PERSONALIZED FEDERATED FINE-TUNING FOR HET EROGENEOUS DATA: A TWO-LEVEL LOW RANK ADAPTATION APPROACH

Anonymous authors

Paper under double-blind review

ABSTRACT

We study the personalized federated fine-tuning task with heterogeneous client data in the context of foundation models, where clients collaboratively finetune a foundation model (e.g., BERT, GPT) without sharing their local data, achieving personalized models simultaneously. While recent efforts have applied parameter-efficient fine-tuning techniques like low-rank adaptation (LoRA) or training prompts in federated settings, they often overlook data heterogeneity and model personalization. The primary challenge is that a single common adapter or prompt learner may not suffice for the diverse data of all clients. To address this issue, we propose PF2LoRA, a new personalized federated fine-tuning algorithm based on a novel two-level low rank adaptation framework on top of LoRA. Given the pretrained foundation model whose weight is frozen, our algorithm aims to learn two levels of adaptation simultaneously: the first level aims to learn a common adapter for all clients, while the second level fosters individual client personalization. This framework explicitly accommodates variations in adapter matrix ranks across clients and introduces minimal additional memory overhead, as the second-level adaptation comprises a small number of parameters compared to the first level. Our experiments on natural language understanding and generation tasks demonstrate that PF2LoRA significantly outperforms existing federated fine-tuning methods.

030 031 032

033

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

1 INTRODUCTION

034 Federated learning (FL) (McMahan et al., 2017a; Kairouz et al., 2021) has emerged as a crucial 035 paradigm for enabling collaborative training of machine learning models across distributed clients while preserving data privacy (McMahan et al., 2017b; Geyer et al., 2017). FL is particularly impor-037 tant in some scenarios that involve sensitive data, such as healthcare (Brisimi et al., 2018; Sheller 038 et al., 2020), finance (Yang et al., 2019), and mobile devices (Bonawitz et al., 2019). However, in the context of foundation models like BERT (Devlin et al., 2018) and GPT (Radford et al., 2018), traditional FL algorithms face significant challenges due to the complexity of these models. It re-040 quires huge computing resources when directly fine-tuning model parameters on the heterogeneous 041 data distributed across different clients. 042

To address the issue of fine-tuning foundation models, many parameter-efficient fine-tuning (PEFT) methods such as prompt tuning (Lester et al., 2021) and low-rank adaptation (LoRA) (Hu et al., 2021) have been explored, where LoRA freezes the original pre-trained parameters $W \in \mathbb{R}^{m \times n}$ of the foundation model while fine-tuning additional low rank matrices $B \in \mathbb{R}^{m \times r}$ and $A \in \mathbb{R}^{r \times n}$, $r \ll \min(m, n)$. This technique enables fine-tuning large models with a reduced number of trainable parameters, making them more suitable for resource-constrained devices. This paper specifically focuses on LoRA in the context of federated learning for heterogeneous data.

A natural method to perform low rank adaptation in federated learning is to adopt the same rank r of matrices A and B across different clients. This method is referred to as homogeneous LoRA (HOMLoRA), but it does not accommodate the personalized requirement of clients with heterogeneous data distributions. Recent work HETLoRA (Cho et al., 2024) highlights the importance of heterogeneous rank configurations to enable personalized federated learning, which proposed

Client 1 n \mathbf{h} Client 2 Client M Second level First level Pretrained LLM adapter adapter $D_k = 0$ retrained Veights $W \in \mathbb{R}$ $r_{\rm I}$ B_{i} $\mathcal{N}(0,$ Frozen Fine-tuning Fine-tuning ¢ Client 3 Client M Х The second (a) Locally fine-tune for Client k(b) Communication

Figure 1: Illustration of the two-level low-rank adaptation framework. The first level learns a common adapter $\{A, B\}$ for all clients, and the common adapter is synchronized by averaging across all the clients at every communication round. The second level aims to learn a client-specific and lightweight adapter $\{C_k, D_k\}$ for a specific client $k \in [1, M]$, while the lightweight adapters introduce negligible additional memory overhead.

077 "matrix truncation", "local rank self-pruning", and "sparsity-weighted aggregation" to learn various 078 ranks r_k for the heterogeneous data from clients. However, this approach suffers from two main 079 drawbacks: (1) The initial rank for any client is fixed and in the range of predefined minimal and maximal ranks, which is independent of client data. However, it is possible that the clients learning 081 difficult tasks are assigned with smaller ranks and do not have the capacity to learn their correspond-082 ing tasks well. (2) There are many hyperparameters which need to be tuned, including the minimal and maximal values of rank, the pruning parameter, and the sparsity parameter. It remains unclear 084 how to perform personalized federated fine-tuning such that the adapter is dependent on the data and 085 the procedure has a small number of tuning parameters.

In this paper, we propose PF2LoRA, a novel personalized federated fine-tuning algorithm that ex-087 plicitly incorporates heterogeneous ranks into the problem formulation. Our approach introduces a two-level low-rank adaptation framework on top of LoRA. The first level learns a common adapter for all clients $x = \{B \in \mathbb{R}^{m \times r}, A \in \mathbb{R}^{r \times \hat{n}}\}$, while the second level facilitates individual client personalization by learning client-specific and lightweight adapter y_k for k-th client, where 090 $y_k = \{D_k \in \mathbb{R}^{m \times \tilde{r}}, C_k \in \mathbb{R}^{\tilde{r} \times n}, 0 < \tilde{r} < r\}$ and $1 \le k \le M$. We formulate the proposed two-091 level low-rank adaptation framework as a bilevel optimization problem: we aim to learn a common 092 adapter x to minimize the loss function given the fact the individual client adapters $\{y_k\}_{k=1}^M$ can achieve the best performance when the common adapter x is given. The two-level LoRA frame-094 work explicitly accommodates variations in adapter matrix ranks across clients, which essentially 095 circumvents the rank pruning, matrix truncation, and zero-padding in HETLoRA for the alignment 096 of adapters. The introduced client-specific adapter actually enables the personalized adaptation for 097 heterogeneous data. Besides, the whole framework increases negligible additional memory over-098 head, as the second-level low rank adaptation comprises a small number of parameters compared to 099 the first level. Our main contribution is listed as follows:

100

054

055

056

059

060

061

062

063

064

065

067

068

069

076

102 103 • We propose a novel bilevel formulation for personalized fine-tuning on heterogeneous data, and develop a two-level low rank adaptation framework to efficiently fine-tune foundation model in the scenario of federated learning. The main workflow of our framework is illustrated in Figure 1.

Through extensive experiments on various natural language understanding and generation tasks, we demonstrate that PF2LoRA significantly outperforms existing federated fine-tuning baselines, providing a robust and efficient solution for personalized federated learning with foundation models. For example on GLUE benchmark, PF2LoRA achieves

25.6%, 2.33%, 15.34%, and 2.53% higher performance than state-of-the-art baseline HET-LoRA on MNLI, SST-2, QQP, QNLI dataset, respectively. In addition, through extensive ablation studies, we show that our proposed two-level adaptation framework achieves the highest performance across various data heterogeneity levels and outperforms baseline methods even if they use more trainable parameters.

116

108

110

111

2 RELATED WORK

Parameter-efficient Fine-Tuning. There are various categories of parameter-efficient fine-tuning 117 (PEFT) techniques, where only a subset of parameters of the pretrained foundation model or a small 118 number of additional parameters are updated to adapt to the target task. The first line of work 119 includes bias update or head-tuning (Lee et al., 2019; Zaken et al., 2021; Lawton et al., 2023; Wei 120 et al., 2021) and weight masking (Zhao et al., 2020; Sung et al., 2021; Xu et al., 2021). The second 121 line of work considers adapters (Houlsby et al., 2019; He et al., 2021a), prompt tuning (Lester 122 et al., 2021; Li & Liang, 2021) and low rank matrix adaptation (Hu et al., 2021). Different from 123 these works, our paper focuses on designing new federated learning algorithms based on low rank 124 adaptation with heterogeneous data, where the local client data is decentralized and not shared to 125 other clients.

126 Federated Learning with Fine-tuning. The PEFT framework has been recently incorporated in 127 the FL framework (Babakniya et al., 2023; Zhang et al., 2024; 2023b; Cho et al., 2024; Wang et al., 128 2023). However, most of them do not consider the data heterogeneity in the context of foundation 129 models. To the best of our knowledge, HETLoRA (Cho et al., 2024) is the only work which allows 130 data-independent heterogeneous ranks for each clients by a fixed rank initialization, zero-padding, 131 truncation, self-pruning and sparsity regularization. In contrast, our work promotes data-dependent heterogeneous ranks of local clients by an explicit bilevel modeling and reduce the number of tuning 132 hyperparameters. 133

134 135

3 PRELIMINARIES

136 137 138

139

140

147

148

In this section, we introduce a few parameter-efficient fine-tuning methods in the context of (federated) foundation model learning. It includes LoRA, HOMLORA, HETLORA. Due to limited space, we describe a variant of the personalized federated learning algorithm in the context of low rank adaptation method, namely Per-FedAvg-LoRA, in Appendix A.

Low-rank adaptation (LoRA). LoRA is a technique designed to efficiently fine-tune large pretrained models by injecting trainable low-rank matrices into each layer of a foundation model (Hu et al., 2021). Formally, consider a pre-trained model where the original weight matrix is denoted as $W_0 \in \mathbb{R}^{m \times n}$. The model update ΔW during the fine-tuning can be approximated by multiplication of two low-rank matrices $B \in \mathbb{R}^{m \times r}$ and $A \in \mathbb{R}^{r \times n}$. The updated weight matrix W is then given by:

$$W = W_0 + \Delta W = W_0 + BA. \tag{1}$$

This decomposition allows the model to learn adaptations for down-stream tasks while keeping the
 majority of the original weights frozen, thereby maintaining the pre-trained knowledge and signifi cantly reducing memory and computational overhead.

152 **HOMLoRA**. When considering LoRA in the scenario of federated learning, a natural extension is 153 refereed to as HOMLoRA, which adopts homogeneous rank r across all the clients. Assume that 154 M clients participate in federated learning at every communication round. The objective function of each client k is $f_k(\cdot)$, and the goal is to find a common adapter $x = \{A \in \mathbb{R}^{m \times r}, B \in \mathbb{R}^{r \times n}\}$ that 155 performs well across all the clients. It aims to solve the optimization problem: $\min_x \frac{1}{M} \sum_{k=1}^{M} f_k(x)$. Specifically, each client locally updates their adapters for I steps by Adam (or SGD) using their 156 157 local data, and the server aggregates the adapters from each local clients $\{A_k^t, B_k^t\}_{k=1}^{M}$ (k is the 158 local client id) at iteration t when t is a multiple of I, where I is the number of local updates per round: $A^t = \frac{1}{M} \sum_{k=1}^{M} A_k^t$, $B^t = \frac{1}{M} \sum_{k=1}^{M} B_k^t$. Then the server broadcasts the aggregated adapters back to each client. HoMLoRA can be regarded as a direct extension of FedAvg (McMahan et al., 159 160 161 2017a) in the context of LoRA (Hu et al., 2021).

HETLoRA. Recently, Cho et al. (2024) proposed a heterogeneous LoRA method, namely HET-LoRA, which is able to learn heterogeneous low rank matrices for different clients. The main technical components contain four parts: (1) a fixed rank initialization: where the r_k is fixed for k-th client and $r_{min} \le r_k \le r_{max}$; (2) distribution via truncation, wherein at each communication round, each client truncates the global adapter matrices to align dimensions $A_k^t = A_{irk,..}^t, B_k^t = B_{i..rk}^t$;

(3) local training with self-pruning, which introduces the regularization term (with a penalty coefficient λ) to induce adapter sparsity (with a sparsity factor γ), and it dynamically reduces the r_k by pruning unimportant columns in B_k^t (or rows in A_k^t); (4) sparsity-weighted aggregation, wherein each communication round, to aggregate the adapter matrices with different rank $r_{min} \leq r_k \leq r_{max}$, the server reconstructs $\{A_k^t, B_k^t\}$ by zero-padding them to r_{max} .

