \right] \right) \cdot \nabla_{\mathbf{a}_{i}^{k}} (\mathbf$$

$$= \nabla_{\boldsymbol{\alpha}_i^k} f^k(\boldsymbol{\alpha}^k, \psi_m^k)$$

where (1) uses independence of each sampling for Φ^k ; (2) is because *n* is not infinite and the GS function is continuous and derivable with respect to α_i ; (3) uses the property that the elements of probability vector sum to 1.

2) The bounded variance of variance-reduced policy gradient:

 $= \mathbb{E}_{\{\Phi^{k,r} \sim \mathrm{GS}(\boldsymbol{\alpha}^{k})\}_{r=1}^{I}} \left[\left\| \hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}^{k},\psi_{m}^{k}) - \nabla_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}^{k},\psi_{m}^{k}) \right\|^{2} \right]$

 $\operatorname{Var}_{\{\Phi^{k,r} \sim \operatorname{GS}(\boldsymbol{\alpha}^{k})\}_{r=1}^{I}} \left[\hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}^{k}, \psi_{m}^{k}) \right]$

$$= \mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}} \left\{ \left\| \frac{1}{I-1} \sum_{r=1}^{I} \left[\left(\mathcal{L}(\Phi^{k,r}, \psi_{m}^{k}) - \frac{1}{I} \sum_{w=1}^{I} \mathcal{L}(\Phi^{k,w}, \psi_{m}^{k}) \right) \nabla_{\mathbf{a}_{i}^{k}} \log P(\phi_{i}^{k,r}) \right] - \nabla_{\mathbf{a}_{i}^{k}} f^{k}(\mathbf{a}^{k}, \psi_{m}^{k}) \right\|^{2} \right\}$$

$$= \mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}} \left\{ \left\| \frac{1}{I} \sum_{r=1}^{I} \left[\frac{1}{I-1} \sum_{\substack{w=1\\w \neq r}}^{I} \left(\mathcal{L}(\Phi^{k,r}, \psi_{m}^{k}) - \mathcal{L}(\Phi^{k,w}, \psi_{m}^{k}) \right) \nabla_{\mathbf{a}_{i}^{k}} \log P(\phi_{i}^{k,r}) - \nabla_{\mathbf{a}_{i}^{k}} f^{k}(\mathbf{a}^{k}, \psi_{m}^{k}) \right] \right\|^{2} \right\}$$

$$\stackrel{(1)}{=} \frac{1}{I^2} \sum_{r=1}^{I} \mathbb{E}_{\Phi^{k,r}} \left[\left\| \frac{1}{I-1} \sum_{\substack{w=1\\w \neq r}}^{I} \left(\mathcal{L}(\Phi^{k,r}, \psi_m^k) - \mathcal{L}(\Phi^{k,w}, \psi_m^k) \right) \nabla_{\boldsymbol{\alpha}_i^k} \log P(\phi_i^{k,r}) - \nabla_{\boldsymbol{\alpha}_i^k} f^k(\boldsymbol{\alpha}^k, \psi_m^k) \right\| \right]$$

$$\stackrel{(2)}{\leq} \frac{1}{I^{2}(I-1)^{2}} \sum_{r=1}^{I} \mathbb{E}_{\Phi^{k,r}} \left[\left\| \sum_{\substack{w=1\\w\neq r}}^{I} \left(\mathcal{L}(\Phi^{k,r},\psi_{m}^{k}) - \mathcal{L}(\Phi^{k,w},\psi_{m}^{k}) \right) \nabla_{\boldsymbol{\alpha}_{i}^{k}} \log P(\phi_{i}^{k,r}) \right\| \right]$$

$$\stackrel{(3)}{\leq} \frac{4G^{2}}{I^{2}(I-1)^{2}} \sum_{r=1}^{I} \mathbb{E}_{\Phi^{k,r}} \left[\left\| \sum_{\substack{w=1\\w=1}}^{I} \nabla_{\boldsymbol{\alpha}_{i}^{k}} \log P(\phi_{i}^{k,r}) \right\|^{2} \right] \stackrel{(4)}{\leq} \frac{4G^{2}N}{I(I-1)\tau^{2}\nu^{2}} \stackrel{(5)}{\leq} \frac{8G^{2}N}{I^{2}\tau^{2}\nu^{2}}$$

 $w \neq r$

where (1) uses

$$\begin{split} & \mathbb{E}_{\left\{\Phi^{k,r}\right\}_{r=1}^{I}} \left[\frac{1}{I-1} \sum_{\substack{w=1\\w \neq r}}^{I} \left(\mathcal{L}(\Phi^{k,r},\psi_{m}^{k}) - \mathcal{L}(\Phi^{k,w},\psi_{m}^{k}) \right) \nabla_{\mathbf{a}_{i}^{k}} \log P(\phi_{i}^{k,r}) - \nabla_{\mathbf{a}_{i}^{k}} f^{k}(\mathbf{a}^{k},\psi_{m}^{k}) \right] \right] \\ &= \frac{1}{I-1} \sum_{\substack{w=1\\w \neq r}}^{I} \mathbb{E}_{\Phi^{k,r},\Phi^{k,w},w \neq r} \left[\left(\mathcal{L}(\Phi^{k,r},\psi_{m}^{k}) - \mathcal{L}(\Phi^{k,w},\psi_{m}^{k}) \right) \nabla_{\mathbf{a}_{i}^{k}} \log P(\phi_{i}^{k,r}) \right] - \nabla_{\mathbf{a}_{i}^{k}} f^{k}(\mathbf{a}^{k},\psi_{m}^{k}) \\ &= \frac{1}{I-1} \sum_{\substack{w=1\\w \neq r}}^{I} \mathbb{E}_{\Phi^{k,r}} \left[\mathcal{L}(\Phi^{k,r},\psi_{m}^{k}) \nabla_{\mathbf{a}_{i}^{k}} \log P(\phi_{i}^{k,r}) \right] - \nabla_{\mathbf{a}_{i}^{k}} f^{k}(\mathbf{a}^{k},\psi_{m}^{k}) \\ &= 0 \end{split}$$

and the independence of each sampling for Φ^k . (2) uses inequality $\mathbb{E} \|a - \mathbb{E}a\|^2 \le \mathbb{E} \|a\|^2$; (3) uses Assumption 2; (4) uses $\alpha_{i,j}^{k,r} \ge \nu > 0$ and (5):

$$\nabla_{\boldsymbol{\alpha}_{i}^{k}} \log P(\boldsymbol{\phi}_{i}^{k,r}) \leq \sqrt{N \cdot \max\left\{ \left| \frac{1 - p_{i,j_{i}}^{k,r}}{\tau \boldsymbol{\alpha}_{i,j_{i}}^{k,r}} \right|, \left| - \frac{p_{i,j}^{k,r}}{\tau \boldsymbol{\alpha}_{i,j}^{k,r}} \right| \right\}^{2}} \leq \sqrt{\frac{N}{\tau^{2} \nu^{2}}}$$

 $\frac{1}{I(I-1)} \le \frac{2}{I^2}$

5) is because when $I \ge 2$:

The following lemma shows that the $\mathbb{E}_{\Phi^k \sim GS(\boldsymbol{\alpha}^k)} [\mathcal{L}(\Phi^k, \Psi^k)]$ is L-smooth for $\boldsymbol{\alpha}$. This is crucial for the later convergence analysis of BDPL and Fed-BDPL and is a necessity for convergence proofs.

Lemma 2. L-smooth for $\boldsymbol{\alpha}^k$: At the k-th client, the Φ^k is sampled from probability matrix \boldsymbol{p}^k , and $\boldsymbol{p}^k = GS(\boldsymbol{\alpha}^k), \ \alpha_{i,j} \geq \nu > 0$ for i = 1, ..., n and $j = 1, ..., N, \ \tau > 0$ is the temperature parameter, $\mathbb{E}_{\Phi^k \sim GS(\boldsymbol{\alpha}^k)} \left[\mathcal{L}(\Phi^k, \Psi^k) \right]$ is L-smooth for $\boldsymbol{\alpha}^k$ and $L = \frac{nGN(\tau+1)}{\tau^2\nu^2}$, and then for t-th iteration, the following inequality is satisfied:

$$\begin{split} & \mathbb{E}_{\Phi_{t+1}^{k} \sim GS(\boldsymbol{\alpha}_{t+1}^{k})} \left[\mathcal{L}(\Phi_{t+1}^{k}, \Psi^{k}) \right] - \mathbb{E}_{\Phi_{t}^{k} \sim GS(\boldsymbol{\alpha}_{t}^{k})} \left[\mathcal{L}(\Phi_{t}^{k}, \Psi^{k}) \right] \\ & \leq \left\langle \nabla_{\boldsymbol{\alpha}^{k}} \mathbb{E}_{\Phi_{t}^{k} \sim GS(\boldsymbol{\alpha}_{t}^{k})} \left[\mathcal{L}(\Phi_{t}^{k}, \Psi^{k}) \right], \boldsymbol{\alpha}_{t+1}^{k} - \boldsymbol{\alpha}_{t}^{k} \right\rangle + \frac{L}{2} \left\| \boldsymbol{\alpha}_{t+1}^{k} - \boldsymbol{\alpha}_{t}^{k} \right\|^{2} \end{split}$$

Proof. The objective function:

$$\mathbb{E}_{\Phi^k \sim \mathrm{GS}(\boldsymbol{\alpha}^k)} \left[\mathcal{L}(\Phi^k, \Psi^k) \right] = \sum_{\phi_1^k \sim \mathrm{GS}(\boldsymbol{\alpha}_1^k)} \cdots \sum_{\phi_n^k \sim \mathrm{GS}(\boldsymbol{\alpha}_n^k)} \left(\mathcal{L}(\Phi^k, \Psi^k) \prod_{i=1}^n P(\phi_i^k) \right)$$

We can compute the Hessian of the objective function, $\forall i', i'' \in 1, \dots, n \text{ and } j', j'' \in 1, \dots, N$: If $i' \neq i''$:

$$\begin{aligned} \frac{\partial^2}{\partial \alpha_{i',j'} \partial \alpha_{i'',j''}} \mathbb{E}_{\Phi^k \sim GS(\mathbf{a}^k)} \left[\mathcal{L}(\Phi^k, \Psi^k) \right] \\ = \sum_{\phi_1^k \sim \mathbf{p}_1^k} \cdots \sum_{\phi_n^k \sim \mathbf{p}_n^k} \left(\mathcal{L}(\Phi^k, \Psi^k) \frac{\partial^2}{\partial \alpha_{i',j'} \partial \alpha_{i'',j''}} \prod_{i=1}^n P(\phi_i^k) \right) \\ = \sum_{\phi_1^k \sim \mathbf{p}_1^k} \cdots \sum_{\phi_{i'-1}^k \sim \mathbf{p}_{i'-1}^k} \sum_{\phi_{i'+1}^{k'-1} \phi_{i'+1}^k \sim \mathbf{p}_{i'+1}^k} \cdots \sum_{\phi_{i''-1}^k \phi_{i''+1}^k \sim \mathbf{p}_{i''}^k} \sum_{i=1}^n P(\phi_i^k) \right) \\ = \sum_{\phi_1^k \sim \mathbf{p}_i^k} \sum_{\phi_{i'}^k \sim \mathbf{p}_{i''}^k} \left(\mathcal{L}(\Phi^k, \Psi^k) \frac{\partial^2}{\partial \alpha_{i',j'} \partial \alpha_{i'',j''}} \prod_{i=1}^n P(\phi_i^k) \right) \right) \\ = \sum_{\phi_1^k \sim \mathbf{p}_1^k} \cdots \sum_{\phi_{i'-1}^k \sim \mathbf{p}_{i'-1}^k} \sum_{\phi_{i'+1}^{k'+1} \sim \mathbf{p}_{i'+1}^k} \cdots \sum_{\phi_{i''-1}^k \sim \mathbf{p}_{i''+1}^k} \sum_{i=1}^n P(\phi_i^k) \right) \\ = \sum_{\phi_1^k \sim \mathbf{p}_1^k} \cdots \sum_{\phi_{i'-1}^k \sim \mathbf{p}_{i'-1}^k \phi_{i'+1}^k \sim \mathbf{p}_{i'+1}^k} \cdots \sum_{\phi_{i''-1}^k \sim \mathbf{p}_{i''-1}^k \phi_{i''+1}^k \cdots \phi_{i''+1}^k \cdots \sum_{\phi_{i''-1}^k \sim \mathbf{p}_{i''}^k} \cdots \sum_{\phi_{i''-1}^k \sim \mathbf{p}_{i''+1}^k} \sum_{i=1}^n P(\phi_i^k) \right) \\ = \sum_{\phi_1^k \sim \mathbf{p}_1^k} \cdots \sum_{\phi_{i'-1}^k \sim \mathbf{p}_{i'-1}^k \phi_{i'+1}^k \otimes \mathbf{p}_{i''+1}^k \cdots \sum_{\phi_{i''-1}^k \sim \mathbf{p}_{i''+1}^k} \cdots \sum_{\phi_{i''-1}^k \sim \mathbf{p}_{i''+1}^k \cdots \sum_{\phi_{i''-1}^k \sim \mathbf{p}_{i''+1}^k} \cdots \sum_{\phi_{i''-1}^k \sim \mathbf{p}_{i''+1}^k \cdots \sum_{\phi_{i''-1}^k \sim \mathbf{p}_{i''+1}^k} \cdots \sum_{\phi_{i''-1}^k \sim \mathbf{p}_{i''+1}^k \cdots \sum_{\phi_{i''+1}^k \sim \mathbf{p}_{i''+1}^k \cdots \sum_{\phi_{i''-1}^k \sim \mathbf{p}_{i''+1}^k \cdots \sum_{\phi_{i'''+1}^k \sim \mathbf{p}_{i''+1}^k \cdots \sum_{\phi_{i''-1}^k \sim \mathbf{p}_{i''+1}^k \cdots \sum_{\phi_{i$$

Then:

$$\frac{\partial^{2} \left[P(\phi_{i'}^{k}) P(\phi_{i''}^{k}) \right]}{\partial \alpha_{i',j'} \partial \alpha_{i'',j''}} = \begin{cases} \frac{1 - p_{i',j'_{i}}^{k}}{\tau \alpha_{i',j'_{i}}^{k}} \cdot \frac{1 - p_{i'',j''_{i}}^{k}}{\tau \alpha_{i'',j''_{i}}^{k}}, \text{ if } j' = j'_{i} \text{ and } j'' = j''_{i} \\ \frac{1 - p_{i',j'_{i}}^{k}}{\tau \alpha_{i',j'_{i}}^{k}} \cdot \left(- \frac{p_{i'',j''_{i}}}{\tau \alpha_{i'',j''_{i}}^{k}} \right), \text{ if } j' = j'_{i} \text{ and } j'' \neq j''_{i} \\ - \frac{p_{i',j'_{i}}^{k}}{\tau \alpha_{i'',j'}^{k}} \cdot \frac{1 - p_{i'',j''_{i}}^{k}}{\tau \alpha_{i'',j''_{i}}^{k}}, \text{ if } j' \neq j'_{i} \text{ and } j'' \neq j''_{i} \\ - \frac{p_{i',j'_{i}}}{\tau \alpha_{i',j'}^{k}} \cdot \left(- \frac{p_{i'',j''_{i}}}{\tau \alpha_{i'',j''_{i}}^{k}} \right), \text{ if } j' \neq j'_{i} \text{ and } j'' \neq j''_{i} \end{cases}$$

$$(25)$$

$$\left| \frac{\partial^2 \left[P(\phi_{i'}^k) P(\phi_{i''}^k) \right]}{\partial \alpha_{i',j'} \partial \alpha_{i'',j''}} \right| \le \frac{1}{\tau^2 \nu^2}$$
(26)

Further, based on Assumption 2:

Based on $\alpha_{i,j}^k \ge \nu > 0$:

$$\begin{split} \left| \frac{\partial^2}{\partial \alpha_{i',j'} \partial \alpha_{i'',j''}} \mathbb{E}_{\Phi^k \sim GS(\mathbf{a}^k)} \left[\mathcal{L}(\Phi^k, \Psi^k) \right] \right| \\ &\leq \sum_{\phi_1^k \sim \mathbf{p}_1^k} \cdots \sum_{\phi_{i'-1}^k \sim \mathbf{p}_{i'-1}^k} \sum_{\phi_{i'+1}^k \sim \mathbf{p}_{i'+1}^k} \cdots \sum_{\phi_{i''-1}^k \sim \mathbf{p}_{i''-1}^k} \sum_{\phi_{i''+1}^k \sim \mathbf{p}_{i''+1}^k} \cdots \sum_{\phi_{i''+1}^k \sim \mathbf{p}_{i''+1}^k} \cdots \sum_{\phi_n^k \sim \mathbf{p}_n^k} \left(\mathcal{L}(\Phi^k, \Psi^k) \prod_{\substack{i = 1 \\ i \neq i' \\ i \neq i''}}^n P(\phi_i^k) \right) \cdot \frac{1}{\tau^2 \nu^2} \\ &\leq \frac{G}{\tau^2 \nu^2} \\ \text{If } i' = i'': \\ &\frac{\partial^2}{\partial \alpha_{i',j'} \partial \alpha_{i',j''}} \mathbb{E}_{\Phi^k \sim GS(\mathbf{a}^k)} \left[\mathcal{L}(\Phi^k, \Psi^k) \right] \\ &= \sum_{\phi_1^k \sim \mathbf{p}_1^k} \cdots \sum_{\phi_{i'-1}^k \sim \mathbf{p}_{i'+1}^k} \sum_{\phi_{i'+1}^k \sim \mathbf{p}_{i'+1}^k} \cdots \sum_{\phi_n^k \sim \mathbf{p}_n^k} \left(\mathcal{L}(\Phi^k, \Psi^k) \frac{\partial^2 \left[P(\phi_{i'}^k) \right]}{\partial \alpha_{i',j'} \partial \alpha_{i',j''}} \prod_{\substack{i = 1 \\ i \neq i'}}^n P(\phi_i^k) \right) \\ &= 0 \quad \text{if } i \neq i' \quad \text{we set if } i \neq i'' \end{split}$$

Similar to the analysis in case $i' \neq i''$, we can get:

$$\left| \frac{\partial^2 \left[P(\phi_{i'}^k) \right]}{\partial \alpha_{i',j'} \partial \alpha_{i',j''}} \right| \le \max\left\{ p, 1-p \right\} \cdot \frac{(\tau+1)}{\tau^2 \nu^2} \le \frac{\tau+1}{\tau^2 \nu^2}$$

$$\left|\frac{\partial^2}{\partial \alpha_{i',j'}} \mathbb{E}_{\Phi^k \sim \mathrm{GS}(\boldsymbol{\alpha}^k)} \left[\mathcal{L}(\Phi^k, \Psi^k)\right]\right| \le \frac{(\tau+1)\,G}{\tau^2 \nu^2} \tag{27}$$

Finally, with $H(\alpha)$ denoting the Hessian matrix of $\mathbb{E}_{\Phi^k \sim GS(\alpha^k)} [\mathcal{L}(\Phi^k, \Psi^k)]$, we can get:

$$\|H(\boldsymbol{\alpha})\|_{2} \leq \|H(\boldsymbol{\alpha})\|_{F} \leq \sqrt{n(n-1)N^{2}\left(\frac{G}{\tau^{2}\nu^{2}}\right)^{2} + nN^{2}\left(\frac{(\tau+1)G}{\tau^{2}\nu^{2}}\right)^{2}} \leq \frac{nGN(\tau+1)}{\tau^{2}\nu^{2}}$$
(28)

1025 According to Lemma 1.2.2 in (Nesterov et al., 2018), $\mathbb{E}_{\Phi^k \sim GS(\boldsymbol{\alpha}^k)} \left[\mathcal{L}(\Phi^k, \Psi^k) \right]$ is L-smooth for $\boldsymbol{\alpha}^k$ and $L = \frac{nGN(\tau+1)}{\tau^2 \nu^2}$,

$$\frac{1}{T}\sum_{t=0}^{T-1} \left\|\nabla_{\boldsymbol{\alpha}^{k}}F^{k}(\boldsymbol{\alpha}_{t}^{k},\Psi^{k})\right\|^{2} \leq \frac{4G}{T\eta} + \frac{2\eta nL\sigma_{\psi}^{2}}{B^{k}} + \frac{2\eta nL\sigma_{\alpha}^{2}}{I^{2}}$$
(29)

Proof. According to Lemma 2:

$$\begin{split} & \mathbb{E}_{\Phi_{t+1}^{k}\sim \mathrm{GS}(\boldsymbol{\alpha}_{t+1}^{k})}\left[\mathcal{L}(\Phi_{t+1}^{k},\Psi^{k})\right] - \mathbb{E}_{\Phi_{t}^{k}\sim \mathrm{GS}(\boldsymbol{\alpha}_{t}^{k})}\left[\mathcal{L}(\Phi_{t}^{k},\Psi^{k})\right] \\ & \leq \left\langle \nabla_{\boldsymbol{\alpha}^{k}}\mathbb{E}_{\Phi_{t}^{k}\sim \mathrm{GS}(\boldsymbol{\alpha}_{t}^{k})}\left[\mathcal{L}(\Phi_{t}^{k},\Psi^{k})\right], \boldsymbol{\alpha}_{t+1}^{k} - \boldsymbol{\alpha}_{t}^{k}\right\rangle + \frac{L}{2}\left\|\boldsymbol{\alpha}_{t+1}^{k} - \boldsymbol{\alpha}_{t}^{k}\right\|^{2} \\ & \leq \sum_{i=1}^{n}\left[\left\langle \nabla_{\boldsymbol{\alpha}_{i}^{k}}F^{k}(\boldsymbol{\alpha}_{t}^{k},\Psi^{k}), -\eta\cdot\hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}}f^{k}(\boldsymbol{\alpha}_{t}^{k},\mathcal{B}_{t}^{k})\right\rangle + \frac{L\eta^{2}}{2}\left\|\hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}}f^{k}(\boldsymbol{\alpha}_{t}^{k},\mathcal{B}_{t}^{k})\right\|^{2}\right] \end{split}$$