Then HETLoRA updates the common adapter by aggregating the local adapters with an aggregation weight. In particular, the update rule is $A^{t+1} = \sum_{k=1}^{M} \|\Delta W_k^t\| A_k^t / \sum_{k=1}^{M} \|\Delta W_k^t\|$ and $B^{t+1} = \sum_{k=1}^{M} \|\Delta W_k^t\| B_k^t / \sum_{k=1}^{M} \|\Delta W_k^t\|$, $\Delta W_k^t = B_k^t A_k^t$.

176 However, the performance of HETLoRA heavily depends on (1) the fixed rank initialization, which 177 is independent of data and may cause underfitting or overfitting issues, and (2) the proper setting for 178 a set of hyperparameters, including r_{\min} , r_{\max} , γ , and λ .

To address these issues, we propose a new two-level low-rank adaptation framework for personalized fine-tuning in the next subsection.

4 A NEW TWO-LEVEL ADAPTATION FOR PERSONALIZED FEDERATED FINE-TUNING

As we discussed in Section 3, HOMLoRA uses only one common adapter $x = \{B \in \mathbb{R}^{m \times r}, A \in \mathbb{R}^{m \times r}\}$ 186 $\mathbb{R}^{r \times n}$ across all the clients, which is insufficient to learn from the heterogeneous data in federated 187 learning. Therefore, we introduce a client-specific adapter for every client k with $y_k = \{D_k \in D_k \in N\}$ 188 $\mathbb{R}^{m \times \tilde{r}}, C_k \in \mathbb{R}^{\tilde{r} \times n}, 0 < \tilde{r} < r, 1 \le k \le M$. We emphasize that the newly introduced adapter 189 has a much smaller rank \tilde{r} than that in the common adapter. Empirically, we usually set $\tilde{r} = \frac{r}{4}$ or 190 $\frac{r}{2}$, which means the trainable parameters in the client-specific adapter are only $\frac{1}{4}$ or $\frac{1}{2}$ of those in 191 the common adapter. Thus the new adapter is lightweight and incurs negligible additional memory 192 overhead. 193

Different from equation 1, we incorporate both the common and client-specific adapters. In particular, the adapter for the k-th client can be parameterized as,

$$W_k = W_0 + BA + D_k C_k,$$

(2)

where W_k is the adapter for k-th client, A, B are common adapters for all clients, and (C_k, D_k) are client-specific adapters for k-th client. Since the original weight W_0 is frozen, the trainable parameters in the model are A, B, C_k , D_k for the client k. Different than the HETLoRA whose local client matrix rank is predefined and independent of data, our specific parameterization equation 2 explicitly encourages each adapter W_k for the k-th client to vary over k: it can have different ranks in the range $(r - \tilde{r}, r + \tilde{r})$ and the specific rank is automatically determined by the training data.

We formalize our two-level adaptation framework for personalized federated fine-tuning as the following bilevel optimization problem:

206 207

208

196

181 182

183

185

$$\min_{x} \Phi(x) \coloneqq \frac{1}{M} \sum_{k=1}^{M} f_k(x, y_k^*(x)), \quad \text{(UL)}$$

s.t., $y_k^*(x) \in \arg\min_{y_k} f_k(x, y_k), \quad \text{(LL)}$
(3)

209 210

211 where $f_k(x, y_k) := \mathbb{E}_{\xi \sim \mathcal{D}_k} F_k(x, y_k; \xi)$ is the loss function for the k-th client, F_k the individual 212 loss function for a sample ξ from the k-th client, and \mathcal{D}_k is the data on client k. The upper-level 213 (UL) learns a common adapter x for all the clients upon a set of the best client-specific adapters 214 $\{y_k^*(x) \mid 1 \le k \le M\}$ for given x defined by the lower-level problem. Given the common adapter, 215 the lower-level (LL) aims to locally search the optimal client-specific adapter to fit its respective 214 data, which in fact fosters individual client personalization.

238

Alg	orithm 1 Two-level Adaptation for Personalized Fine-Tuning
1:	Input: $\alpha, \eta, I, T, M, \mathcal{D}_k$
2:	for $k \in \{1, \dots, M\}$ in parallel do
3:	for $t = 0, 1,, T - 1$ do
4:	Sample $\pi_k^t, \xi_k^t, \tilde{\xi}_k^t, \zeta_k^t$ independently from distribution \mathcal{D}_k
5:	$y_k^{t+1} = y_k^t - \alpha \nabla_y F_k(x_k^t, y_k^t; \pi_k^t)$
6:	$x_{L}^{t+1} = x_{L}^{t} - \eta \nabla_{x} F_{k}(x_{L}^{t}, y_{L}^{t+1}; \xi_{L}^{t}) + \alpha \eta \nabla_{xy} F_{k}(x_{L}^{t}, y_{L}^{t}; \zeta_{L}^{t}) \nabla_{y} F_{k}(x_{L}^{t}, y_{L}^{t+1}; \tilde{\xi}_{L}^{t})$
7:	$\inf_{k=0}^{k} t \% I == 0 \text{ then }$
8:	$x^{t+1} = \frac{1}{M} \sum_{k=1}^{M} x_k^{t+1}$
9:	$x_{l}^{t+1} = x^{t+1}$
10:	endif
11:	end for
12:	end for

Algorithm Design. Now we consider solving equation 3 efficiently in personalized federated learning. At the beginning of every communication round, i.e., (t% I = 0), each client k receives the averaged common adapter x_k^t from the server, and starts running its local steps. We run one step SGD for the lower-level problem to approximately find the minimizer of the lower-level problem (line 5 in Algorithm 1).

Define $\Phi_k(x) = f_k(x, y_k^*(x))$, then the gradient of the function $\Phi_k(x_k^t)$ in terms of x_k^t , namely hypergradient (Ghadimi & Wang, 2018), can be calculated by chain rule approximately as follows:

$$\nabla \Phi_k(x_k^t) \approx \nabla \widehat{\Phi}_k(x_k^t) = \nabla_x f_k(x_k^t, y_k^{t+1}) - \alpha \nabla_{xy} f_k(x_k^t, y_k^t) \nabla_y f_k(x_k^t, y_k^{t+1}), \tag{4}$$

where \approx is due to the fact that y_k^{t+1} is only an approximation to the optimal solution $y_k^*(x_k^t)$. Therefore, we use the stochastic version of $\nabla \widehat{\Phi}_k(x_k^t)$ to update the common adapter x on client k, as described in line 6 of Algorithm 1.

In fact, Adam or AdamW can also be used to update the upper-level variable based on the stochastic gradient information to replace the SGD update as in line 6. Empirically, we adopt AdamW as the upper-level optimizer (line 6) and SGD as the lower-level optimizer (line 5) to fine-tune a language model. One can refer to Algorithm 1 for more details, where line 5 is used to update the client-specific adapter, line 6 is used to update the common adapter, and line 8 corresponds to the synchronization of the common adapter.

Detailed Implementation for Language Models. Next, we present the details about how to implement our framework in large language models like RoBERTa (Liu et al., 2019), DeBERTa (He et al., 2020), and GPT-2 (Radford et al., 2019). All of these models are developed based on the transformer (Vaswani et al., 2017) architecture, which consists of multiple layers (e.g., *L*) of self-attention modules and feed-forward networks. The self-attention module contains query ($W_q \in \mathbb{R}^{m \times n}$), key ($W_k \in \mathbb{R}^{m \times n}$) and value ($W_v \in \mathbb{R}^{m \times n}$) weight matrices, where LoRA is typically applied to W_q and W_v .

When initializing the language model for each client, we first load the pretrained weights into the model and then locate all self-attention modules in every layer. For each W_0^i matrix in *i*-th selfattention layer, we initialize two new matrices $B^i \in \mathbb{R}^{m \times r}$, $A^i \in \mathbb{R}^{r \times n}$ (*r* is typically set as 8 or 4), and then merge the multiplication of matrices B^i , A^i and W_0^i to form new merged weight matrix $W_0^i + B^i A^i$. After the model finishes initialization, we freeze all the pretrained parameters, i.e., $W_0 = \{W_0^0, ..., W_0^{L-1}\}$ is frozen, and only matrices $B = \{B^0, ..., B^{L-1}\}$, $A = \{A^0, ..., A^{L-1}\}$ are allowed to be trained. Here $\{B, A\}$ matrices form the common adapter *x* in equation 3. The above description outlines the standard model initialization procedure of HOMLORA.

Now, let us discuss the differences in model initialization between our framework and HOMLoRA. In addition to the initialization procedures mentioned above, our framework requires two matrices: $D_k^i \in \mathbb{R}^{m \times \tilde{r}}$ and $C_k^i \in \mathbb{R}^{\tilde{r} \times n}$ (\tilde{r} is typically set as $\frac{r}{4}$ or $\frac{r}{2}$) to be initialized in *i*-th self-attention layer on client k. Then the new matrix can be parameterized as $W_0^i + B^i A^i + D_k^i C_k^i$. Note that $D_k = \{D_k^0, ..., D_k^{L-1}\}$ and $C_k = \{C_k^0, ..., C_k^{L-1}\}$ along with $\{B, A\}$ are trainable parameters of the model. Here $\{D_k, C_k\}$ matrices form the client-specific adapter y_k in equation 3. Refer to PyTorch-style pseudocode 2 in Appendix B for implementation details.

5 EXPERIMENTS

271 272

In this section, we first conduct extensive experiments with our algorithm and baselines on two major 273 natural language tasks, i.e., natural language understanding (NLU) and natural language generation 274 (NLG), where NLU experiments include the text classification on GLUE benchmark (Wang et al., 275 2018) and question answering task on SQuAD v1 (Rajpurkar et al., 2016) and v2 (Rajpurkar et al., 276 2018). NLG experiments are performed on E2E NLG Challenge dataset (Novikova et al., 2017) and WebNLG dataset (Gardent et al., 2017). Then we execute the ablation studies to explore (1) the 278 performance comparison when other baselines have more trainable parameters than ours in section 5.3; (2) the impact of data heterogeneity on PF2LoRA and baselines in Appendix F.1; and (3) the 279 importance of bilevel optimization in our framework in Appendix F.2. In addition, we analyze the 280 training stability of HETLoRA and PF2LoRA in Appendix G. We compare with four baselines, 281 including Centralized LoRA, Homogeneous LoRA (HOMLoRA), Personalized Federated Average 282 LoRA (Per-FedAvg-LoRA), and Heterogeneous LoRA (HETLoRA).

283 284 285

5.1 NATURAL LANGUAGE UNDERSTANDING

286 287 5.1.1 ROBERTA ON TEXT CLASSIFICATION

288 Model. In this section, we adopt RoBERTa (Devlin et al., 2018) model to perform the personalized 289 federated fine-tuning on the NLU task, i.e., text classification of GLUE benchmark. RoBERTa (Liu et al., 2019) is an enhancement of BERT (Devlin et al., 2018) designed to improve its performance 290 on general NLU tasks. It is commonly used and remains a competitive performance in NLU. We 291 take the pre-trained RoBERTa base (125M parameters) and RoBERTa large (355M parameters) 292 from HuggingFace library (Wolf et al., 2020) and apply the LoRA technique to the model. Other 293 baselines inject only common adapters into each attention-layer of the pretrained RoBERTa. In contrast, PF2LoRA injects common adapters and client-specific adapters into the pretrained model. 295 For baselines Centralized LoRA, HOMLora, and Per-FedAvg-LoRA, we initialize the rank $r_k = 8$ 296 across all the clients. HETLoRA initializes the client rank r_k , such that $r_{min} \leq r_1 \leq r_2 \leq$ 297 $\cdots \leq r_M \leq r_{max}$. In all the experiments, we tune the best r_{min} and r_{max} , and initially assign 298 $r_k = r_{min} + \frac{(r_{max} - r_{min})(k-1)}{M}$. PF2LoRA sets the rank r_k of the common adapter to 8 and the rank 299 \tilde{r}_k of the client-specific adapter to 2. The number of trainable parameters of RoBERTa base/large 300 corresponding to the initial rank are listed in Table 7 in Appendix C.1. We can observe that the 301 number of trainable parameters in PF2LoRA is slightly increased. Note that HETLoRA uses a 302 different rank for matrices on different clients, leading to different number of trainable parameters in 303 each client, we count the average trainable parameters of the clients. Experimental results regarding baselines with more trainable parameters will be discussed in Section 5.3. 304

Dataset. We follow the non-i.i.d. partitioning protocol in (Karimireddy et al., 2020) to split each dataset into heterogeneous client datasets with varying label distributions. Specifically, for a similarity parameter $s \in [0, 1]$, each client's local dataset is composed of two parts. The first $(100 \times s)\%$ is comprised of i.i.d. samples from the complete dataset, and the remaining $100 \times (1 - s)\%$ of data is sorted by label. We select five classification datasets for text classification, including CoLA, MNLI, SST-2, QQP, QNLI, from GLUE benchmark. The data summary information is presented in Table 8 in Appendix C.1.