Both sides take expectations for $\mathbb{E}_{\{\Phi_t^r \sim GS(\boldsymbol{\alpha}_t)\}_{r=1}^I}$ and $\mathbb{E}_{\mathcal{B}_t}$ at the same time, we abbreviate $\mathbb{E}_{\{\Phi^r \sim GS(\boldsymbol{\alpha})\}_{r=1}^I}$ as $\mathbb{E}_{\{\Phi^r\}_{r=1}^I}$:

$$\mathbb{E}_{\{\Phi^{r}\}_{r=1}^{I}} \mathbb{E}_{\mathcal{B}_{t}} \left\{ \mathbb{E}_{\Phi_{t+1}^{k} \sim \mathrm{GS}(\boldsymbol{\alpha}_{t+1}^{k})} \left[\mathcal{L}(\Phi_{t+1}^{k}, \Psi^{k}) \right] - \mathbb{E}_{\Phi_{t}^{k} \sim \mathrm{GS}(\boldsymbol{\alpha}_{t}^{k})} \left[\mathcal{L}(\Phi_{t}^{k}, \Psi^{k}) \right] \right\}$$

$$\leq \mathbb{E}_{\{\Phi^{r}\}_{r=1}^{I}} \mathbb{E}_{\mathcal{B}_{t}} \left\{ \sum_{i=1}^{n} \left[\left\langle \nabla_{\boldsymbol{\alpha}_{i}^{k}} F^{k}(\boldsymbol{\alpha}_{t}^{k}, \Psi^{k}), -\eta \cdot \hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\rangle + \frac{L\eta^{2}}{2} \left\| \hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\|^{2} \right] \right\}$$

$$= \sum_{i=1}^{n} \underbrace{\mathbb{E}_{\{\Phi^{r}\}_{r=1}^{I}} \mathbb{E}_{\mathcal{B}_{t}} \left[\left\langle \nabla_{\boldsymbol{\alpha}_{i}^{k}} F^{k}(\boldsymbol{\alpha}_{t}^{k}, \Psi^{k}), -\eta \cdot \hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\rangle + \frac{L\eta^{2}}{2} \left\| \hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\|^{2} \right]}_{a)}$$

1055 For a):

$$\begin{split} & \mathbb{E}_{\left\{\Phi^{r}\right\}_{r=1}^{I}} \mathbb{E}_{\mathcal{B}_{t}} \left[\left\langle \nabla_{\boldsymbol{\alpha}_{t}^{k}} F^{k}(\boldsymbol{\alpha}_{t}^{k}, \Psi^{k}), -\eta \cdot \hat{\nabla}_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\rangle + \frac{L\eta^{2}}{2} \left\| \hat{\nabla}_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\|^{2} \right] \\ &= \left\langle \nabla_{\boldsymbol{\alpha}_{t}^{k}} F^{k}(\boldsymbol{\alpha}_{t}^{k}, \Psi^{k}), -\eta \cdot \mathbb{E}_{\left\{\Phi^{r}\right\}_{r=1}^{I}} \mathbb{E}_{\mathcal{B}_{t}} \left[\hat{\nabla}_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right] \right\rangle + \frac{L\eta^{2}}{2} \mathbb{E}_{\left\{\Phi^{r}\right\}_{r=1}^{I}} \mathbb{E}_{\mathcal{B}_{t}} \left[\left\| \hat{\nabla}_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right] \right\rangle \\ &+ \frac{L\eta^{2}}{2} \left\{ \left\| \mathbb{E}_{\left\{\Phi^{r}\right\}_{r=1}^{I}} \mathbb{E}_{\mathcal{B}_{t}} \left[\hat{\nabla}_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right] \right\|^{2} + \operatorname{Var}_{\mathcal{B}_{t}, \left\{\Phi^{r}\right\}_{r=1}^{I}} \left[\hat{\nabla}_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right] \right\} \\ & \left(\stackrel{2}{=} -\eta \cdot \left\langle \nabla_{\boldsymbol{\alpha}_{t}^{k}} F^{k}(\boldsymbol{\alpha}_{t}^{k}, \Psi^{k}), \nabla_{\boldsymbol{\alpha}_{t}^{k}} F^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\rangle + \frac{L\eta^{2}}{2} \left\| \nabla_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\|^{2} \\ &+ \frac{L\eta^{2}}{2} \mathbb{E}_{\left\{\Phi^{r}\right\}_{r=1}^{I} \mathbb{E}_{\mathcal{B}_{t}} \left[\left\| \hat{\nabla}_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\rangle + \frac{L\eta^{2}}{2} \left\| \nabla_{\boldsymbol{\alpha}_{t}^{k}} F^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\|^{2} \\ &+ \frac{L\eta^{2}}{2} \mathbb{E}_{\left\{\Phi^{r}\right\}_{r=1}^{I} \mathbb{E}_{\mathcal{B}_{t}} \left[\left\| \hat{\nabla}_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\rangle + \frac{L\eta^{2}}{2} \left\| \nabla_{\boldsymbol{\alpha}_{t}^{k}} F^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\|^{2} \\ &+ \frac{L\eta^{2}}{2} \mathbb{E}_{\left\{\Phi^{r}\right\}_{r=1}^{I} \mathbb{E}_{\mathcal{B}_{t}} \left[\left\| \hat{\nabla}_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) - \mathbb{E}_{\left\{\Phi^{r}\right\}_{r=1}^{I} \mathbb{E}_{\mathcal{B}_{t}} \left[\hat{\nabla}_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right] \right\|^{2} \\ &+ L\eta^{2} \cdot \mathbb{E}_{\left\{\Phi^{r}\right\}_{r=1}^{I} \mathbb{E}_{\mathcal{B}_{t}} \left[\left\| \mathbb{E}_{\mathcal{B}_{t}} \left[\hat{\nabla}_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right] - \mathbb{E}_{\left\{\Phi^{r}\right\}_{r=1}^{I} \mathbb{E}_{\mathcal{B}_{t}} \left[\hat{\nabla}_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right] \right\|^{2} \\ &+ L\eta^{2} \cdot \mathbb{E}_{\mathcal{B}_{t}} \operatorname{Var}_{\left\{\Phi^{r}\right\}_{r=1}^{I}} \left[\mathbb{E}_{\mathcal{B}_{t}} \left[\hat{\nabla}_{\boldsymbol{\alpha}_{t}^{k}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right] \right] \end{aligned}$$

$$\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\end{array} & \left(\frac{L\eta^{2}}{2} - \eta\right) \left\| \nabla_{\boldsymbol{\alpha}_{i}^{k}} F^{k}(\boldsymbol{\alpha}_{t}^{k}, \Psi^{k}) \right\|^{2} + L\eta^{2} \cdot \frac{\sigma_{\psi}^{2}}{B} + L\eta^{2} \cdot \mathbb{E}_{\mathcal{B}_{t}} \operatorname{Var}_{\{\Phi^{r}\}_{r=1}^{I}} \left[\hat{\nabla}_{\boldsymbol{\alpha}_{i}} f^{k}(\boldsymbol{\alpha}_{t}^{k}, \Psi^{k}) \right] \\
\end{array} \\
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\end{array} & \left(\begin{array}{ll}
\end{array} & \left(\begin{array}{ll}
\end{array} & \left(\frac{L\eta^{2}}{2} - \eta\right) \left\| \nabla_{\boldsymbol{\alpha}_{i}^{k}} F^{k}(\boldsymbol{\alpha}_{t}^{k}, \Psi^{k}) \right\|^{2} + L\eta^{2} \cdot \frac{\sigma_{\psi}^{2}}{B} + L\eta^{2} \cdot \frac{\sigma_{\alpha}^{2}}{I^{2}} \\
\end{array} \\
\end{array} \\
\begin{array}{ll}
\end{array} & \left(\begin{array}{ll}
\end{array} & \left(\begin{array}{ll}
\end{array} & \left(\frac{L\eta^{2}}{2} - \eta\right) \left\| \nabla_{\boldsymbol{\alpha}_{i}^{k}} F^{k}(\boldsymbol{\alpha}_{t}^{k}, \Psi^{k}) \right\|^{2} + L\eta^{2} \cdot \frac{\sigma_{\psi}^{2}}{B} + L\eta^{2} \cdot \frac{\sigma_{\alpha}^{2}}{I^{2}} \\
\end{array} \\
\end{array} \\
\end{array} \\$$

where (1) uses the independence between Ψ sampling and Φ sampling; (2) use the unbiasedness of stochastic gradient and variance-reduced policy gradient in Assumption 1 and Lemma 1; (3) use the inequality $||a + b||^2 \le 2 ||a||^2 + 2 ||b||^2$; (4) and (5) use the bounded variance of stochastic gradient and variance-reduced policy gradient in Assumption 1 and Lemma 1. Then:

$$\mathbb{E}_{\Phi_{t+1}\sim \mathrm{GS}(\boldsymbol{a}_{t+1}^k)}\left[\mathcal{L}(\Phi_{t+1}^k,\Psi^k)\right] - \mathbb{E}_{\Phi_t^k\sim \mathrm{GS}(\boldsymbol{a}_t^k)}\left[\mathcal{L}(\Phi_t^k,\Psi^k)\right]$$

$$\leq \sum_{i=1}^{n} \left[\left(\frac{L\eta^2}{2} - \eta \right) \left\| \nabla_{\mathbf{a}_i^k} F^k(\mathbf{a}_t^k, \Psi^k) \right\|^2 + L\eta^2 \cdot \frac{\sigma_{\psi}^2}{B} + L\eta^2 \cdot \frac{\sigma_{\alpha}^2}{I^2} \right]$$

We combine the gradient of α_i for each prompt token:

$$\mathbb{E}_{\Phi_{t+1}^k \sim \mathrm{GS}(\boldsymbol{\alpha}_{t+1})} \left[\mathcal{L}(\Phi_{t+1}^k, \Psi^k) \right] - \mathbb{E}_{\Phi_t^k \sim \mathrm{GS}(\boldsymbol{\alpha}_t^k)} \left[\mathcal{L}(\Phi_t^k, \Psi^k) \right] \\ \leq \left(\frac{L\eta^2}{2} - \eta \right) \left\| \nabla_{\boldsymbol{\alpha}^k} F^k(\boldsymbol{\alpha}_t^k, \Psi^k) \right\|^2 + nL\eta^2 \cdot \frac{\sigma_{\psi}^2}{B} + nL\eta^2 \cdot \frac{\sigma_{\alpha}^2}{I^2} \right]$$

where $\boldsymbol{\alpha}^k = (\boldsymbol{\alpha}_1^k, \cdots, \boldsymbol{\alpha}_i^k, \cdots \boldsymbol{\alpha}_n^k).$