312 Experiment Details. We run federated fine-tuning algorithms across 8 clients (NVIDIA A100 313 GPU), where all the clients participate in the training process, while the first client (with ID = 0) 314 also implements the parameter aggregation and distribution at every communication round. Central-315 ized LoRA, HOMLoRA and HETLoRA use the AdamW optimizer to update the common adapter. Per-FedAvg-LoRA adopts SGD to implement one-step update and AdamW to update the common 316 adapter. PF2LoRA uses SGD to update the client-specific adapter and AdamW to update the com-317 mon adapter. The learning rates for all methods are tuned and the best choices of learning rate for 318 each baseline can be found in Table 6 in Appendix C.1. For fair comparison, we keep the batch size 319 $\mathcal{B} = 16$, and communication interval I = 10 for all the federated baselines. The communication 320 rounds R are set according to the dataset size, {CoLA: 50, MNLI: 300, SST-2: 100, QQP: 300, 321 QNLI: 100}, and we keep the same R for all the baselines in a dataset. 322

We first execute federated fine-tuning on each client's training data and then evaluate the model on each client's test data. The final test results are the average of each client's result. We use

325	Table 1: Roberta-base results on GLUE benchmark. We report "Matthew's correlation" for CoLA
326	and "Accuracy" for MNLI, SST-2, QQP and QNLI. Higher value means "better performance".

una ricea	ilde j 101 111 (EI, 001 2	, रूरा धाव	Qi (Li, ingi	ier varae m	cuilo octie	i periorinan
	Method	CoLA	MNLI	SST-2	QQP	QNLI
	Centralized LoRA	56.85	83.48	93.58	86.97	89.70
	HOMLoRA	50.75	70.56	92.47	79.61	85.45
	Per-FedAvg-LoRA	51.11	74.73	90.56	81.26	78.59
	HETLoRA	53.76	73.33	93.67	81.49	91.86
	PF2LoRA	54.19	92.14	95.85	93.99	94.18

324

327 328

330 331

334 335

336

337

338

"Matthews's correlation" to measure the performance on CoLA and "Accuracy" to measure the performance on MNLI,SST-2, QQP, QNLI. The results are presented in Table 1 for RoBERTa base and Table 14 in Appendix D for RoBERTa large, where the heterogeneity level s = 0.3 is set for CoLA and s = 0.9 for MNLI, SST-2, QQP, and QNLI. PF2LoRA outperforms significantly all the baselines including Centralized LoRA on datasets MNLI, SST-2, QQP, QNLI. On dataset of CoLA, PF2LoRA performs better than other federated baselines and close to Centralized LoRA.

369

5.1.2 DEBERTA ON QUESTION ANSWERING

Model. DeBERTa (He et al., 2021b) is an enhanced transformer encoder. It improves the understanding capability for text compared to BERT and RoBERTa, and thus can perform better in more
sophisticated natural language tasks, such as question-answering, and sentiment analysis. We adopt
DeBERTa v3 with 86M parameters in the question-answering task SQuAD v1 and v2.

Dataset. SQuAD v1/v2 is a reading comprehension dataset consisting of 100k+/150k+(v1/v2)348 question-answering pairs extracted from Wikipedia articles. Each sample consists of a passage, a 349 question, and an answer, where the answers in SQuAD v1 can be derived from the given passage, 350 but SQuAD v2 includes some questions that do not have an answer in the passage, thus it serves as 351 a more challenging benchmark for reading comprehension. These questions involve a wide range of 352 topics, e.g., history, science and technology, geography and places. We construct the heterogeneous 353 data based on the question topics. There are 442 unique topics in the training set for both datasets, 354 but only 48 (35) topics for SQuAD v1 (v2) test set, and the topics in the training and test set are 355 totally different. To guarantee consistency of data distribution between the training and test set, we 356 uniformly sample 80% from the original training set as the new training set and regard the rest as the 357 test set. Then we use the same way mentioned in Section 5.1.1 to construct the heterogeneous data 358 with heterogeneity parameter s = 0.5. Exact Match (EM) score and F1 score are two commonly 359 used metrics to evaluate the quality of answers that models provide.

360 **Experiment Details**. Considering the complexity of the question-answering task, we run federated 361 fine-tuning across 4 clients (NVIDIA A100 GPU) with the same heterogeneity parameter s = 0.5, 362 communication rounds R = 200, communication interval I = 10 for SQuAD v1/v2. The optimizer 363 for different baselines follows the settings in Section 5.1.1. The batch size \mathcal{B} is fixed as 16 for all the baselines for fair comparison. The best learning rate settings for baselines are listed in Table 9 364 in Appendix C.2. The initial rank settings for all the baselines can be found in Table 10 in Appendix 365 C.2. The test results of PF2LoRA and other baselines are shown in Table 2. PF2LoRA exhibits the 366 highest EM score and F1 score among all federated baselines. For example, PF2LoRA outperforms 367 the best baseline by 4.08% in terms of EM score and 2.20% in terms of F1 score on SQuAD v1. 368

370	Table 2. Deber	ta-v3 results on t	SOuAD
371	1000 2. 00001	SOuAD v1.0	SOuAD v2.0
372	Method	(EM/F1)	(EM/F1)
374	Centralized LoRA ¹	68.72/83.36	44.56/53.31
375	HOMLoRA	68.57/82.99	42.53/52.70
376	Per-FedAvg-LoRA	68.80/83.08	43.15/53.16
377	HETLORA	68.64/83.28	44.53/54.69
	PF2LoRA	71.01/85.11	44.95/54.71

3783795.2 NATURAL LANGUAGE GENERATION

For NLG tasks, we follow LoRA (Hu et al., 2021) to use GPT-2 medium model for federated finetuning on WebNLG and E2E NLG Challenge dataset.

Model. GPT-2 (Radford et al., 2019) is an advanced language model developed by OpenAI. It builds on the success of the original GPT model, and has been widely applied in natural language understanding and generation. We use GPT-2 medium with 345M parameters and GPT2-XL with 1.5 Billion parameters for NLG tasks.

Dataset. WebNLG dataset is a benchmark for evaluating natural language generation systems. It fo-cuses on generating coherent and contextually relevant text from structured data (e.g., RDF triples). It includes various domains such as sports, cities, universities, hotels and more. E2E NLG Chal-lenge dataset is a NLG dataset especially focusing on restaurants domain. It emphasizes generating natural, human-like text from structured data (including attributes like restaurant name, food type, price range and rating). For WebNLG, we find that the text style and feature vary with the domains, so we construct the heterogeneous data based on the entry domains. There are 10 domains in the training set and test set. We split the domain into 8 (the number of clients) groups, and make sure that the domains of training and test set on a client are the same. E2E NLG Challenge dataset col-lects information of 34 restaurants in the training set and 18 restaurants in the test set. We split all the restaurants into 8 (the number of clients) groups by the name, and make sure that the restaurant names in the test set that a client receives are covered by its training set.

Table 3: GPT-2 generation results on WebNLG dataset.

	0			
Method	BLEU ↑	MET \uparrow	TER \downarrow	ROUGE-L↑
Centralized LoRA	0.6031	0.7807	0.5900	0.4169
HOMLoRA	0.5141	0.7271	0.5697	0.4736
Per-FedAvg-LoRA	0.5152	0.7219	0.5746	0.4740
HETLoRA	0.5196	0.7219	0.5746	0.4740
PF2LoRA	0.5261	0.7301	0.5733	0.4769

Table 4: GPT2-XL generation results on WebNLG dataset

Tuble 1. Of 12 AL generation results on webtilles dataset.					
Method	BLEU ↑	MET \uparrow	TER \downarrow	ROUGE-L↑	
HOMLoRA	0.5768	0.7771	0.6103	0.3967	
Per-FedAvg-LoRA	0.5783	0.7783	0.6157	0.3972	
HETLoRA	0.5763	0.7789	0.6164	0.3922	
PF2LoRA	0.5881	0.7832	0.6198	0.3978	

Table 5: The comparison results with more trainable parameters in baselines. We report "Matthew's correlation" for CoLA and "Accuracy" for MNLI, SST-2, QQP and QNLI. Higher value means "better performance".

Method	Initial Rank	# Parameters	CoLA	MNLI	SST-2	QQP	QNLI
HOMLoRA	$r_k = 12$	0.44M	52.01	73.82	92.63	80.11	86.27
Per-FedAvg-LoRA	$r_{k} = 12$	0.44M	52.35	78.62	89.65	81.12	81.41
HETLoRA	$r_{max} = 16, r_{min} = 8$	0.43M	53.43	79.32	94.83	81.71	92.12
PF2LoRA	$r_k = 8, \tilde{r}_k = 2$	0.37M	54.19	92.14	95.85	93.99	94.18

Experiment Details. We follow the procedures in LoRA (Hu et al., 2021) to implement language generation, including (1) fine-tuning the language model, (2) generating outputs for text data using beam search, (3) decoding the outputs, and (4) evaluating the generated outputs. The NLG experiments are run across 8 clients (NVIDIA A100 GPU), where each client fine-tunes the adapter on the

¹Note that the results do not exactly match the LoRA results reported in Table 2 in (Zhang et al., 2023a).
The reason is that the test data used in our experiment is different and more difficult. The test data is a subset of the original training data, which contains much more topics (442 topics) than that in the original test data (48 topics).

432 data of specific domains (WebNLG) or restaurants (E2E NLG Challenge), and then generates indi-433 vidual outputs for the client test data during the evaluation phase. We use metrics including BLEU, 434 NIST, METEOR (MET), TER, ROUGE-L, CIDEr to measure the quality of generated texts.

435 The total communication rounds R are set to 200 for WebNLG and 300 for E2E NLG Challenge, 436 and communication interval is fixed as I = 10 for both datasets. The optimizer setting follows the 437 previous Section 5.1.1. The batch size $\mathcal{B} = 4$, and beam search width bw = 10 are kept for all the 438 baselines. We tune the best step size for each baseline, and the details are summarized in Table 11 439 and Table 12 in Appendix C.3. In addition, the rank initialization for each algorithm and the number 440 of trainable parameters are summarized in Table 13 in Appendix C.3. The test results for GPT-2 441 medium are presented in Tables 3 and 15 in Appendix E, and the results of GPT2-XL are presented 442 in Table 4. We can see that PF2LoRA achieves the best performance in almost all metrics compared to other federated fine-tuning baselines. For example, PF2LoRA on GPT-2 medium achieves 1.25%443 higher BLEU score than HETLoRA on WebNLG and 3.85% higher BLUE score than HETLoRA 444 on E2E NLG Challenge. 445

- 446 5.3 ABLATION STUDIES
- 447

We execute the ablation studies to explore (1) the performance comparison when other baselines 448 have more trainable parameters than ours. (2) the impact of data heterogeneity on PF2LoRA and 449 baselines. (3) the importance of bilevel optimization in our framework. Due to the space limitation, 450 the details about (2) and (3) have been deferred to the Appendix. Refer to Appendix F.1 for the 451 impact of data heterogeneity and Appendix F.2 for the importance of bilevel optimization. 452

453 Baselines with More Trainable Parameters. The lightweight client-specific adapters introduce additional trainable parameters. For fair comparison with other baselines, we consider to increase the 454 number of trainable parameters in other baselines. Specifically, we increase the initial rank r_k (from 455 8 to 12) for baselines HOMLoRA and Per-FedAvg-LoRA in the text classification experiments. 456 Note that HETLoRA has different rank initialization $r_{min} \leq r_k \leq r_{max}$ for different client k, so 457 we count the average trainable parameters of the clients. we can also control the number of trainable 458 parameters by specifying r_{min} and r_{max} . We specify $r_{min} = 5, r_{max} = 12$ in CoLA dataset and 459 $r_{min} = 8, r_{max} = 12$ in other four text classification datasets. The number of trainable parameters 460 of each baseline and the corresponding test score in each dataset are summarized in Table 5. Even if 461 other algorithms have more trainable parameters than our method, PF2LoRA still demonstrates the 462 best performance. PF2LoRA, with negligible additional trainable parameters, significantly improves 463 the performance in personalized federated learning. 464

THEORETICAL JUSTIFICATION 6

In this section, we provide the theoretical justification for the Algorithm 1 in an simplified scenario: we consider the single machine case (M = 1) and assume we have access to the deterministic gradient oracle. In this case the algorithm reduces the following formulation: 470

$$\min_{x} \Phi(x) \coloneqq f(x, y^*(x)), \qquad \text{(UL)}$$

s.t., $y^*(x) \in \arg\min_{x} f(x, y), \qquad \text{(LL)},$ (5)

The update of Algorithm 1 in the single machine case with deterministic gradient reduces to the 474 following update rule: 475

476 477

479

484

465

466 467

468

469

471 472 473

$$y^{t+1} = y^{t} - \alpha \nabla_{y} F(x^{t}, y^{t})$$

$$x^{t+1} = x^{t} - \eta [\nabla_{x} F(x^{t}, y^{t+1}) + \alpha \nabla_{xy} F(x^{t}, y^{t}) \nabla_{y} F(x^{t}, y^{t+1})].$$
(6)

480 We will establish the convergence of the update rule equation 6 under the following assumptions.

481 **Assumption 6.1.** (i) f are bounded below, $\Phi(x_0) - \min_x \Phi(x) \leq \Delta$; (ii) f is μ -strongly convex in 482 terms of y for given x; (iii) f is continuously differentiable and $L_{f,1}$ -smooth jointly in (x, y); (iv) f 483 is twicely differentiable and $\nabla^2 f$ is $L_{f,2}$ -Lipschitz jointly in (x, y).

Remark: These assumptions are standard in the bilevel optimization literature (Kwon et al., 2023; 485 Ji et al., 2021).

Theorem 6.2 (Convergence Guarantees). Suppose Assumption 6.1 holds. Define the smoothness parameter $L_{\Phi} = L_{f,1} + \frac{L_{f,1}^2}{\mu}$, and choose $\alpha = \frac{1}{4L_{f,1}}, \eta =$ $\min\left(\frac{\mu^2}{5L_{f,1}^3\sqrt{(\frac{4L_{f,1}}{\mu} - \frac{\mu}{4L_{f,1}})}}, \frac{1}{8L_{\Phi}}, \sqrt{\frac{1}{16N}}, \sqrt[3]{\frac{1}{81NL_{\Phi}}}\right), \text{ and } N = \frac{25L_{f,1}^4(\frac{4L_{f,1}}{\mu} + 1)}{16\mu^2}.$ Then, we

have $\frac{1}{T}\sum_{t=0}^{T-1} \|\nabla \Phi(x^t)\|^2 \leq O(1/T)$, where T is the total number of iterations.

Remark: Theorem 6.2 provides a convergence guarantee with O(1/T) convergence rate for the squared gradient norm. It means that it requires $O(1/\epsilon^2)$ gradient or Hessian-vector product evaluations for finding an ϵ -stationary point (i.e., finding a x such that $\|\nabla \Phi(x)\| \le \epsilon$). This complexity matches the convergence rate of gradient descent for smooth nonconvex function. In addition, compared with existing double-loop bilevel optimization algorithms such as Ji et al. (2021), our update rule equation 6 is an single-loop bilevel algorithm and hence is easy to implement in practice.

7 CONCLUSION

In this paper, we presented PF2LoRA, a novel personalized federated fine-tuning algorithm for het erogeneous data based on a two-level low-rank adaptation framework, where the first level aims to
 learns a common adapter for all the clients and the second level fosters individual client personal ization. Our approach addresses the limitations of existing methods, such as data-independent rank
 initialization and excessive hyperparameter tuning. Through comprehensive experiments on NLU
 and NLG tasks, PF2LoRA demonstrated significant performance improvements over state-of-the-art
 baselines, with negligible additional memory overhead.

- References
- Sara Babakniya, Ahmed Roushdy Elkordy, Yahya H Ezzeldin, Qingfeng Liu, Kee-Bong Song, Mostafa El-Khamy, and Salman Avestimehr. Slora: Federated parameter efficient fine-tuning of language models. *arXiv preprint arXiv:2308.06522*, 2023.
- Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir
 Ivanov, Chloe Kiddon, Jakub Konečný, Stefano Mazzocchi, Brendan McMahan, et al. Towards
 federated learning at scale: System design. *Proceedings of machine learning and systems*, 1:
 374–388, 2019.
- Theodora S Brisimi, Ruidi Chen, Theofanie Mela, Alex Olshevsky, Ioannis Ch Paschalidis, and Wei
 Shi. Federated learning of predictive models from federated electronic health records. *International journal of medical informatics*, 112:59–67, 2018.
- Yae Jee Cho, Luyang Liu, Zheng Xu, Aldi Fahrezi, and Gauri Joshi. Heterogeneous low rank approximation for federated fine-tuning of on-device foundation models. *arXiv preprint arXiv:2401.06432*, 2024.
 - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning: A meta learning approach. *arXiv preprint arXiv:2002.07948*, 2020.
- 532
 533 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pp. 1126–1135. PMLR, 2017.
- Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez Beltrachini. The webnlg challenge: Generating text from rdf data. In *Proceedings of the 10th international conference on natural language generation*, pp. 124–133, 2017.
- ⁵³⁹ Robin C Geyer, Tassilo Klein, and Moin Nabi. Differentially private federated learning: A client level perspective. *arXiv preprint arXiv:1712.07557*, 2017.

540 541	Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. <i>arXiv preprint arXiv:1802.02246</i> , 2018.
542 543 544	Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. <i>arXiv preprint arXiv:1510.00149</i> , 2015.
545 546 547	Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning. <i>arXiv preprint arXiv:2110.04366</i> , 2021a.
548 549	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert with disentangled attention. <i>arXiv preprint arXiv:2006.03654</i> , 2020.
550 551 552 553	Pengcheng He, Jianfeng Gao, and Weizhu Chen. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. <i>arXiv preprint arXiv:2111.09543</i> , 2021b.
554 555 556	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, An- drea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In <i>International conference on machine learning</i> , pp. 2790–2799. PMLR, 2019.
557 558 559 560	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> , 2021.
561 562	Kaiyi Ji, Junjie Yang, and Yingbin Liang. Bilevel optimization: Convergence analysis and enhanced design. In <i>International conference on machine learning</i> , pp. 4882–4892. PMLR, 2021.
563 564 565 566	Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. <i>Foundations and trends</i> ® <i>in machine learning</i> , 14(1–2):1–210, 2021.
568 569 570	Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In <i>International conference on machine learning</i> , pp. 5132–5143. PMLR, 2020.
571 572 573	Jeongyeol Kwon, Dohyun Kwon, Stephen Wright, and Robert D Nowak. A fully first-order method for stochastic bilevel optimization. In <i>International Conference on Machine Learning</i> , pp. 18083–18113. PMLR, 2023.
575 576 577	Neal Lawton, Anoop Kumar, Govind Thattai, Aram Galstyan, and Greg Ver Steeg. Neural archi- tecture search for parameter-efficient fine-tuning of large pre-trained language models. <i>arXiv</i> <i>preprint arXiv:2305.16597</i> , 2023.
578 579	Jaejun Lee, Raphael Tang, and Jimmy Lin. What would elsa do? freezing layers during transformer fine-tuning. <i>arXiv preprint arXiv:1911.03090</i> , 2019.
580 581 582	Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. <i>arXiv preprint arXiv:2104.08691</i> , 2021.
583 584	Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. <i>arXiv</i> preprint arXiv:2101.00190, 2021.
585 586 587 588	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> , 2019.
589 590 591	Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In <i>Artificial intelligence and statistics</i> , pp. 1273–1282. PMLR, 2017a.
592 593	H Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially private recurrent language models. <i>arXiv preprint arXiv:1710.06963</i> , 2017b.

622

623

624

632

633

634

635

- Jekaterina Novikova, Ondřej Dušek, and Verena Rieser. The e2e dataset: New challenges for end-to-end generation. *arXiv preprint arXiv:1706.09254*, 2017.
 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don't know: Unanswerable questions
 for squad. *arXiv preprint arXiv:1806.03822*, 2018.
- Gregory C Reinsel and Raja P Velu. *Multivariate reduced-rank regression*. Springer, 1998.
- Micah J Sheller, Brandon Edwards, G Anthony Reina, Jason Martin, Sarthak Pati, Aikaterini Kotrot sou, Mikhail Milchenko, Weilin Xu, Daniel Marcus, Rivka R Colen, et al. Federated learning in
 medicine: facilitating multi-institutional collaborations without sharing patient data. *Scientific reports*, 10(1):12598, 2020.
- 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
 4
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.
 Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv* preprint arXiv:1804.07461, 2018.
 - Boxin Wang, Yibo Jacky Zhang, Yuan Cao, Bo Li, H Brendan McMahan, Sewoong Oh, Zheng Xu, and Manzil Zaheer. Can public large language models help private cross-device federated learning? *arXiv preprint arXiv:2305.12132*, 2023.
- Colin Wei, Sang Michael Xie, and Tengyu Ma. Why do pretrained language models help in down stream tasks? an analysis of head and prompt tuning. *Advances in Neural Information Processing Systems*, 34:16158–16170, 2021.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art
 natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pp. 38–45, 2020.
 - Runxin Xu, Fuli Luo, Zhiyuan Zhang, Chuanqi Tan, Baobao Chang, Songfang Huang, and Fei Huang. Raise a child in large language model: Towards effective and generalizable fine-tuning. *arXiv preprint arXiv:2109.05687*, 2021.
- Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept and applications. ACM Transactions on Intelligent Systems and Technology (TIST), 10(2):1–19, 2019.
- Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: Simple parameter-efficient fine-tuning
 for transformer-based masked language-models. *arXiv preprint arXiv:2106.10199*, 2021.
- Jianyi Zhang, Saeed Vahidian, Martin Kuo, Chunyuan Li, Ruiyi Zhang, Tong Yu, Guoyin Wang, and
 Yiran Chen. Towards building the federatedgpt: Federated instruction tuning. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6915–6919. IEEE, 2024.
- Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In *International Conference on Learning Representations*. Openreview, 2023a.

648 649 650 651	Zhuo Zhang, Yuanhang Yang, Yong Dai, Qifan Wang, Yue Yu, Lizhen Qu, and Zenglin Xu. Fed- petuning: When federated learning meets the parameter-efficient tuning methods of pre-trained language models. In <i>Annual Meeting of the Association of Computational Linguistics 2023</i> , pp. 9963–9977. Association for Computational Linguistics (ACL), 2023b.
652	Mangija Zhao, Tao Lin, Eai Mi, Magin Jaggi, and Hingiah Sabijitza. Magking as an affaiant altarna
653 654	tive to finetuning for pretrained language models. <i>arXiv preprint arXiv:2004.12406</i> , 2020.
655	
656	
657	
659	
650	
660	
661	
660	
662	
003	
665	
666	
667	
669	
660	
670	
671	
672	
673	
674	
675	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	
701	

702 DETAILS OF PER-FEDAVG-LORA А 703

704 Per-FedAvg-LoRA. Per-FedAvg-LoRA is built upon a well-known personalized federated learning 705 approach called Per-FedAvg (Fallah et al., 2020), with the trainable model parameters being low 706 rank matrices such as in LoRA. Per-FedAvg is a typical personalized federated learning algorithm, 707 which incorporates Model-Agnostic Meta-Learning (MAML) (Finn et al., 2017) to FedAvg algo-708 rithm (McMahan et al., 2017a) to enable models quickly adapting to heterogeneous data. When it 709 is applied to low rank adaptation, we can get a new variant, namely Per-FedAvg-LoRA. The goal 710 of Per-FedAvg-LoRA is to find a common adapter x which can perform well after it is updated by one-step gradient descent on each client. In particular, Per-FedAvg-LoRA is trying to solve the 711 following formulation using the FedAvg algorithm: 712 $\min_{x} \frac{1}{M} \sum_{k=1}^{M} f_k(x - \alpha \nabla f_k(x)),$

713 714

715 716

717

718 719

where $\alpha > 0$ is the step size. Note that Per-FedAvg-LoRA uses adapter matrices with homogeneous rank across all the clients.