We let $\eta \leq \frac{1}{L}$, then both sides accumulate with respect to $t = 0, 1, \dots, T-1$ and divide by T:

$$\begin{array}{ll} 1103 & & \frac{1}{T} \sum_{t=0}^{T-1} \left(\eta - \frac{L\eta^2}{2} \right) \left\| \nabla_{\pmb{\alpha}^k} F^k(\pmb{\alpha}_t^k, \Psi^k) \right\|^2 \\ 1105 & & \leq \frac{1}{T} \sum_{t=0}^{T-1} \left[\mathbb{E}_{\Phi_t^k \sim \mathrm{GS}(\pmb{\alpha}_t^k)} \left[\mathcal{L}(\Phi_t^k, \Psi^k) \right] - \mathbb{E}_{\Phi_{t+1}^k \sim \mathrm{GS}(\pmb{\alpha}_{t+1}^k)} \left[\mathcal{L}(\Phi_{t+1}^k, \Psi^k) \right] \right] + \frac{nL\eta^2 \sigma_{\psi}^2}{B} + \frac{nL\eta^2 \sigma_{\alpha}^2}{I^2} \\ 1108 & & \text{Then,} \end{array}$$

$$\begin{aligned} &\frac{1}{T}\sum_{t=0}^{T-1} \left\|\nabla_{\boldsymbol{\alpha}^{k}}F^{k}(\boldsymbol{\alpha}_{t}^{k},\Psi^{k})\right\|^{2} \\ &\leq \frac{\mathbb{E}_{\Phi_{0}^{k}\sim\mathrm{GS}(\boldsymbol{\alpha}_{0}^{k})}\left[\mathcal{L}(\Phi_{0}^{k},\Psi^{k})\right] - \mathbb{E}_{\Phi_{T}^{k}\sim\mathrm{GS}(\boldsymbol{\alpha}_{T}^{k})}\left[\mathcal{L}(\Phi_{T}^{k},\Psi^{k})\right]}{T}\frac{2}{2\eta - L\eta^{2}} + \frac{2\eta}{2 - L\eta}\left(\frac{nL\sigma_{\psi}^{2}}{B} + \frac{nL\sigma_{\alpha}^{2}}{I^{2}}\right) \\ &\leq \frac{2\left(\mathbb{E}_{\Phi_{0}^{k}\sim\mathrm{GS}(\boldsymbol{\alpha}_{0}^{k})}\left[\mathcal{L}(\Phi_{0}^{k},\Psi^{k})\right] - \mathbb{E}_{\Phi_{T}^{k}\sim\mathrm{GS}(\boldsymbol{\alpha}_{T}^{k})}\left[\mathcal{L}(\Phi_{T}^{k},\Psi^{k})\right]\right)}{T\eta} + \frac{2\eta nL\sigma_{\psi}^{2}}{B} + \frac{2\eta nL\sigma_{\alpha}^{2}}{I^{2}} \\ &\operatorname{According to Assumption 2}, \mathbb{E}_{\Phi_{0}^{k}\sim\mathrm{GS}(\boldsymbol{\alpha}_{N}^{k})}\left[\mathcal{L}(\Phi_{0}^{k},\Psi^{k})\right] - \inf_{t}\mathbb{E}_{\Phi_{0}^{k}\sim\mathrm{GS}(\boldsymbol{\alpha}_{N}^{k})}\left[\mathcal{L}(\Phi_{T}^{k},\Psi^{k})\right] < 2G, \text{ then:} \end{aligned}$$

 $\Phi_t^k \sim \mathrm{GS}(\boldsymbol{a}_t^k) \left[\boldsymbol{\sim} \left(\boldsymbol{\Psi}_t, \boldsymbol{\Psi} \right) \right]$

$$\frac{1}{T}\sum_{t=0}^{T-1} \left\|\nabla_{\boldsymbol{\alpha}^{k}}F^{k}(\boldsymbol{\alpha}_{t}^{k},\Psi^{k})\right\|^{2} \leq \frac{4G}{T\eta} + \frac{2\eta nL\sigma_{\psi}^{2}}{B} + \frac{2\eta nL\sigma_{\alpha}^{2}}{I^{2}}$$

Remark 2. The convergence term of BDPL consists of three parts. The first term is typical of first-order optimization algorithms converging to non-convex functions. The second term is the stochasticity due to random mini-batch gradients. The first term and second term can be combined when $\eta = \frac{c}{\sqrt{T}}$ and $c = \sqrt{\frac{2B^k G}{nL\sigma_\psi^2}}$:

$$\frac{4G}{T\eta} + \frac{2\eta n L \sigma_{\psi}^2}{B^k} = \frac{4}{\sqrt{T}} \sqrt{\frac{2n L \sigma_{\psi}^2 G}{B}},$$

which can decrease as the number of iterations and mini-batch size increase Reddi et al. (2016); and the third term is the stochasticity due to prompt sampling, which decreases with the number of prompt samples.

1134 1135 1136 1137 Remark 3. According to the definition of Gumbel-Softmax (2), there is randomness about $\boldsymbol{u} = \{u_i\}_{i=1}^n \sim \text{Uniform}(\mathbf{0}, \mathbf{1}_n) \text{ in } \nabla_{\boldsymbol{\alpha}} F^k(\boldsymbol{\alpha}_t^k, \Psi^k) \text{ i.e. } \boldsymbol{p}^k = \text{GS}(\boldsymbol{\alpha}^k, \boldsymbol{u}) \text{ in fact and } \Phi^k \sim \boldsymbol{p}^k, \text{ we can further discuss the result in Lemma 3:}$

$$\mathbb{E}_{\boldsymbol{u}}\left\|\nabla_{\boldsymbol{\alpha}^{k}}F^{k}(\boldsymbol{\alpha}_{t}^{k},\Psi^{k})\right\|^{2} = \int_{\boldsymbol{0}}^{\boldsymbol{1}}\left\|\nabla_{\boldsymbol{\alpha}^{k}}F^{k}(\boldsymbol{\alpha}_{t}^{k},\Psi^{k})\right\|^{2}\boldsymbol{1}_{n}d\boldsymbol{u} \leq \left\|\nabla_{\boldsymbol{\alpha}^{k}}F^{k}(\boldsymbol{\alpha}_{t}^{k},\Psi^{k})\right\|^{2}$$

1140 Considering the randomness of u, we can still obtain the convergence of BDPL. In subsequent 1141 analyzes of Fed-BDPL, we can obtain similar result for u.

1142 1143 **Corollary 3.** Convergence rate of BDPL: Let $\eta = \min\left\{\frac{1}{L}, \frac{1}{\sqrt{T}}\right\}$, we can get the following con-1144 vergence rate for BDPL:

 $\frac{1}{T}\sum_{t=0}^{T-1} \left\|\nabla_{\pmb{\alpha}^k}F^k(\pmb{\alpha}_t^k,\Psi^k)\right\|^2 \le \mathcal{O}(\frac{1}{\sqrt{T}})$

(30)

1138 1139

1148 1149

Proof. Convergence rate:

$$\frac{1}{T}\sum_{t=0}^{T-1} \left\|\nabla_{\boldsymbol{\alpha}^{k}}F^{k}(\boldsymbol{\alpha}_{t}^{k},\Psi^{k})\right\|^{2} = \frac{4G}{\sqrt{T}} + \frac{1}{\sqrt{T}} \cdot \frac{2nL\sigma_{\psi}^{2}}{B} + \frac{1}{\sqrt{T}} \cdot \frac{2nL\sigma_{\alpha}^{2}}{I^{2}} = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

1156 A.4 CONVERGENCE ANALYSIS OF FEDERATED PROMPT TUNING

¹¹⁵⁷ Definition in the Federated Black-box Discrete Prompt Learning:

Data slicing: The data slices in each client are defined as $D^k = \{D^k\}_{k=1}^K = \{\Psi^k, Y^k\}_{k=1}^K$, where $\Psi^k = \{\psi_m^k\}_{m=1}^{M^k}$ denote the input sentence and $Y^k = \{y_m^k\}_{m=1}^{M^k}$ denote the label, k is the index of the client, and m is the index of the sample, there are totally M^k samples in the dataset D^k . The clients hold the sampling probability vector $\boldsymbol{q} = \{q^{[k]}\}_{k=1}^K = \{\frac{M^k}{M}\}_{k=1}^K$. K_* is the number of clients selected and U_t is the set of corresponding clients. **Average parameter** $\boldsymbol{\alpha}$:

1168

1169

1170

1171 1172

1173

1174 1175 1176

1178

1187

Average loss:

$$\mathbb{E}_{\Phi_t \sim \mathrm{GS}(\boldsymbol{\alpha}_t)} \left[\mathcal{L}(\Phi_t, \Psi) \right] = \sum_{k=1}^K q^{[k]} \cdot \mathbb{E}_{\Phi_t^k \sim \mathrm{GS}(\boldsymbol{\alpha}_t^k)} \left[\mathcal{L}(\Phi_t^k, \Psi^k) \right]$$

 $\boldsymbol{\alpha}_{i,(t)} = \frac{1}{K_*} \sum_{k \in II_*} \boldsymbol{\alpha}_{i,(t)}^k$

 $oldsymbol{lpha}_t = rac{1}{K_*} \sum_{k \in II_*} oldsymbol{lpha}_t^k$

1177 Average gradient:

$$\nabla_{\boldsymbol{\alpha}_{i}} F(\boldsymbol{\alpha}, \Psi^{k}) \stackrel{def}{=} \sum_{k=1}^{K} q^{[k]} \cdot \nabla_{\boldsymbol{\alpha}_{i}} F^{k}(\boldsymbol{\alpha}, \Psi^{k})$$
(31)

$$\nabla_{\boldsymbol{\alpha}_{i}^{k}} F(\boldsymbol{\alpha}^{k}, \Psi^{k}) \stackrel{def}{=} \sum_{k=1}^{K} q^{[k]} \cdot \nabla_{\boldsymbol{\alpha}_{i}^{k}} F^{k}(\boldsymbol{\alpha}^{k}, \Psi^{k})$$
(32)

1183
1184
1185
1185
1186

$$\nabla_{\boldsymbol{\alpha}_{i}^{k}}F^{*}(\boldsymbol{\alpha}^{k},\Psi^{k}) \stackrel{def}{=} \frac{1}{K_{*}} \sum_{k \in U_{t}} \nabla_{\boldsymbol{\alpha}_{i}^{k}}F^{k}(\boldsymbol{\alpha}^{k},\Psi^{k})$$
(33)

$$\nabla_{\boldsymbol{\alpha}_{i}^{k}} f^{*}(\boldsymbol{\alpha}^{k}, \mathcal{B}^{k}) \stackrel{def}{=} \frac{1}{K_{*}} \sum_{k \in U_{t}} \nabla_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}^{k}, \mathcal{B}^{k})$$
(34)

$$\hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}}f^{*}(\boldsymbol{\alpha}^{k}, \mathcal{B}^{k}) \stackrel{def}{=} \frac{1}{K_{*}} \sum_{k \in U_{*}} \hat{\nabla}_{\boldsymbol{\alpha}_{i}}f^{k}(\boldsymbol{\alpha}^{k}, \mathcal{B}^{k})$$
(35)

 $\boldsymbol{\alpha}_{i,(t+1)} = \boldsymbol{\alpha}_{i,(t)} - \eta \cdot \hat{\nabla}_{\boldsymbol{\alpha}^{k}} f^{*}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}^{k})$

Note that when $t \neq sE$, although the above some definitions don't exist in algorithm and experiment, we can still calculate and analyze them. In order to facilitate the analysis of the iteration process, we assume that they exist.