(7)

В **PyTorch-style Pseudocode for PF2LoRA**

720 721 722

723

724

725

726

727

728

In this section, we show the PyTorch-style pseudocode for PF2LoRA. Our two-level low rank adapter framework can be derived by slightly modifying the LoRA module and integrating it into federating learning. When creating low rank adapters, we need to initialize two types of adapters, i.e., common adapters and the client adapters. The initial rank dimension for the common adapter is typically set to r, while for the client adapter, it is set to $\frac{r}{2}$. In addition, we require two different optimizers to update the common and client adapters. The common adapter is updated using AdamW, and the client adapter is updated using SGD. It's important to note that hypergradient calculation is necessary when updating the common adapter. Besides, our framework can be easily plugged into multiple language models, such as RoBERTa, DeBERTa and GPT-2, and others.

729 730 731

С **EXPERIMENT SETUP**

732 733 734

735

C.1 ROBERTA ON TEXT CLASSIFICATION

We use grid search to find the best learning rate for each algorithm in the range of $\{1.0 \times 10^{-4}, 5.0 \times 10$ 736 10^{-4} , 1.0×10^{-3} , 2.0×10^{-3} , 5.0×10^{-3} }. For algorithm Per-FedAvg-LoRA, we search the optimal 737 learning rate for one-step update and the common adapter update, respectively. For PF2LoRA, we 738 also search for the best learning rate for the client-specific adapter update and the common adapter 739 update. The selected learning rates for each algorithm are listed in Table 6, where we use slash 740 to separate two learning rates for Per-FedAvg-LoRA and PF2LoRA, with the former learning rate 741 being for the common adapter. For HETLoRA, we fix the sparsity parameter $\gamma = 0.99$ across all 742 the datasets and set the penalty factor $\lambda = 1.0 \times 10^{-3}$ for CoLA dataset, and $\lambda = 5.0 \times 10^{-3}$ for 743 MNLI, SST-2, QQP, and QNLI. The rank initialization and the number of trainable parameters are 744 summarized in Table 7. The details of the text classification datasets are summarized in Table 8. 745

Table 6: Learning rate setting for RoBERETa model on GLUE benchmark. We use slash to separate 747 two learning rates for Per-FedAvg-LoRA and PF2LoRA. For Per-FedAvg-LoRA, the former one is 748 the learning rate for the common adapter, the latter one is the learning rate for one-step SGD. For 749 PF2LoRA, the former one is the learning rate for the common adapter, the latter one is the learning 750 rate for the client-specific adapter.

(51		1 1				
	Method	CoLA	MNLI	SST-2	QQP	QNLI
752	Centralized LoRA	1.0×10^{-3}				
753	HOMLoRA	1.0×10^{-3}	1.0×10^{-3}	2.0×10^{-3}	1.0×10^{-3}	1.0×10^{-3}
155	Per-FedAvg-LoRA	$2.0 \times 10^{-3}/1.0 \times 10^{-2}$	$1.0 \times 10^{-3}/1.0 \times 10^{-3}$	$1.0 \times 10^{-3}/1.0 \times 10^{-3}$	$1.0 \times 10^{-3}/1.0 \times 10^{-3}$	$2.0 \times 10^{-3}/1.0 \times 10^{-3}$
754	HETLoRA	5.0×10^{-3}	2.0×10^{-3}	2.0×10^{-3}	2.0×10^{-3}	2.0×10^{-3}
755	PF2LoRA	$2.0 \times 10^{-3}/1.0 \times 10^{-4}$	$1.0 \times 10^{-3}/1.0 \times 10^{-3}$			

756 Algorithm 2 PyTorch-style Pseudocode for PF2LoRA model_name: the name of pretrained model 1 758 2 # lr_in, lr_out: the learning rate for client and common adapter 759 3 # T: the total number of communication rounds, I: communication interval 760 4 # r: low rank parameter 761 5 # train_dataloader 762 763 7 import torch.distributed as dist 764 8 dist.init_process_group() 9 target_modules = ["query", "value"] 765 10 pretrained_model = LLM_Model.from_pretrained(model_name) 766 ii model = get_peft_model(pretrained_model, target_modules, r) 767 12 optimizer_outer = AdamW(model.common_adpter.parameters(), lr_in) 768 13 optimizer_inner = SGD (model.client_adpter.parameters(), lr_out) 769 14 15 step = 0770 16 for epoch_idx in range(total_epochs) 771 17 for data_batch in train_dataloader: 772 inner_batch, outer_batch = data_batch 18 773 update_client_adapter(model, inner_batch, optimizer_inner) 19 774 20 update_common_adapter(model, outer_batch, optimizer_outer) 775 21 if step % I == 0: 22 dist.reduce(model.common_adapter.parameters(), dst=0, 776 op=self.dist.ReduceOp.SUM) 777 24 average(model.common_adpter.parameters()) 778 dist.broadcast(model.common_adapter.parameters(), src=0) 25 779 26 step += 1 27 # 28 def get_peft_model(pretrained_model, target_modules, r) 781 for module_name, _ in pretrained_model.named_modules(): 29 782 if module_name in target_modules: 30 783 target_module= pretrained_model.get_submodule(module_name 31 784) 32 create_and_replace(target_module, r) 785 33 786 34 def create_and_replace(target_module, r) 787 35 if isinstance(target_module, Linear): 788 36 target_module.initialize_common_adapter(r) 37 target_module.initialize_client_adapter(r/2) 38 target_module.set_trainable_params() 791 792 793 Table 7: Trainable parameters of RoBERTa-base/large. 794

	# Trainable Parameters
Method	(base/large)
HOMLoRA ($r_k = 8$)	0.30M/0.79M
Per-FedAvg-LoRA ($r_k = 8$)	0.30M/0.79M
HETLORAS ($r_{max} = 12, r_{min} = 8$)	0.35M/0.94M
PF2LoRA ($r_k = 8, \tilde{r}_k = 2$)	0.37M/0.99M

C.2 DEBERTA ON QUESTION ANSWERING

803

We search for the optimal learning rate from the range of $\{1.0 \times 10^{-4}, 5.0 \times 10^{-4}, 1.0 \times 10^{-3}, 2.0 \times 10^{-3}, 5.0 \times 10^{-3}, 1.0 \times 10^{-2}\}$ for each algorithm on SQuAD v1 and v2 dataset. Refer to Table 9 for detailed learning rate settings. The rank initialization and the number of trainable parameters for different algorithms are presented in Table 10. For PF2LoRA, we initialize the rank of client-specific adapter $\tilde{r}_k = \frac{r_k}{2} = 4$, and we set the best value of $r_{min} = 6$, $r_{max} = 14$ for HETLoRA. In addition, HETLoRA uses the sparsity parameter $\gamma = 0.99$ and the penalty factor $\lambda = 5.0 \times 10^{-3}$ on both SQuAD v1 and v2 datasets.

81	0
81	1
81	2

Table 8: The summary of GLUE benchmark.

Corpus	# Train	# Test	# Lable	Metrics
CoLA	8.5k	1k	2	Matthew's correlation
MNLI	393k	20k	3	Accuracy
SST-2	67k	1.8k	2	Accuracy
QQP	364k	391k	2	Accuracy
QNLI	108k	5.7k	2	Accuracy

Table 9: Learning rate choices for question-answering dataset SQuAD v1/v2.

Method	SQuAD v1	SQuAD v2
Centralized LoRA	1.0×10^{-3}	5.0×10^{-4}
HOMLoRA	1.0×10^{-3}	5.0×10^{-4}
Per-FedAvg-LoRA	$2.0 \times 10^{-3}/1.0 \times 10^{-3}$	$1.0 \times 10^{-3}/1.0 \times 10^{-3}$
HETLoRA	5.0×10^{-3}	5.0×10^{-3}
PF2LoRA	$2.0 \times 10^{-3}/1.0 \times 10^{-2}$	$1.0 imes 10^{-3} / 1.0 imes 10^{-2}$

C.3 GPT-2 ON WEBNLG AND E2E NLG CHALLENGES

The optimal learning rates for each algorithm on WebNLG and E2E NLG Challenges are turned from the range $\{1.0 \times 10^{-4}, 5.0 \times 10^{-4}, 1.0 \times 10^{-3}, 2.0 \times 10^{-3}, 3.0 \times 10^{-3}, 4.0 \times 10^{-3}, 5.0 \times 10^{-3}\}, 10^{-3}, 10^{-$ and the learning rate settings are summarized in Table 11. For the rank initialization, we follow LoRA paper (Hu et al., 2021) and choose a small rank $r_k = 4$ for Centralized LoRA, HOMLoRA, and Per-FedAvg-LoRA. We turn the best parameters and set $r_{min} = 6, r_{max} = 12$ for HET-LoRA. PF2LoRA uses the same $r_k = 4$ for the common adapter and $\tilde{r_k} = 2$ for the client-specific adapter. The detailed rank settings and the number of trainable parameters are shown in Table 13. HETLoRA sets the sparsity parameter $\gamma = 0.99$ and the penalty factor $\lambda = 5.0 \times 10^{-4}$ on both WebNLG and E2E NLG Challenge datasets.

D SUPPLEMENTARY EXPERIMENTAL RESULTS FOR TEXT CLASSIFICATION

This section provides experimental results about RoBERTa large model on GLUE benchmark. The comparison results with other baselines are shown in Table 14. We can observe that PF2LoRA achieves higher classification performance. For example, PF2LoRA outperforms HETLoRA by 3.88%, 22.24%, 2.99%, 13.89% and 2.69% on the five datasets, respectively.

E SUPPLEMENTARY EXPERIMENTAL RESULTS FOR E2E NLG CHALLENGE

This section provides experimental results for E2E NLG dataset in Table 15. Compared to other federated baselines, our approach demonstrates the best performance on four metrics (BLEU, NIST, ROUGE-L, CIDEr) of five.

859	Method	Rank initialization	# Trainable parameters
860	Centralized LoRA	$r_k = 8$	0.30M
861	HOMLoRA	$r_k = 8$	0.30M
000	Per-FedAvg-LoRA	$r_k = 8$	0.30M
002	HETLoRA	$r_{min} = 6, r_{max} = 14$	0.30M
863	PF2LoRA	$r_k = 8, \tilde{r}_k = 4$	0.44M

Table 11: Learning rate choices for GPT-2 medium on NLG dataset WebNLG and E2E NLG Challenge.

867	Method	WebNLG		E2E NLG Challenge	
868	Centralized LoRA	1.0×10^{-1}	10^{-3}	1.0×10^{-3}	
869	HOMLoRA	1.0×10^{-1}	10^{-3}	1.0×10^{-3}	
870	Per-FedAvg-LoRA	2.0×10^{-3}	1.0×10^{-4}	$2.0 \times 10^{-3}/2.0 \times 10^{-3}$	
871	HETLoRA	$2.0 \times$	10^{-3}	2.0×10^{-3}	
872	PF2LoRA	$2.0 \times 10^{-3}/$	1.0×10^{-3}	$3.0 \times 10^{-3} / 5.0 \times 10^{-4}$	
873					
874			~~~~		
875	Table 12: Learning ra	ate choices for	GPT2-XL o	n NLG dataset WebNLG.	
876	Method		Web	NLG	
877	HOML	oRA	$1.0 \times$	< 10 ⁻³	
878	Per-Fed	IAvg-LoRA	1.0×10^{-3}	$/1.0 \times 10^{-4}$	
879	HETLO	\mathbf{RA} 1.0 >		$< 10^{-3}$	
880	PF2LoF	RA	1.0×10^{-3}	1.0×10^{-4}	
881					
882	Table 13: Rank i	nitialization ar	nd trainable r	parameters for GPT-2.	
883	Method	Rank initialization		# Trainable parameters	
884	Centralized LoRA	$r_k = 4$		0.39M	
885	HOMLoRA	$r_k = 4$		0.39M	
886	Per-FedAvg-LoRA	$r_k = 4$		0.39M	
887	HETLORA	$r_{min} = 6, r$	$r_{max} = 12$	0.81M	
888	PF2LoRA	$r_k = 4,$	$\tilde{r}_k = 2$	0.59M	

F MORE ABLATION STUDIES.