Average mini-batch stochastic variance-reduced policy gradient descent:

The following lemma shows that the local gradient is biased compared to the global gradient due to the fact that the distribution of data on different clients may be different (heterogeneity), and the bias of the gradient can be analyzed using the bias about α (Haddadpour & Mahdavi (2019)).

Lemma 4. Bound bias between local and average α . Let $\alpha_{i,j} \geq \nu > 0$ for i = 1, ..., n and $j = 1, ..., N, \eta$ is the learning rate, let \mathbb{E} represents $\mathbb{E}_{\left\{\Phi_t^{k,r} \sim GS(\boldsymbol{\alpha}_t^k)\right\}_{r=1}^{I}}$ and $\mathbb{E}_{\mathcal{B}_t^k}$, the bias between the local and the average α can be bounded:

$$\frac{1}{T}\sum_{t=0}^{T-1}\sum_{k=1}^{K}q^{[k]} \cdot \mathbb{E}\left\|\boldsymbol{\alpha}_{i,(t)} - \boldsymbol{\alpha}_{i,(t)}^{k}\right\|^{2}$$

 $\leq E\eta^2 \left(\frac{2\sigma_\psi^2}{B} + \frac{2\sigma_\alpha^2}{I^2}\right) \left(1 + \frac{1}{K_*}\right) + \frac{2\eta^2 E^2 \lambda}{T} \left(1 + \frac{1}{K_*}\right) \sum_{t=0}^{T-1} \left\|\nabla_{\mathbf{a}_t^k} F(\mathbf{a}_t^k, \Psi^k)\right\|^2$

Proof. Define:

$$t_c \triangleq \lfloor rac{t}{E}
floor E$$
 $oldsymbol{lpha}_{i,(t_c)} = rac{1}{K_*} \sum_{k \in U_{t_c}} oldsymbol{lpha}_{i,(t_c)}^k$

Then, for $t_c + 1 \le t < t_c + E$:

 $\boldsymbol{\alpha}_{i,(t)}^{k} = \boldsymbol{\alpha}_{i,(t_{c})} - \sum_{\alpha=t}^{t-1} \eta \cdot \hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}} f^{k}(\boldsymbol{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k})$ (36)

$$\boldsymbol{\alpha}_{i,(t)} = \boldsymbol{\alpha}_{i,(t_c)} - \frac{1}{K_*} \sum_{k \in U_{\rho}} \sum_{\rho=t_c}^{t-1} \eta \cdot \hat{\nabla}_{\boldsymbol{\alpha}_i^k} f^k(\boldsymbol{\alpha}_{\rho}^k, \mathcal{B}_{\rho}^k)$$
(37)

For the k-th client, let \mathbb{E} represents $\mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}}$ and $\mathbb{E}_{\mathcal{B}_{t}^{k}}$, and we take \mathbb{E} for $\left\|\boldsymbol{\alpha}_{i,(t)} - \boldsymbol{\alpha}_{i,(t)}^{k}\right\|^{2}$: $\mathbb{E}\left\|\boldsymbol{\alpha}_{k}\right\| = \boldsymbol{\alpha}_{k}^{k} \left\|\boldsymbol{\alpha}_{k}\right\|^{2}$

$$\begin{split} & \mathbb{E} \left\| \mathbf{\alpha}_{i,(t)}^{i} - \mathbf{\alpha}_{i,(t)}^{i} \right\| \\ & = \mathbb{E} \left\| \mathbf{\alpha}_{i,(t_{c})} - \frac{1}{K_{*}} \sum_{k \in U_{\rho}} \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) - \mathbf{\alpha}_{i,(t_{c})} + \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) \right\|^{2} \\ & = \mathbb{E} \left\| \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) - \frac{1}{K_{*}} \sum_{k \in U_{\rho}} \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) \right\|^{2} \\ & = \mathbb{E} \left\| \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) - \frac{1}{K_{*}} \sum_{k \in U_{\rho}} \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) \right\|^{2} \\ & = \mathbb{E} \left\| \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) \right\|^{2} + \mathbb{E} \left\| \frac{1}{K_{*}} \sum_{k \in U_{\rho}} \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) \right\|^{2} \\ & = 2\mathbb{E} \left\| \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) - \mathbb{E} \left[\sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) \right] \right\|^{2} + 2 \left\| \mathbb{E} \left[\sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) \right] \right\|^{2} \\ & = 2\mathbb{E} \left\| \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) - \mathbb{E} \left[\sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) \right] \right\|^{2} \\ & = 2\mathbb{E} \left\| \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) - \mathbb{E} \left[\sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) \right] \right\|^{2} \\ & = 2\mathbb{E} \left\| \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) - \mathbb{E} \left[\sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) \right] \right\|^{2} \\ & = 2\mathbb{E} \left\| \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) - \mathbb{E} \left[\sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) \right] \right\|^{2} \\ & = 2\mathbb{E} \left\| \sum_{\rho=t_{c}}^{t-1} \eta \cdot \hat{\nabla}_{\mathbf{\alpha}_{i}^{k}} f^{k}(\mathbf{\alpha}_{\rho}^{k}, \mathcal{B}_{\rho}^{k}) - \mathbb{E} \left[\sum_$$

$$\begin{aligned} &+ 2\mathbb{E} \left\| \frac{1}{K_{*}} \sum_{k \in U_{\mu}} \sum_{p=k_{*}}^{t-1} \eta \cdot \hat{\nabla}_{\alpha_{*}^{k}} f^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{B}_{\mu}^{k}) - \mathbb{E} \left[\frac{1}{K_{*}} \sum_{k \in U_{\mu}} \sum_{p=k_{*}}^{t-1} \eta \cdot \hat{\nabla}_{\alpha_{*}^{k}} f^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{B}_{\mu}^{k}) \right] \right\|^{2} \\ &+ 2\left\| \mathbb{E} \left[\frac{1}{K_{*}} \sum_{k \in U_{\mu}} \sum_{p=k_{*}}^{t-1} \eta \cdot \hat{\nabla}_{\alpha_{*}^{k}} f^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{B}_{\mu}^{k}) - \sum_{p=k_{*}}^{t-1} \eta \cdot \nabla_{\alpha_{*}^{k}} F^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{\Psi}^{k}) \right\|^{2} \\ &+ 2\mathbb{E} \left\| \frac{1}{K_{*}} \sum_{k \in U_{\mu}} \sum_{p=k_{*}}^{t-1} \eta \cdot \hat{\nabla}_{\alpha_{*}^{k}} f^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{B}_{\mu}^{k}) - \sum_{p=k_{*}}^{t-1} \eta \cdot \nabla_{\alpha_{*}^{k}} F^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{\Psi}^{k}) \right\|^{2} \\ &+ 2\mathbb{E} \left\| \frac{1}{K_{*}} \sum_{k \in U_{\mu}} \sum_{p=k_{*}}^{t-1} \eta \cdot \nabla_{\alpha_{*}^{k}} f^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{B}_{\mu}^{k}) - \frac{1}{K_{*}} \sum_{k \in U_{\mu}} \sum_{p=k_{*}}^{t-1} \eta \cdot \nabla_{\alpha_{*}^{k}} F^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{\Psi}^{k}) \right\|^{2} \\ &+ 2\mathbb{E} \left\| \frac{1}{K_{*}} \sum_{k \in U_{\mu}} \sum_{p=k_{*}}^{t-1} \eta \cdot \nabla_{\alpha_{*}^{k}} F^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{B}_{\mu}^{k}) - \nabla_{\alpha_{*}^{k}} F^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{\Psi}^{k}) \right\|^{2} \\ &+ 2\mathbb{E} \left\| \frac{1}{K_{*}} \sum_{k \in U_{\mu}} \sum_{p=k_{*}}^{t-1} \eta \cdot \left[\hat{\nabla}_{\alpha_{*}^{k}} f^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{B}_{\mu}^{k}) - \nabla_{\alpha_{*}^{k}} F^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{\Psi}^{k}) \right] \right\|^{2} + 2 \left\| \frac{1}{\mu} \sum_{p=k_{*}}^{t-1} \eta \cdot \nabla_{\alpha_{*}^{k}} F^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{\Psi}^{k}) \right\|^{2} \\ &+ 2\mathbb{E} \left\| \frac{1}{K_{*}} \sum_{k \in U_{\mu}} \sum_{p=k_{*}}^{t-1} \eta \cdot \left[\hat{\nabla}_{\alpha_{*}^{k}} f^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{B}_{\mu}^{k}) - \nabla_{\alpha_{*}^{k}} F^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{\Psi}^{k}) \right\|^{2} \\ &+ 2\mathbb{E} \left\| \frac{1}{K_{*}} \sum_{k \in U_{\mu}} \sum_{p=k_{*}}^{t-1} \eta \cdot \left[\hat{\nabla}_{\alpha_{*}^{k}} f^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{B}_{\mu}^{k}) - \nabla_{\alpha_{*}^{k}} F^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{\Psi}^{k}) \right\|^{2} \\ &+ 2\mathbb{E} \left\| \frac{1}{K_{*}} \sum_{k \in U_{\mu}} \sum_{p=k_{*}}^{t-1} \eta \cdot \left[\hat{\nabla}_{\alpha_{*}^{k}} f^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{B}_{\mu}^{k}) - \nabla_{\alpha_{*}^{k}} F^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{\Psi}^{k}) \right]^{2} \\ &+ 2\mathbb{E} \left\| \frac{1}{K_{*}} \sum_{k \in U_{\mu}} \sum_{p=k_{*}}^{t-1} \eta \cdot \left[\hat{\nabla}_{\alpha_{*}^{k}} f^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{B}_{\mu}^{k}) - \nabla_{\alpha_{*}^{k}} F^{k} (\boldsymbol{\alpha}_{\mu}^{k}, \boldsymbol{\Psi}^{k}) \right]^{2} \\ &+ 2\mathbb{E} \left\| \frac{1}{K_{*}} \sum_{k \in U_{\mu}} \sum_{p=k_{*}}^{t-1} \mathbb{E} \left\| \left\| \hat{\nabla}_{\alpha_{*}^{k}} f^{k} (\boldsymbol{\alpha}_{\mu}^$$

where (1) use inequality $||a - b||^2 \le 2 ||a||^2 + 2 ||b||^2$; (2) use $\mathbb{E} ||x - \mathbb{E} [x]||^2 = \mathbb{E} ||x||^2 - ||\mathbb{E} [x]||^2$; (3) use the unbiasedness of stochastic gradient and variance-reduced policy gradient in Assumption 1 and Lemma 1; (4) use the independence of mini-batch and prompt sampling in each client; (5) use inequality $||\sum_{z=1}^{Z} a_z||^2 \le Z \sum_{z=1}^{Z} ||a_z||^2$; (6) use the bounded variance of stochastic gradient and variance-reduced policy gradient in Assumption 1 and Lemma 1. 1296 Then, we sum both sides with respect to $k \in U_{\rho}$:

$$\begin{split} & \sum_{k \in U_{\rho}} \mathbb{E} \left\| \boldsymbol{\alpha}_{i,(t)} - \boldsymbol{\alpha}_{i,(t)}^{k} \right\|^{2} \\ & \leq 2\eta^{2} \sum_{\rho=t_{c}}^{t-1} \left(\frac{2\sigma_{\psi}^{2}}{B^{k}} + \frac{2\sigma_{\alpha}^{2}}{I^{2}} \right) (K_{*} + 1) + 2\eta^{2} (t - t_{c}) \left(1 + \frac{1}{K_{*}} \right) \sum_{k \in U_{\rho}} \sum_{\rho=t_{c}}^{t-1} \left\| \nabla_{\boldsymbol{\alpha}_{i}^{k}} F^{k}(\boldsymbol{\alpha}_{\rho}^{k}, \Psi^{k}) \right\|^{2} \end{split}$$