F.1 THE IMPACT OF HETEROGENEITY LEVELS

Heterogeneity level is regarded as an important factor in federated learning. In this section, we explore the impact of various heterogeneity levels on the performance of algorithms. We run PF2LoRA and other baselines on text classification datasets SST-2 and QNLI with three different heterogeneity levels s = 0.6, 0.9, 1.0. The accuracy results are shown in Table 16. PF2LoRA performs consistently well on different heterogeneity levels, and HETLoRA follows. The performance of HOM-LoRA and Per-FedAvg-LoRA decreases significantly as the heterogeneity level increases. Especially, PF2LoRA outperforms other baselines in a large margin in the case of very high heterogeneity, e.g., 4.35% higher than HETLoRA and 13.87% higher than HOMLoRA on SST-2 dataset.

Next, we further study the impact of relatively lower heterogeneity levels on the algorithms. We run PF2LoRA and other federated baselines on CoLA dataset in the heterogeneity levels of s = 0.2, s = 0.3 and s = 0.4, and the results of "Matthew's correlation" are summarized in Table 17. PF2LoRA outperforms all the baselines consistently in various heterogeneity levels. For example, PF2LoRA surpasses the best baseline HETLoRA by 4.36%, 0.8% and 12.15% in heterogeneity levels of s = 0.2, s = 0.3, s = 0.4 respectively. Therefore, our algorithm PF2LoRA demonstrates the high robustness to heterogeneity levels.

909 910

911

889 890

891 892 893

894

864

F.2 PERFORMANCE WITH/WITHOUT BILEVEL OPTIMIZATION

We conduct an ablation study to verify the effect of bilevel optimization. Instead of applying bilevel
 optimization in equation 3, we update parameters in the common and client-specific adapters simul taneously.

915
916
917
$$\min_{x,y_k} \frac{1}{M} \sum_{k=1}^M f_k(x,y_k),$$
917 (8)

$$f_k(x, y_k) \coloneqq \mathbb{E}_{\xi \sim \mathcal{D}_k} F_k(x, y_k; \xi),$$

Table 14: Roberta-large results on GLUE benchmark. We report "Matthew's correlation" for CoLA and "Accuracy" for MNLI, SST-2, OQP and QNLI. Higher value means "better performance".

· · · · · · · · · · · · · · · · · · ·		 				I · · ·
ihod (Μ	 CoLA	MNLI	SST-2	QQP	QNLI
tralized LoRA	C	 57.32	84.71	93.67	88.43	90.27
MLoRA	Н	51.71	74.51	93.33	79.76	89.63
FedAvg-LoRA	Pe	51.20	75.68	92.64	81.83	79.49
ΓLoRA	Н	54.15	76.38	94.53	82.55	92.31
LoRA	Pl	56.25	93.37	97.36	94.02	94.79
MLoRA FedAvg-LoRA FLoRA	H Pe H Pl	51.52 51.71 51.20 54.15 56.25	74.51 75.68 76.38 93.37	93.33 92.64 94.53 97.36	79.76 81.83 82.55 94.02	90. 89. 79. 92. 94.

Table 15: GPT-2 generation results on E2E dataset

		8			
method	BLEU↑	NIST ↑	MET \uparrow	ROUGE-L↑	CIDEr ↑
Centralized LoRA	0.6833	8.5321	0.4642	0.7046	2.4023
HOMLoRA	0.5585	7.0986	0.4349	0.6095	1.8327
Per-FedAvg-LoRA	0.5683	7.1190	0.4327	0.6109	1.8984
HETLORA	0.5505	7.0088	0.4093	0.5697	1.7167
PF2LoRA	0.5717	7.1621	0.4321	0.6111	1.9088

where \mathcal{D}_k is the data on client k. Specifically, we keep the optimizer settings mentioned in Section 5.1.1, where a SGD optimizer is applied to updating the client-specific adapter and an AdamW optimizer to the common adapter. The difference is that we do not use the hypergradient equation 4 to update the common adapter, instead update it by $x_k^{t+1} = x_k^t - \eta \nabla_x F_k(x_k^t, y_k^t; \xi_k^t)$. We execute our "two-level low rank adaptation" framework without bilevel optimization on text classification of GLUE benchmark. For fair comparison, we keep the same hyperparameter settings as that in Section 5.1.1, including heterogeneity level, learning rates, communication rounds, communication interval and initial rank dimension on the same dataset. The comparison results are shown in Figure 2, where we can see that the framework with bilevel optimization (BO) always performs better than that without BO, especially on harder classification task, such as CoLA dataset.

946 947 948

918

938

939

940

941

942

943

944

945

G STABILITY ANALYSIS

949 950 951

952

953

954

955

956

957

958

959

960

Despite that HETLoRA is a strong baseline which performs usually well on heterogeneous data. However, we empirically observe that the training process of HETLoRA is not as stable as ours and Centralized LoRA in Figure 3, where the training loss and perplexity (ppl) are averaged across all the clients. A possible and reasonable explanation is that HETLoRA adopts dynamical rank pruning and matrices truncation which directly change the intrinsic structure of local adapters, leading to unstable training. On the one hand, pruning removes some columns or rows from the original weights, which can degrade the model performance and require some steps of fine-tuning to recover the performance (Han et al., 2015). On the other hand, each client is required to truncate the common adapter matrices to align the matrices' dimensions at each communication round, which inevitably loses some potentially important information. In contrast, our method circumvents the alignment issue of adapter matrices by assigning a uniform rank r_k to the common adapter and uniform \tilde{r}_k to all the client-specific adapters.

961 962 963

964

Table 16: Results in different heterogeneity levels. We use "Accuracy" to measure the performance here, and higher value means "better performance".

•			SST-2			QNLI	
	Methods	s=0.6	s=0.9	s=1.0	s=0.6	s=0.9	s=1.0
	HOMLoRA	92.66	92.47	83.49	86.62	85.45	67.32
	Per-FedAvg-LoRA	90.80	90.56	85.29	85.32	78.59	50.48
	HETLoRA	93.74	93.67	91.11	89.28	91.86	89.09
_	PF2LoRA	94.12	95.85	95.07	92.87	94.18	93.64

"better performance".

CoLA Methods s=0.2 s=0.3 s=0.4HOMLoRA 52.91 50.75 43.17 51.11 44.44 53.48 Per-FedAvg-LoRA HETLoRA 53.86 53.76 45.03 PF2LoRA 56.20 54.19 50.50

Table 17: Matthew's correlation on CoLA in different heterogeneity levels. Higher value means



Figure 2: Performance comparison with/without bilevel optimization (BO). We report "Matthew's correlation" for CoLA and "Accuracy" for MNLI, SST-2, QQP and QNLI. Higher score means "better performance"



Figure 3: The averaged training loss and perplexity on natural language generation task of WebNLG.

1015 1016 1017

1018

972

973

974

975

976

977

978

979

980 981 982

983

984

985

986

987 988 989

994

995 996

997

998

999 1000

1001

1002

1003

1004

1008

1009

1011

1012

1013 1014

Η **GENERATED RESULT OF NLU**

1019 GENERATED EXAMPLES FOR E2E NLG CHALLENGE H.1 1020

1021 Table 18 and 19 show the generated examples of algorithm HETLoRA and PF2LoRA. The federated fine-tuning experiments are run across 8 clients on E2E NLG Challenges, where we construct the 1023 heterogeneous data by the "name" of restaurants, thus each client has different meta-information from different restaurants. There are 18 restaurants in the test set distributed in 8 clients. We show 1024 the generated examples based given context information on each client, while multiple references 1025 are provided to evaluate the quality of generated contents. We compare the generated contents from HETLoRA and PF2LoRA. In most cases, PF2LoRA can generate more complete and logically
coherent sentences. For example, the generated contents on client 4 and client 7, HETLoRA misses
some important information (highlighted in green). The examples on client 1, 2, 3 and 4, PF2LoRA
produces more grammatically coherent sentences than HELORA.

	Table 18: The generated examples for E2E NLG Challenges
	Client 0
Context	name : blue spice — type : pub — food : english — area : riverside — family friendly : yes — near : rainbow vegetarian ca
References	in riverside, near the rainbow vegetarian café, you can find a family friendly pub called blue spice. if you like english food there is a family - friendly pub called blue spice near the rainbow vegetarian café in riverside. the blue spice is a child - friendly, english pub located in riverside area, near rainbow vegetarian café. blue spice is located near rainbow vegetarian café in the riverside area and is a kid friendly pub that serves analish food.
	there is a pub called blue spice which serves english food, is kid friendly, and is in riverside near rainbow vegetarian café.
	blue spice is a child - friendly pub near rainbow vegetarian café in the riverside area.
	blue spice near rainbow vegetarian café in riverside is a pub serving english meals and child friendly the blue spice is a public it is located near rainbow vegetarian café in the area of riverside this is a family
	friendly pub
	serving english food . an english serving child friendly pub in riverside is blue spice near rainbow vegetarian café
	there is a pub that provides food and is children friendly, near rainbow vegetarian café and the riverside and is
	called blue spice. situated near the rainbow vegetarian café in the riverside area of the city, the blue spice publics ideal if you fancy
	traditional english food whilst out with the kids.
HETLoRA	blue spice is a pub near rainbow vegetarian café in the riverside area . it is family friendly
PF2LoRA	blue spice is a family friendly pub that serves english food. it is located in the riverside area
	near the rainbow vegetarian café .
	Client 1
Context	name : the cricketers — type : coffee shop — customer rating : low — family friendly : no — near : ranch
References	the cricketers is a coffee shop with a low customer rating, located near ranch. it is not family - friendly.
HETLoRA PF2LoRA	city centre coffee shop, the cricketers, is not family - friendly and has a low customer rating. it is located near ranch. north of ranch, there is a coffee shop called the cricketers. it is not family - friendly and has a low customer rating.
	Client 2
Context	name : the mill — type : restaurant — food : english — price : moderate — customer rating : 3 out of 5 — area : riverside — family friendly : yes — near : café rouge
References	the riverside area has restaurant near the café rouge that is both in the moderate price range and kid friendly called the mill it has a 3 out of 5 customer rating and serves english food
	the riverside area near café rouge has a restaurant that is kids - friendly . it has a price range in the mill . i give
	the food a 3 out of 5. the mill is a kids friendly restaurant that has moderate prices and serves english food it has a 3 out of 5 customer
	rating and is located in the riverside area near the café rouge.
HETLoRA	the mill is a moderately priced english restaurant near café rouge in the riverside area. it is kid friendly and has
PF2LoRA	a customer rating of 3 out of 5. the mill is a moderately priced restaurant in the riverside area near café rouge. it serves english food and is kid
	friendly. it has a customer rating of 3 out of 5.
	Client 3
Context	name : the phoenix — type : pub — food : french — price : £ 20 - 25 — customer rating : high — area : riverside — family friendly : no — near : crowne plaza hotel
References	a pub that is not kid friendly is located in the riverside area near crowne plaza hotel. it is named the phoenix
	, has french food and price range of £ 20 - £ 30 and a high customer rating. the phoenix, which is a pub that is not kid friendly, is near crowne plaza hotel and serves french food in the price
	range of \pounds 20 - 25 in the riverside area . it has a high customer rating .
HETLoRA	the phoenix is a pub near the crowne plaza hotel in the riverside area. it has a high customer rating and a price range of $(20, 25)$, it is not hid friendly.
PF2LoRA	the phoenix is a pub in the riverside area near the crowne plaza hotel. it serves french food with a price range of $\frac{1}{2}$
	\pounds 20 - 25 and has a high customer rating . it is not kid friendly .