1304 We take the expectation of both sides about $\mathbb{E}_{U_{\rho}}$:

$$\mathbb{E}_{U_{\rho}}\left[\sum_{k\in U_{\rho}}\mathbb{E}\left\|\boldsymbol{\alpha}_{i,(t)}-\boldsymbol{\alpha}_{i,(t)}^{k}\right\|^{2}\right]$$

$$=K_{*}\sum_{k=1}^{K}q^{[k]}\cdot\mathbb{E}\left\|\boldsymbol{\alpha}_{i,(t)}-\boldsymbol{\alpha}_{i,(t)}^{k}\right\|^{2}$$

$$\stackrel{(1)}{\leq}2\eta^{2}\sum_{\rho=t_{c}}^{t-1}\left(\frac{2\sigma_{\psi}^{2}}{B^{k}}+\frac{2\sigma_{\alpha}^{2}}{I^{2}}\right)(K_{*}+1)+2\eta^{2}(t-t_{c})\left(K_{*}+1\right)\sum_{k=1}^{K}\sum_{\rho=t_{c}}^{t-1}q^{[k]}\cdot\left\|\nabla_{\boldsymbol{\alpha}_{i}^{k}}F^{k}(\boldsymbol{\alpha}_{\rho}^{k},\Psi^{k})\right\|^{2}$$

where (1) is because we sample client set U_{ρ} uniformly at random where client k is sampled with probability $q^{[k]}$ for $1 \le k \le K$ with replacement, and we define $U_{\rho} = \{k_1, ..., k_b, ..., k_{K_*}\}$, then

$$\begin{split} \begin{bmatrix} 1317 \\ 1318 \\ 1319 \\ 1320 \\ 1320 \\ 1321 \\ 1322 \\ 1323 \\ 1324 \\ 1325 \\ 1326 \\ 1326 \\ 1326 \\ 1327 \\ 1328 \\ \end{bmatrix} & = \sum_{b=1}^{K_*} \sum_{k=1}^{K} q^{[k]} \cdot \left\| \nabla_{\boldsymbol{\alpha}_i^k} F^k(\boldsymbol{\alpha}_{\rho}^k, \Psi^k) \right\|^2 \\ & = \sum_{b=1}^{K_*} \sum_{k=1}^{K} q^{[k]} \cdot \left\| \nabla_{\boldsymbol{\alpha}_i^k} F^k(\boldsymbol{\alpha}_{\rho}^k, \Psi^k) \right\|^2 \\ & = K_* \sum_{k=1}^{K} q^{[k]} \cdot \left\| \nabla_{\boldsymbol{\alpha}_i^k} F^k(\boldsymbol{\alpha}_{\rho}^k, \Psi^k) \right\|^2 \end{aligned}$$

Thus:

$$\begin{split} & \sum_{k=1}^{K} q^{[k]} \cdot \mathbb{E} \left\| \mathbf{\alpha}_{i,(t)} - \mathbf{\alpha}_{i,(t)}^{k} \right\|^{2} \\ & \leq 2\eta^{2} \sum_{\rho=t_{c}}^{t-1} \left(\frac{2\sigma_{\psi}^{2}}{B^{k}} + \frac{2\sigma_{\alpha}^{2}}{I^{2}} \right) \left(1 + \frac{1}{K_{*}} \right) + 2\eta^{2} (t - t_{c}) \left(1 + \frac{1}{K_{*}} \right) \sum_{k=1}^{K} \sum_{\rho=t_{c}}^{t-1} q^{[k]} \cdot \left\| \nabla_{\mathbf{\alpha}_{i}^{k}} F^{k}(\mathbf{\alpha}_{\rho}^{k}, \Psi^{k}) \right\|^{2} \end{split}$$

We take
$$v = \rho - t_c$$
, $\gamma = t - 1 - t_c$, and add iteration on both sides:

$$\begin{aligned} & \begin{array}{ll} 1338 & & \sum_{t=0}^{T-1} \sum_{k=1}^{K} q^{[k]} \cdot \mathbb{E} \left\| \mathbf{\alpha}_{i,(t)} - \mathbf{\alpha}_{i,(t)}^{k} \right\|^{2} \\ 1341 & & \\ 1342 & & \\ 1342 & & \\ 1343 & & \\ 1344 & \\ 1344 & \\ 1344 & \\ 1344 & \\ 1345 & & \\ 1346 & & \\ & & \\ \sum_{s=0}^{\lfloor \frac{T-1}{E} \rfloor} \sum_{\gamma=0}^{E-1} 2\eta^{2} \sum_{\upsilon=0}^{\gamma} \left(\frac{2\sigma_{\psi}^{2}}{B^{k}} + \frac{2\sigma_{\alpha}^{2}}{I^{2}} \right) \left(1 + \frac{1}{K_{*}} \right) \\ 1348 & & \\ 1348 & & \\ 1349 & & + \sum_{s=0}^{\lfloor \frac{T-1}{E} \rfloor} \sum_{\gamma=0}^{E-1} 2\eta^{2} (\gamma+1) \left(1 + \frac{1}{K_{*}} \right) \sum_{k=1}^{K} \sum_{\upsilon=0}^{\gamma} q^{[k]} \cdot \left\| \nabla_{\mathbf{\alpha}_{i}^{k}} F^{k}(\mathbf{\alpha}_{sE+\upsilon}^{k}, \Psi^{k}) \right\|^{2} \end{aligned}$$

$$\begin{array}{ll} 1350 \\ 1351 \\ 1351 \\ 1352 \\ 1352 \\ 1352 \\ 1353 \\ 1354 \\ 1355 \\ 1355 \\ 1356 \\ 1356 \\ 1356 \\ 1356 \\ 1356 \\ 1356 \\ 1357 \\ 1358 \\ 1358 \\ 1359 \\$$

where (1) use $0 \le \gamma \le E - 1$ and $1 + ... + E = \frac{E(E+1)}{2}$. Finally multiply both sides simultaneously by $\frac{1}{T}$:

$$\frac{1}{T} \sum_{t=0}^{T-1} \sum_{k=1}^{K} q^{[k]} \cdot \mathbb{E} \left\| \boldsymbol{\alpha}_{i,(t)} - \boldsymbol{\alpha}_{i,(t)}^{k} \right\|^{2} \\
\leq (E+1)\eta^{2} \left(\frac{2\sigma_{\psi}^{2}}{B^{k}} + \frac{2\sigma_{\alpha}^{2}}{I^{2}} \right) \left(1 + \frac{1}{K_{*}} \right) + \frac{2\eta^{2}E^{2}}{T} \left(1 + \frac{1}{K_{*}} \right) \sum_{k=1}^{K} \sum_{t=0}^{T-1} q^{[k]} \cdot \left\| \nabla_{\boldsymbol{\alpha}_{i}^{k}} F^{k}(\boldsymbol{\alpha}_{t}^{k}, \Psi^{k}) \right\|^{2} \\
\stackrel{(1)}{\leq} (E+1)\eta^{2} \left(\frac{2\sigma_{\psi}^{2}}{B^{k}} + \frac{2\sigma_{\alpha}^{2}}{I^{2}} \right) \left(1 + \frac{1}{K_{*}} \right) + \frac{2\eta^{2}E^{2}\lambda}{T} \left(1 + \frac{1}{K_{*}} \right) \sum_{t=0}^{T-1} \left\| \sum_{k=1}^{K} q^{[k]} \cdot \nabla_{\boldsymbol{\alpha}_{i}^{k}} F^{k}(\boldsymbol{\alpha}_{t}^{k}, \Psi^{k}) \right\|^{2}$$

$$\begin{array}{l} 1370 \\ 1371 \\ 1372 \\ 1373 \\ 1373 \\ 1374 \end{array} = (E+1)\eta^2 \left(\frac{2\sigma_{\psi}^2}{B^k} + \frac{2\sigma_{\alpha}^2}{I^2}\right) \left(1 + \frac{1}{K_*}\right) + \frac{2\eta^2 E^2 \lambda}{T} \left(1 + \frac{1}{K_*}\right) \sum_{t=0}^{T-1} \left\| \nabla_{\boldsymbol{\alpha}_i^k} F(\boldsymbol{\alpha}_t^k, \Psi^k) \right\|^2 \\ 1373 \\ 1374 \end{array}$$
 where (1) use Assumption 3. \Box

where (1) use Assumption 3.

Theorem 1. Suppose Assumption 1, 2 and 3 hold, for $t = 0, 1, ..., T - 1, B = \min \{B^1, ..., B^K\}$ where B^k is the local mini-batch size. $\alpha_{i,j} \ge \nu > 0$ for i = 1, ..., n and j = 1, ..., N, I is the sampling times for prompt. $\sigma_{\alpha}^2 = \frac{8G^2N}{\tau^2\nu^2}$ is the variance of the variance-reduced policy gradient, σ_{ψ}^2 is the variance of the stochastic gradient, $\mathbb{E}_{\Phi^k \sim GS(\boldsymbol{\alpha}^k)} \left[\mathcal{L}(\Phi^k, \Psi^k) \right]$ is L-smooth for $\boldsymbol{\alpha}^k$ and $L = \frac{nGN(\tau+1)}{\tau^2\nu^2}$, and η satisfies the following inequality:

 $0 < \eta \le \eta^* = \frac{-\lambda L + \sqrt{\lambda^2 L^2 + 8L^2 E\left(1 + \frac{1}{K_*}\right)}}{\lambda^2 L^2 + 8L^2 E\left(1 + \frac{1}{K_*}\right)}$

$$8L^2E\left(1+\frac{1}{K_*}\right)$$

(38)

Then, the Fed-BDPL's full gradient $\nabla_{\alpha} F(\alpha_t, \Psi^k)$ satisfies the following inequality:

$$\frac{1}{T} \sum_{t=0}^{T-1} \left\| \nabla_{\boldsymbol{\alpha}} F(\boldsymbol{\alpha}_{t}, \Psi^{k}) \right\|^{2} \\ \leq \frac{4G}{\eta T} + \frac{2(E+1)L^{2}\eta^{2}n\sigma_{\psi}^{2}(1+\frac{1}{K_{*}}) + 2nL\eta\sigma_{\psi}^{2}}{B} + \frac{2(E+1)L^{2}\eta^{2}n\sigma_{\alpha}^{2}(1+\frac{1}{K_{*}}) + 2nL\eta\sigma_{\alpha}^{2}}{I^{2}}$$