Table 19: The generated examples for E2E NLG Challenges (continued). Client 4 Context name : the punter — type : resturm — food : italian — price : cheap — customer rating : average — area : rit — family friendly : no — near : rainbow vegetarian café References bello and welcome to the punter , we serve the finest italian food around and have an average customer rating is a restaurant serving italian food for adults can be found on the riverside near rainbow vegetarian café . It is not finally - friendly and has an average customer rating. PEIZLORA Client 5 Client 5 Client 5 Client 5 Client 5 Client 5 Client 4 Client 4 Client 5			
Table 19: The generated examples for E2E NLG Challenges (continued). Client 4 Context name : the punter — type : restaurant — food : italian — price : cheap — customer rating : average — area : ri — food : italian — price : cheap — customer rating italian food around and have an average customer rating italian food arouth the riverside our price range is cheap for such good food at the moment we are not family - friendly. a restaurant wer in a rear trainbow vegetarian café and our area is the riverside our price range is cheap for such good food at the moment we are not family - friendly. a restaurant wer in the riverside area near rainbow vegetarian café . It is not family - friendly and has an average customer rating. PELORA Net putter : byee : pub — food : japanese — price : less than £ 20 — customer rating mills in sol family - friendly : no — near : raja indian cusine Met wallts is fast food with pub on side raja indian cusine Met wallts is fast food with pub on side raja indian cusine in the city centre with price less than £ 20 HETLORA Met wallts is a pub that serves japanese food near the raja indian cusine : ht valuts is japanese with pub on side raja indian cusine in the city centre with price less than £ 20 HETLORA Met wallts is a pub that serves japanese orate is a varage customer rating and is family - fr			
Table 19: The generated examples for E2E NLG Challenges (continued). Client 4 Context name : the punter — type : restaurant — food : italian — price : chcap — customer rating : average — area : rit — family friendly in 0 — near : rainbow vegetarian café and our area is the riverside our price range is the cheap for such good food at the moment we are not family - friendly. References hello and welcome to the punter , we serve the finest italian food around and have an average customer rating this very good food at the moment we are not family - friendly. a restaurant we are near ratinbow vegetarian café and our area is the riverside our price range is the cheap for such good food at the moment we are not family - friendly. a restaurant her punter is located in the riverside area near rainbow vegetarian café . It is not family - friendly and has an average customer rating. PF2LoRA Client 5 Client 5 <td colsp<="" th=""><th></th><th></th></td>	<th></th> <th></th>		
Table 19: The generated examples for E2E NLG Challenges (continued). Client 4 Context name : the punter — type : restaurant — food : italian — price : cheap — customer rating : average — area : ri — family friendly : no — near : rianbow vegetarian café References Hello and welcome to the punter , we serve the fines i italian food around and have an average customer rating the is very good food at the moment we are not family - friendly . a restaurant serving italian food for adults can be found on the riverside near rainbow vegetarian café . it is not family - friendly and has an average customer rating . PF2LoRA the italian restaurant the punter is located in the riverside area near rainbow vegetarian café . it is not family - friendly and has an average customer rating . PF2LoRA the italian restaurant the punter is located in the riverside area . it has an average customer rating and is not family - friendly . Client 5 name : the vaults — type : pub — food : japanese — price : less than £ 20 — customer rating : average the vaults must be that is less than 20 pounds and has an average customer rating and is family - friendly . the vaults is japanese with pub on side raja indian cuisine in the city centre with price less than £ 20 HETLORA the vaults is japanese with pub on side raja indian cuisine in the city centre with average no less than £ 20 Her vaults is japanese with pub on side raja indian cuisine in the city centre with average no less than £ 20 Her vaults is japanese with out on side raja indian cuisin			
Integenerated examples for EZE NLS Challenges (continued). Client 4 Client 4 Context name : the punter — type : restaurant — food : italian — price : cheap — customer rating : average — area : rin — family friendly : no — near : rainbow vegetarian café and our area is the riverside our price range is ' cheap for such good to at the moment we are not family - friendly . a restaurant twe are near rainbow vegetarian café in the riverside near rainbow vegetarian café . it is not family - friendly and has an average customer rating . PF2LoRA Net load to the moment we are not family - friendly . a restaurant serving italian food for adults can be found on the riverside near rainbow vegetarian café . it is not family - friendly and has an average customer rating . Client 5 Client 6 <th></th> <th></th>			
Chert 4 Context array = the purter — type : restaurant — foot : italian — price : cheap — customer rating : average — area : ri — family friendly : no — near : rainbow vegetarian café References hello and welcome to the punter, we serve the finest italian food around and have an average customer rating th is very good for a restaurant we are near rainbow vegetarian café and our area is the riverside our price range is ' cheap for such good food at the moment we are not family - friendly . a restaurant serving italian food for adults can be found on the riverside near rainbow vegetarian café . it is not family - friendly and has an average customer rating . HFTLORA the italian restaurant the punter is located in the riverside area near rainbow vegetarian café . it is not family - friendly and has an average customer rating . HFZLORA the units is a cheap italian restaurant near the rainbow vegetarian café in the riverside area . it has an average customer rating and is not family - friendly : Client 5 Context name : the vaults — type : pub — food : japanese — price : less than £ 20 — customer rating : average the vaults is fast food with pub on side raja indian cuisine for on , the vaults is fast food with pub on side raja indian cuisine the vaults pub that is less than 20 pounds and has an average customer rating and is family - friendly . the vaults is apanese with pub on side raja indian cuisine the vaults is apanese with pub on side raja indian cuisine the vaults is apanese with pub on side raja indian cuisine the city centre with average outsomer rating . HFTLORA the vaults is a pub that serves japanese food . it is located in the city centre with average customer rating . HFTLORA the vaults is a pub that serves japanese food . It is located in the city centre near raja indian cuisine . It is not family - friendly and has a price range of less than £ 20. It is not family - friendly . Client 6 Context		Table 19: The generated examples for E2E NLG Challenges (continued).	
Context name : the punter — type : restautant — tood : name — the : cheap — costoner name : average — area : the — family friendly : no — near : ninbow vegetarian café References hello and welcome to the punter , we serve the finest italian food around and have an average customer rating the server good for a restaurant we are near rainbow vegetarian café and our area is the riverside our price range is cheap for such good food at the moment we are not rainbow vegetarian café in the riverside near rainbow vegetarian café . It is not family - friendly and has an average customer rating . PF2LoRA the titalian restaurant the punter is located in the riverside area near rainbow vegetarian café . It is not family - friendly and has an average customer rating . PF2LoRA the punter is a cheap talian restaurant near the rainbow vegetarian café in the riverside area . It has an average customer rating and is not family - friendly . Client 5 Context name : the vaults — type : pub — food : japanese — price : less than £ 20 — customer rating : average — area : (ity centre — family friendly : no — near : raja indian cuisine the vaults is fast food with pub on side raja indian cuisine in the city centre with price less than £ 20 average the vaults is japanese with pub on side raja indian cuisine in the city centre with average no less than £ 20 HETLoRA the vaults is a pub that serves japanese food . It is located in the city centre near raja indian cuisine . The vaults is a pub that serves japanese food . It is located in the city centre near raja indian cuisine . The vaults is a japanese pub located in the city centre near raja indian cuisine . The vaults is a japanese pub located in the city centre near raja indian cuisine . The vaults is a japanese pub located in the city centre near raja indian cuisine . The vaults is a japanese pub located in the city centre near raja indian cuisine . The vaults is a japanese pub located in the city centre near raja indian cuisine . The vaults is a japanese pub loc	Contout	Client 4	
References hello and velcome to the punter, we serve the finest italian food around and have an average customer rating the is very good for a restaurant we are near rainbow vegetarian café and our area is the riverside our price range is veheap for such good food at the moment we are not family - friendly. a restaurant serving italian food for adults can be found on the riverside near rainbow vegetarian café . it is not family - friendly and has an average customer rating. PF2LoRA the italian restaurant he punter is located in the riverside area near rainbow vegetarian café . it is not family - friendly and has an average customer rating . PF2LoRA the italian restaurant hear in the rainbow vegetarian café in the riverside area . it has an average customer rating and is not family - friendly . Context name : the vaults — type : pub — food : japanese — price : less than £ 20 — customer rating : average the vaults pub that is less than 20 pounds and has an average customer rating and is family - friendly . they serve japanese food near the center of the city and also near the raja indian cuisine . the vaults is japanese with pub on side raja indian cuisine in the city centre with average no less than £ 20 Average the vaults is apanese with pub on side raja indian cuisine . PF2LoRA the vaults is a pub that serves japanese food . it is located in the city centre with average an less than £ 20 HETLoRA the vaults is japanese with pub on side raja indian cuisine . it has an average customer rating. the vaults is a pabe that serves japanese food . it is located in the city centre with average no less than £ 20. HET	Context	- family friendly : no - near : rainbow vegetarian café	
HETLORA the italian restaurant the punter is located in the riverside area near rainbow vegetarian café . it is not family - friendly and has an average customer rating . PF2LORA Client 5 Context name : the vaults — type : pub — food : japanese — price : less than £ 20 — customer rating : average — area : city centre — family friendly : no — near : raja indian cuisine References no , the vaults is fast food with pub on side raja indian cuisine in the city centre with price less than £ 20 average the vaults pub that is less than 20 pounds and has an average customer rating and is family - friendly . they serve japanese food near the center of the city and also near the raja indian cuisine . HETLORA the vaults is japanese with pub on side raja indian cuisine in the city centre with average no less than £ 20 HETLORA the vaults is japanese with pub on side raja indian cuisine in the city centre mat raja indian cuisine . it is not family - friendly and has a price range of less than £ 20. it has an average customer rating and a price range of less than £ 20. it has an average customer rating and a price range of less than £ 20. it is not family - friendly . PF2LORA name : the waterman — type : pub — food : italian — price : high — area : riverside — family friendly : yes — near : raja indian cuisine in the riverside area near raja indian cuisine . it is children friendly . near raja indian cuisine in the riverside area near raja indian cuisine . it is children friendly . Client 6 Context name : the waterman is a high price range italian pub in the riverside area ara raja indian cuisine	References	hello and welcome to the punter , we serve the finest italian food around and have an average customer rating this very good for a restaurant we are near rainbow vegetarian café and our area is the riverside our price range is we cheap for such good food at the moment we are not family - friendly . a restaurant serving italian food for adults can be found on the riverside near rainbow vegetarian café . the punter average ratings , and cheap prices	
PF2LoRA PF2LoRA Client 5 Context name : the valuts — type : pub — food : japanese — price : less than £ 20 — customer rating : average means is an ight price range in the rainbow vegetarian café in the riverside area . it has an average customer rating : average — area : city centre — family friendly : no — near : raja indian cuisine References no , the vaults is fast food with pub on side raja indian cuisine in the city centre with price less than £ 20 average the vaults pub that is less than 20 pounds and has an average customer rating and is family - friendly . they serve japanese food near the center of the city and also near the raja indian cuisine . the vaults is japanese with pub on side raja indian cuisine in the city centre with average no less than £ 20 HETLORA HETLORA He vaults is a pub that serves japanese food . it is located in the city centre with average no less than £ 20 HETLORA HETLORA He vaults is a japanese pub located in the city centre mear raja indian cuisine . it is not family - friendly and has an fice range of less than £ 20. It has an average customer rating and a price range of less than £ 20. it is not family - friendly and has field are care raja indian cuisine . it is not family - friendly and has an field are raja indian cuisine . it is not family - friendly and has an field are range in the riverside area near raja indian cuisine . it is children friendly and has negotive family friendly - friendly . Meterences the waterman is a high price range italian pub in the riverside area near raja indian cuisine . it is children friendly . the waterman is an expensive family friendly establishment located near raja indian cuisine . located near raja indian cuisine in the riverside area , the waterman is an elite, but family friendly ustablished pub . HETLORA the waterman is a night price dialian pub near raja indian cuisine . it is children friendly pub	HETLoRA	the italian restaurant the punter is located in the riverside area near rainbow vegetarian café . it is not	