Proof. According to Lemma 2:

$$\mathbb{E}_{\Phi_{t+1}\sim \mathrm{GS}(\boldsymbol{\alpha}_{t+1})} \left[\mathcal{L}(\Phi_{t+1}, \Psi) \right] - \mathbb{E}_{\Phi_t \sim \mathrm{GS}(\boldsymbol{\alpha}_t)} \left[\mathcal{L}(\Phi_t, \Psi) \right]$$

$$\leq \left\langle \nabla_{\boldsymbol{\alpha}} \mathbb{E}_{\Phi_t \sim \mathrm{GS}(\boldsymbol{\alpha}_t)} \left[\mathcal{L}(\Phi_t, \Psi) \right], \boldsymbol{\alpha}_{t+1} - \boldsymbol{\alpha}_t \right\rangle + \frac{L}{2} \|\boldsymbol{\alpha}_{t+1} - \boldsymbol{\alpha}_t\|^2$$

$$\sum_{n=1}^{n} \left[\left\langle \mathbf{x}_{t+1} - \mathbf{x}_{t+1} \right\rangle - \left\langle \mathbf{x}_{t+1} - \mathbf{x}_{t+1} \right\rangle - \left\langle \mathbf{x}_{t+1} - \mathbf{x}_{t+1} \right\rangle \right]$$

$$\leq \sum_{i=1}^{n} \left[\left\langle \nabla_{\boldsymbol{\alpha}_{i}} F(\boldsymbol{\alpha}_{t}, \Psi^{k}), -\eta \cdot \hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}} f^{*}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\rangle + \frac{L\eta^{2}}{2} \left\| \hat{\nabla}_{\boldsymbol{\alpha}_{i}^{k}} f^{*}(\boldsymbol{\alpha}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\|^{2} \right]$$

We take the expectations about $\{\Phi^{k,r}\}_{r=1}^{I}$, \mathcal{B}_{t}^{k} and U_{t} on both sides respectively:

$$\mathbb{E}_{\left\{\Phi^{k,r}\right\}_{r=1}^{I}}\mathbb{E}_{\mathcal{B}_{t}^{k}}\mathbb{E}_{U_{t}}\left\{\mathbb{E}_{\Phi_{t+1}\sim \mathsf{GS}(\boldsymbol{\alpha}_{t+1})}\left[\mathcal{L}(\Phi_{t+1},\Psi)\right]-\mathbb{E}_{\Phi_{t}\sim \mathsf{GS}(\boldsymbol{\alpha}_{t})}\left[\mathcal{L}(\Phi_{t},\Psi)\right]\right\}$$

For a):

where (1) use inequality $2\langle a,b\rangle = ||a||^2 + ||b||^2 - ||a-b||^2$; (2) use the convexity of ℓ_2 norm; (3) use L-smooth in Lemma 2. For b):

$$\begin{aligned}
\begin{aligned}
& \text{1458} & \text{(2)} & \mathbb{E}_{U_{t}} \left\{ \frac{L\eta^{2}}{2} \mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}} \mathbb{E}_{B_{t}^{k}} \left[\left\| \hat{\nabla}_{\mathbf{a}_{t}^{k}} f^{*}(\mathbf{a}_{t}^{k}, \mathcal{B}_{t}^{k}) - \nabla_{\mathbf{a}_{t}^{k}} F^{*}(\mathbf{a}_{t}^{k}, \Psi^{k}) \right\|^{2} \right] \right\} \\
& + \mathbb{E}_{U_{t}} \left[\frac{L\eta^{2}}{2} \left\| \nabla_{\mathbf{a}_{t}^{k}} F^{*}(\mathbf{a}_{t}^{k}, \Psi^{k}) \right\|^{2} \right] \\
& + \mathbb{E}_{U_{t}} \left\{ \frac{L\eta^{2}}{2} \mathbb{E}_{U_{t}} \left\| \nabla_{\mathbf{a}_{t}^{k}} F^{*}(\mathbf{a}_{t}^{k}, \Psi^{k}) \right\|^{2} \right] \\
& + \mathbb{E}_{U_{t}} \left\{ L\eta^{2} \mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}} \mathbb{E}_{B_{t}^{k}} \left[\left\| \hat{\nabla}_{\mathbf{a}_{t}^{k}} f^{*}(\mathbf{a}_{t}^{k}, \mathcal{B}_{t}^{k}) - \mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}} \hat{\nabla}_{\mathbf{a}_{t}^{k}} f^{*}(\mathbf{a}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\|^{2} \right] \right\} \\
& + \mathbb{E}_{U_{t}} \left\{ L\eta^{2} \mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}} \mathbb{E}_{B_{t}^{k}} \left[\left\| \mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}} \hat{\nabla}_{\mathbf{a}_{t}^{k}} f^{*}(\mathbf{a}_{t}^{k}, \mathcal{B}_{t}^{k}) - \mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}} \mathbb{E}_{B_{t}^{k}} \hat{\nabla}_{\mathbf{a}_{t}^{k}} f^{*}(\mathbf{a}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\|^{2} \right] \right\} \\
& + \mathbb{E}_{U_{t}} \left\{ L\eta^{2} \mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}} \mathbb{E}_{B_{t}^{k}} \left[\left\| \mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}} \hat{\nabla}_{\mathbf{a}_{t}^{k}} f^{*}(\mathbf{a}_{t}^{k}, \mathcal{B}_{t}^{k}) - \mathbb{E}_{\{\Phi^{k,r}\}_{r=1}^{I}} \mathbb{E}_{B_{t}^{k}} \hat{\nabla}_{\mathbf{a}_{t}^{k}} f^{*}(\mathbf{a}_{t}^{k}, \mathcal{B}_{t}^{k}) \right\|^{2} \right] \right\} \\
& + \mathbb{E}_{U_{t}} \left\{ L\eta^{2} \mathbb{E}_{\{U_{t}} \left\| \nabla_{\mathbf{a}_{t}^{k}} F^{*}(\mathbf{a}_{t}^{k}, \Psi^{k}) \right\|^{2} + L\eta^{2} \cdot \mathbb{E}_{U_{t}} \left[\frac{\sigma^{2}}{I^{2}} \right] + L\eta^{2} \cdot \mathbb{E}_{U_{t}} \left[\frac{\sigma^{2}}{B} \right] \\
& + \mathbb{E}_{U_{t}} \left\{ \frac{L\eta^{2}}{2} \sum_{k=1}^{K} q^{[k]} \cdot \left\| \nabla_{\mathbf{a}_{t}^{k}} F^{k}(\mathbf{a}_{t}^{k}, \Psi^{k}) \right\|^{2} + \frac{L\eta^{2}\sigma^{2}}{I^{2}} + \frac{L\eta^{2}\sigma^{2}}{B} \\
& = \frac{L\eta^{2}}{2} \sum_{k=1}^{K} q^{[k]} \cdot \left\| \nabla_{\mathbf{a}_{t}^{k}} F^{k}(\mathbf{a}_{t}^{k}, \Psi^{k}) \right\|^{2} + \frac{L\eta^{2}\sigma^{2}}{I^{2}} + \frac{L\eta^{2}\sigma^{2}}{B} \\
& = \frac{L\eta^{2}}{2} \sum_{k=1}^{K} q^{[k]} \cdot \nabla_{\mathbf{a}_{t}^{k}} F^{k}(\mathbf{a}_{t}^{k}, \Psi^{k}) \right\|^{2} + \frac{L\eta^{2}\sigma^{2}}{I^{2}} + \frac{L\eta^{2}\sigma^{2}}{B} \\
& = \frac{L\eta^{2}}{2} \sum_{k=1}^{K} q^{[k]} \cdot \nabla_{\mathbf{a}_{t}^{k}} F^{k}(\mathbf{a}_{t}^{k}, \Psi^{k}) \right\|^{2} \\
& = \frac{L\eta^{2}}{2} \sum_{k=1}^{K} q^{[k]} \cdot \nabla_{\mathbf{a}_{t}^{k}} F^{k}(\mathbf{a}_{t}^{k}, \Psi^{k}) \\
& = \frac{L\eta^{2}}{2} \sum_{k=1}^{K} q^{[k]} \cdot \nabla_{\mathbf{a}_{t}^{k}} F^{k$$

where (1) use $\mathbb{E} \|x - \mathbb{E} [x]\|^2 = \mathbb{E} \|x\|^2 - \|\mathbb{E} [x]\|^2$; (2) use the unbiasedness of stochastic gradient and variance-reduced policy gradient in Assumption 1 and Lemma 1; (3) use $\|a + b\|^2 \le 2 \|a\|^2 + 2 \|b\|^2$; (4) use the bounded variance of stochastic gradient and variance-reduced policy gradient in Assumption 1 and Lemma 1; (5) use Assumption 3. Combining a) and b):

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \left\{ \mathbb{E}_{\Phi_{t+1}\sim GS(\boldsymbol{\alpha}_{t+1})} \left[\mathcal{L}(\Phi_{t+1}, \Psi) \right] - \mathbb{E}_{\Phi_t \sim GS(\boldsymbol{\alpha}_t)} \left[\mathcal{L}(\Phi_t, \Psi) \right] \right\} \\ &\leq \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^n \left\{ -\frac{\eta}{2} \left[\left\| \nabla_{\boldsymbol{\alpha}_i} F(\boldsymbol{\alpha}_t, \Psi^k) \right\|^2 + \left\| \nabla_{\boldsymbol{\alpha}_i^k} F(\boldsymbol{\alpha}_t^k, \Psi^k) \right\|^2 \right] + \frac{\eta}{2} \left[\sum_{k=1}^K q^{[k]} L^2 \cdot \left\| \boldsymbol{\alpha}_{i,(t)} - \boldsymbol{\alpha}_{i,(t)}^k \right\|^2 \right] \right\} \\ &+ \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^n \left\{ \frac{L\eta^2 \lambda}{2} \left\| \sum_{k=1}^K q^{[k]} \cdot \nabla_{\boldsymbol{\alpha}_i^k} F^k(\boldsymbol{\alpha}_t^k, \Psi^k) \right\|^2 + \frac{L\eta^2 \sigma_{\boldsymbol{\alpha}}^2}{I^2} + \frac{L\eta^2 \sigma_{\boldsymbol{\psi}}^2}{B} \right\} \\ & \left(\stackrel{(1)}{\leq} -\frac{\eta}{2} \frac{1}{T} \sum_{t=0}^{T-1} \left\| \nabla_{\boldsymbol{\alpha}} F(\boldsymbol{\alpha}_t, \Psi^k) \right\|^2 + \left[-\frac{\eta}{2} + \frac{\lambda L\eta^2}{2} + \eta^3 L^2 E(1 + \frac{1}{K_*}) \right] \frac{1}{T} \sum_{t=0}^{T-1} \left\| \nabla_{\boldsymbol{\alpha}^k} F(\boldsymbol{\alpha}_t^k, \Psi^k) \right\|^2 \\ &+ \frac{(E+1)L^2 \eta^3 n \sigma_{\boldsymbol{\psi}}^2(1 + \frac{1}{K_*}) + nL\eta^2 \sigma_{\boldsymbol{\psi}}^2}{B} + \frac{(E+1)L^2 \eta^3 n \sigma_{\boldsymbol{\alpha}}^2(1 + \frac{1}{K_*}) + nL\eta^2 \sigma_{\boldsymbol{\alpha}}^2}{I^2} \end{aligned}$$

where (1) use Lemma 4. Then,we can get:

$$\begin{split} &\frac{1}{T}\sum_{t=0}^{T-1} \left\| \nabla_{\alpha} F(\alpha_{t}, \Psi^{k}) \right\|^{2} \\ &\leq \frac{\mathbb{E}_{\Phi_{0} \sim \mathrm{GS}(\alpha_{0})} \left[\mathcal{L}(\Phi_{0}, \Psi) \right] - \mathbb{E}_{\Phi_{T} \sim \mathrm{GS}(\alpha_{t})} \left[\mathcal{L}(\Phi_{T}, \Psi) \right]}{\eta T} \\ &+ \left[-1 + \lambda L \eta + 2\eta^{2} L^{2} E(1 + \frac{1}{K_{*}}) \right] \frac{1}{T} \sum_{t=0}^{T-1} \left\| \nabla_{\alpha^{k}} F(\alpha_{t}^{k}, \Psi^{k}) \right\|^{2} \\ &+ \frac{2(E+1)L^{2} \eta^{2} n \sigma_{\psi}^{2}(1 + \frac{1}{K_{*}}) + 2nL \eta \sigma_{\psi}^{2}}{B} + \frac{2(E+1)L^{2} \eta^{2} n \sigma_{\alpha}^{2}(1 + \frac{1}{K_{*}}) + 2nL \eta \sigma_{\alpha}^{2}}{I^{2}} \end{split}$$

 $-1 + \lambda L\eta + 2\eta^2 L^2 E(1 + \frac{1}{\kappa}) \le 0$ $0 < \eta \le \eta^* = \frac{-\lambda L + \sqrt{\lambda^2 L^2 + 8L^2 E\left(1 + \frac{1}{K_*}\right)}}{8L^2 E\left(1 + \frac{1}{K_*}\right)}$ Finally: $\frac{1}{T}\sum_{i=1}^{T-1}\left\|\nabla_{\pmb{\alpha}}F(\pmb{\alpha}_t,\Psi^k)\right\|^2$ $\leq \frac{4G}{nT} + \frac{2(E+1)L^2\eta^2 n \sigma_{\psi}^2 (1+\frac{1}{K_*}) + 2nL\eta \sigma_{\psi}^2}{B} + \frac{2(E+1)L^2\eta^2 n \sigma_{\alpha}^2 (1+\frac{1}{K_*}) + 2nL\eta \sigma_{\alpha}^2}{I^2}$ **Corollary 1. Convergence rate and complexity** 1) Convergence rate: Let $\eta = \min \left\{ \eta^*, \frac{1}{\sqrt{T}}, \frac{1}{L} \right\}$. Under this condition, the following holds: $\frac{1}{T} \sum_{i=1}^{T-1} \left\| \nabla_{\boldsymbol{\alpha}} F(\boldsymbol{\alpha}_t, \Psi) \right\|^2 = \mathcal{O}(\frac{1}{\sqrt{T}})$ (40)2) Complexity: To guarantee an ϵ -solution, such that $\frac{1}{T}\sum_{t=0}^{T-1} \left\| \nabla_{\alpha} F(\alpha_t, \Psi^k) \right\|^2 \leq \epsilon^2$, the follow-ing condition must hold: $T_{\epsilon} = \mathcal{O}\left(\frac{1}{\epsilon^4}\right)$ (41)*Proof.* From condition of η : $0 < \eta \le \eta^* = \frac{-\lambda L + \sqrt{\lambda^2 L^2 + 8L^2 E\left(1 + \frac{1}{K_*}\right)}}{8L^2 E\left(1 + \frac{1}{K_*}\right)}$ We let $\eta = \min\left\{\eta^*, \frac{1}{\sqrt{T}}, \frac{1}{L}\right\}$ and the following holds: $\frac{1}{T}\sum_{i=1}^{T-1} \left\| \nabla_{\boldsymbol{\alpha}} F(\boldsymbol{\alpha}_t, \Psi^k) \right\|^2$ $\leq \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)(4G) + \mathcal{O}(\frac{1}{\sqrt{T}})\left(\frac{2nL\sigma_{\psi}^{2}}{B} + \frac{2nL\sigma_{\alpha}^{2}}{I^{2}}\right)$ $+\mathcal{O}\left(\frac{1}{\sqrt{T}}\right)\left[\frac{2(E+1)Ln\sigma_{\psi}^{2}(1+\frac{1}{K_{*}})}{B}+\frac{2(E+1)Ln\sigma_{\alpha}^{2}(1+\frac{1}{K_{*}})}{I^{2}}\right]$ $=\mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$ In addition, we considering the iterative complexity, to get a ϵ -solution: $\frac{1}{T}\sum_{\boldsymbol{\alpha}}^{T-1}\left\|\nabla_{\boldsymbol{\alpha}}F(\boldsymbol{\alpha}_t,\Psi^k)\right\|^2 \leq \epsilon^2$

Based on Assumption 2, $\mathbb{E}_{\Phi_0 \sim GS(\boldsymbol{\alpha}_0)} [\mathcal{L}(\Phi_0, \Psi)] - \inf_t \mathbb{E}_{\Phi_t \sim GS(\boldsymbol{\alpha}_t)} [\mathcal{L}(\Phi_t, \Psi)] \leq 2G$ and:

We need to choose:

$$T_{\epsilon} = \left[\frac{4G}{\epsilon^2} + \frac{2(E+1)Ln\sigma_{\psi}^2(1+\frac{1}{K_*}) + 2nL\sigma_{\psi}^2}{B\epsilon^2} + \frac{2(E+1)Ln\sigma_{\alpha}^2(1+\frac{1}{K_*}) + 2nL\sigma_{\alpha}^2}{I^2\epsilon^2} \right]^2$$

1566
1567
$$T_{\epsilon} = \mathcal{O}\left(\frac{1}{\epsilon^4}\right)$$
1568
1569

Corollary 2. The impact of K_* (FedOne): In the FL framework, $T_{\epsilon}K_*$ represents the total number of queries made to the cloud-based LLM service to achieve an ϵ -solution, whose quantity is directly proportional to the cost incurred for utilizing the LLM. The following condition holds:

$$T_{\epsilon}K_* \propto K_* \tag{42}$$

Therefore, $T_{\epsilon}K_*$ is a function of K_* that increase steadily for $K_* = 1, 2, \ldots, K$, indicating that the optimal K_* for query efficiency is $K_* = 1$.

Proof. According to Corollary 1, when $\eta = \min \left\{ \eta^*, \frac{1}{\sqrt{T}}, \frac{1}{L} \right\}$, to get a ϵ -solution:

$$\frac{1}{T}\sum_{t=0}^{T-1} \left\|\nabla_{\boldsymbol{\alpha}}F(\boldsymbol{\alpha}_t, \Psi^k)\right\|^2 \le \epsilon^2$$
(43)

We need to choose:

$$T_{\epsilon} = \left[\frac{4G}{\epsilon^2} + \frac{2(E+1)Ln\sigma_{\psi}^2(1+\frac{1}{K_*}) + 2nL\sigma_{\psi}^2}{B\epsilon^2} + \frac{2(E+1)Ln\sigma_{\alpha}^2(1+\frac{1}{K_*}) + 2nL\sigma_{\alpha}^2}{I^2\epsilon^2}\right]^2$$

The $T_{\epsilon}K_{*}$ proportional to the query time satisfies the following equality:

$$T_{\epsilon}K_{*}$$

$$= \left[\frac{4G}{\epsilon^2} + \frac{2(E+1)Ln\sigma_{\psi}^2(1+\frac{1}{K_*}) + 2nL\sigma_{\psi}^2}{B\epsilon^2} + \frac{2(E+1)Ln\sigma_{\alpha}^2(1+\frac{1}{K_*}) + 2nL\sigma_{\alpha}^2}{I^2\epsilon^2}\right]^2 \cdot K_*$$

$$= \left[\left(\frac{4G}{\epsilon^2} + \frac{2(E+2)Ln\sigma_{\psi}^2}{B\epsilon^2} + \frac{2(E+2)Ln\sigma_{\alpha}^2}{I^2\epsilon^2}\right) \cdot \sqrt{K_*} + \left(\frac{2(E+1)Ln\sigma_{\psi}^2}{B\epsilon^2} + \frac{2(E+1)Ln\sigma_{\alpha}^2}{I^2\epsilon^2}\right) \cdot \frac{1}{\sqrt{K_*}}\right]^2$$

Then, $T_{\epsilon}K_{*}$ is a function of K_{*} that first decreases and then increases and that the optimal K_{*} for query efficiency exists:

$$K_*^{opt} = \frac{\frac{2(E+1)Ln\sigma_{\psi}^2}{B\epsilon^2} + \frac{2(E+1)Ln\sigma_{\alpha}^2}{I^2\epsilon^2}}{\frac{4G}{\epsilon^2} + \frac{2(E+2)Ln\sigma_{\psi}^2}{B\epsilon^2} + \frac{2(E+2)Ln\sigma_{\alpha}^2}{I^2\epsilon^2}} < 1$$

Finally, consider only the effect of K_* and $K_* \ge 1$, we can get:

$$T_{\epsilon}K_* \propto K_*$$

Therefore, $T_{\epsilon}K_*$ is a function of K_* that increase steadily for $K_* = 1, 2, \ldots, K$, indicating that the optimal K_* for query efficiency is $K_* = 1$.

Remark 4. We carefully analyze the optimal K_* by balancing convergence complexity and query time, aiming to achieve the fastest convergence with the fewest queries. Specifically, we fix the convergence accuracy, denoted as ϵ , and determine the number of iterations, T_{ϵ} , required to achieve this accuracy. We then analyze the total number of queries, given by $T_{\epsilon}K_* \propto K_*$. This analysis shows that increasing K_* for $K_* = 1, 2, \dots, K$, the number of activated clients can accelerate convergence by leveraging more data; however, it also increases query overhead, leading to higher communication and computational costs. We rigorously demonstrate that the number of queries required for Fed-BDPL to achieve an ϵ -solution is proportional to K_* . This result emphasizes that optimal query efficiency is achieved when only a single client ($K_* = 1$) is activated per round.

1620 B EXPERIMENT DETAILS

1622 B.1 GLUE DATASET DETAILS

Below is the table for the detail of the dataset we use, and their associated tasks, metrics, and data domains.

Dataset	$ \mathbf{L} $	Train	Dev	Test	Туре	Metrics	Domain
MNLI	3	393K	9.8K	9.8K	NLI	acc.	fiction, reports
QQP	2	364K	40K	391K	paraphrase	F1	Quora
SST-2	2	6.7K	872	1.8K	sentiment	acc.	movie reviews
MRPC	2	3.7K	408	1.7K	paraphrase	F1	news
CoLA	2	8.6K	1K	1K	acceptability	Matthews corr.	books, articles
QNLI	2	105K	5.5K	5.5K	NLI	acc.	Wikipedia
RTE	2	2.5K	277	3K	NLI	acc.	news, Wikipedia

Table 6: The statistics and metrics of seven datasets in GLUE benchmark, |L|: number of classes for classification tasks.