Client 5 Context name : the vaults — type : pub — food : japanese — price : less than £ 20 — customer rating : average — area : city centre — family friendly : no — near : raja indian cuisine References no , the vaults is fast food with pub on side raja indian cuisine in the city centre with price less than £ 20 average the vaults pub that is less than 20 pounds and has an average customer rating and is family - friendly . they serve japanese food near the center of the city and also near the raja indian cuisine . He vaults is japanese with pub on side raja indian cuisine in the city centre with average no less than £ 20 HETLORA the vaults is a pub that serves japanese food . it is located in the city centre near raja indian cuisine . it is not family - friendly and has a price range of less than £ 20 . it has an average customer rating . PF2LORA the vaults is a pub bloated in the city centre near raja indian cuisine . it is not family - friendly and has a price range of less than £ 20 . it is not family - friendly . Ontext name : the waterman — type : pub — food : italian — price : high — area : riverside — family friendly : yes — near : raja indian cuisine , the waterman pub has high prices and facilities for the family . the pub the waterman is a nexpensive family friendly establishment located near raja indian cuisine . I located near raja indian cuisine in the riverside area near raja indian cuisine . I is children friendly . near raja indian cuisine in the riverside area near raja indian cuisine . I conveniently located on the riverside , close to raja indian cuisine , the waterman is an elite , but family friendly . PF2LORA HETLORA th	PF2LoRA	tamily - friendly and has an average customer rating. the punter is a cheap italian restaurant near the rainbow vegetarian café in the riverside area. it has an average customer rating and is not family - friendly.	
Context name : the vaults — type : pub — food : japanese — price : less than £ 20 — customer rating : average — area : city centre — family friendly : no — near : raja indian cuisine References no , the vaults is fast food with pub on side raja indian cuisine in the city centre with price less than £ 20 average the vaults pub that is less than 20 pounds and has an average customer rating and is family - friendly . they serve japanese food near the center of the city and also near the raja indian cuisine . the vaults is japanese with pub on side raja indian cuisine in the city centre with average no less than £ 20 HETLORA the vaults is a pub that serves japanese food . it is located in the city centre with average no less than £ 20 HETLORA the vaults is a pub that serves japanese of less than £ 20. it has an average customer rating . PF2LORA the vaults is a japanese pub located in the city centre near raja indian cuisine . it is an average customer rating and a price range of less than £ 20. it is not family - friendly . Context name : the waterman — type : pub — food : italian — price : high — area : riverside — family friendly : yes — near : raja indian cuisine , the waterman pub has high prices and facilities for the family . . near raja indian cuisine , the waterman pub has high prices and facilities for the family . . near raja indian cuisine in the riverside area , the waterman is a high priced italian pub in the riverside area raja indian cuisine . . located near raja indian cuisine in the riverside area , the waterman is a neite , but family friendly ustablishment located oner raja indian cuisine . . located near raja indian cuisine in the riverside area near raja indian cuisine . it i		Client 5	
References no., the vaults is fast food with pub on side raja indian cuisine in the city centre with price less than £ 20 average the vaults pub that is less than 20 pounds and has an average customer rating and is family - friendly . they serve japanese food near the center of the city and also near the raja indian cuisine . the vaults is japanese with pub on side raja indian cuisine in the city centre with average no less than £ 20 HETLORA the vaults is a pub that serves japanese food . it is located in the city centre with average no less than £ 20 HETLORA the vaults is a japanese pub located in the city centre near raja indian cuisine . it is not family - friendly and has a price range of less than £ 20 . it has an average customer rating . PF2LoRA the vaults is a japanese pub located in the city centre near raja indian cuisine . it is not family - friendly . Client 6 Context name : the waterman — type : pub — food : italian — price : high — area : riverside — family friendly : yes — near : raja indian cuisine . near raja indian cuisine the viverside area , the waterman is a high price children friendly . near raja indian cuisine in the riverside area , the waterman is a high priced children friendly pub serving food . conveniently located on the riverside area a, the waterman is a high priced children friendly nub serving food . conveniently located on the riverside area near raja indian cuisine . it is children friendly . the waterman is a night priced italian pub near raja indian cuisine . the waterman is a night price pub .	Context	name : the vaults — type : pub — food : japanese — price : less than £ 20 — customer rating : average — area : city centre — family friendly : no — near : raja indian cuisine	
the vaults is japanese with pub on side raja indian cuisine in the city centre with average no less than £ 20 HETLORA the vaults is a pub that serves japanese food . it is located in the city centre near raja indian cuisine . it is not family - friendly and has a price range of less than £ 20. it has an average customer rating . PF2LoRA the vaults is a japanese pub located in the city centre near raja indian cuisine . it has an average customer rating and a price range of less than £ 20. it is not family - friendly . PF2LoRA name : the waterman — type : pub — food : italian — price : high — area : riverside — family friendly : yes — near : raja indian cuisine , the waterman must be has high prices and facilities for the family . the pub the waterman is a nexpensive family friendly establishment located near raja indian cuisine . Net evaluts is a high price in the riverside area , the waterman is an eite , but family priendly pub serving food . conveniently located on the riverside , close to raja indian cuisine in the riverside area . it is children friendly . HETLORA the waterman is a high priced italian pub near raja indian cuisine i. it is children friendly pub serving food . Conveniently located on the riverside , close to raja indian cuisine in the riverside area . it is children friendly . PF2LoRA the waterman is a high priced italian pub near raja indian cuisine i. it is children friendly . PF2LoRA the waterman is a high priced italian pub near raja indian cuisine . it is children friendly . Versite in the waterman is an italian pub located in the riverside area a rear r	References	no, the vaults is fast food with pub on side raja indian cuisine in the city centre with price less than $\pounds 20$ average the vaults pub that is less than 20 pounds and has an average customer rating and is family - friendly. they serve japanese food near the center of the city and also near the raja indian cuisine.	
HETLoRA the vaults is a pub that serves japanese food . it is located in the city centre near raja indian cuisine . it is not family - friendly and has a price range of less than £ 20. it has an average customer rating . PF2LoRA the vaults is a japanese pub located in the city centre near raja indian cuisine . it has an average customer rating . PF2LoRA the vaults is a japanese pub located in the city centre near raja indian cuisine . it has an average customer rating . PF2LoRA name : the waterman — type : pub — food : italian — price : high — area : riverside — family friendly : yes near : raja indian cuisine . References the waterman is a high price range italian pub in the riverside area near raja indian cuisine . it is children friendly . near raja indian cuisine , the waterman pub has high prices and facilities for the family . the pub the waterman is an expensive family friendly establishment located near raja indian cuisine . located near raja indian cuisine in the riverside area , the waterman is a high priced children friendly pub serving food . conveniently located on the riverside , close to raja indian cuisine , the waterman is an elite , but family friendly established pub . HETLORA the waterman is a high priced italian pub near raja indian cuisine in the riverside area . it is children friendly . the waterman is a high price ditalian pub near raja indian cuisine in the riverside area . it is children friendly . HETLORA the waterman is a high priced italian pub near raja indian cuisine . it is children friendly and has a high price range .		the vaults is japanese with pub on side raja indian cuisine in the city centre with average no less than $\pounds 20$	
Client 6 Client 6 Context name : the waterman — type : pub — food : italian — price : high — area : riverside — family friendly : yes — near : raja indian cuisine References the waterman is a high price range italian pub in the riverside area near raja indian cuisine . it is children friendly . near raja indian cuisine , the waterman pub has high prices and facilities for the family . near raja indian cuisine , the waterman pub has high prices and facilities for the family . the pub the waterman is an expensive family friendly establishment located near raja indian cuisine . located near raja indian cuisine in the riverside area , the waterman is a high priced children friendly pub serving food . conveniently located on the riverside , close to raja indian cuisine , the waterman is an elite , but family friendly established pub . HETLORA the waterman is a high priced italian pub near raja indian cuisine in the riverside area . it is children friendly . the waterman is a high price range . Client 7 Context name : name : wildwood — type : pub — food : indian — area : city centre — family friendly : yes — near : raja indian cuisine mate : name : name : wildwood — type : pub — food : indian — area : city centre — family friendly : yes — near : raja indian cuisine wildwood is in the city centre and raja indian cuisine the wildwood is a family friendly indian pub .	HETLoRA PF2LoRA	the vaults is a pub that serves japanese food . it is located in the city centre near raja indian cuisine . it is not family - friendly and has a price range of less than $\pounds 20$. it has an average customer rating . the vaults is a japanese pub located in the city centre near raja indian cuisine . it has an average customer rating	
Client 6 Context name : the waterman — type : pub — food : italian — price : high — area : riverside — family friendly : yes — near : raja indian cuisine References the waterman is a high price range italian pub in the riverside area near raja indian cuisine . it is children friendly . near raja indian cuisine , the waterman pub has high prices and facilities for the family . the pub the waterman is an expensive family friendly establishment located near raja indian cuisine . located near raja indian cuisine in the riverside area , the waterman is a high priced children friendly pub serving food . conveniently located on the riverside , close to raja indian cuisine , the waterman is an elite , but family friendly established pub . HETLORA the waterman is a high priced italian pub near raja indian cuisine in the riverside area . it is children friendly . PF2LoRA PF2LoRA the waterman is a nitalian pub located in the riverside area near raja indian cuisine . it is children friendly and has a high price range . Context name : name : wildwood — type : pub — food : indian — area : city centre — family friendly : yes — near : raja indian cuisine References located near the city centre and raja indian cuisine the wildwood is a family friendly indian pub . wildwood also offers indian food to go along with the family friendly pub located near raja indian cuisine References located near the city centre near raja indian cuisine is kid friendly and serves indian food . a pub near raja indian cuisine in the city centre near raja indian cuisine is kid friendly pub located near raja indian cuisine		and a price range of less than £ 20. it is not family - friendly.	
Context name : the waterman — type : pub — food : italian — price : high — area : riverside — family friendly : yes — near : raja indian cuisine References the waterman is a high price range italian pub in the riverside area near raja indian cuisine . it is children friendly . near raja indian cuisine , the waterman pub has high prices and facilities for the family . the pub the waterman is an expensive family friendly establishment located near raja indian cuisine . located near raja indian cuisine in the riverside area , the waterman is a high priced children friendly pub serving food . conveniently located on the riverside , close to raja indian cuisine , the waterman is an elite , but family friendly established pub . HETLORA the waterman is a high priced italian pub near raja indian cuisine in the riverside area , it is children friendly . PF2LoRA Mewaterman is a high price italian pub near raja indian cuisine in the riverside area , it is children friendly . the waterman is an italian pub located in the riverside area near raja indian cuisine . it is children friendly and has a high price range . Client 7 Context name : name : wildwood — type : pub — food : indian — area : city centre — family friendly : yes — near : raja indian cuisine References located near the city centre and raja indian cuisine the wildwood is a family friendly indian pub . wildwood also offers indian food to go along with the family friendly pub located near raja indian cuisine HETLORA aji indian cuisine pub in the city centre near raja indian cuisine is kid friendly and serves indian food . a pub near raja indian cuisine in the city centre called wildwood serves indian food and i		Client 6	
References the waterman is a high price range italian pub in the riverside area near raja indian cuisine . it is children friendly . near raja indian cuisine , the waterman pub has high prices and facilities for the family . the pub the waterman is an expensive family friendly establishment located near raja indian cuisine . located near raja indian cuisine in the riverside area , the waterman is a high priced children friendly pub serving food . conveniently located on the riverside , close to raja indian cuisine , the waterman is an elite , but family friendly established pub . HETLORA the waterman is a high priced italian pub near raja indian cuisine in the riverside area . it is children friendly . the waterman is an italian pub located in the riverside area near raja indian cuisine . it is children friendly and has a high price range . Client 7 Context name : name : wildwood — type : pub — food : indian — area : city centre — family friendly : yes — near : raja indian cuisine References located near the city centre and raja indian cuisine . it is a pub that serves indian food and is family friendly wildwood also offers indian food to go along with the family friendly pub located near raja indian cuisine HETLORA aji indian cuisine pub in the city centre near raja indian cuisine is kid friendly and serves indian food . a pub near raja indian cuisine in the city centre called wildwood serves indian food and is kid friendly .	Context	name : the waterman — type : pub — food : italian — price : high — area : riverside — family friendly : yes — near : raja indian cuisine	
HETLoRA the waterman is a high priced italian pub near raja indian cuisine in the riverside area . it is children friendly . PF2LoRA the waterman is an italian pub located in the riverside area near raja indian cuisine . it is children friendly and has a high price range . Client 7 Context name : name : wildwood — type : pub — food : indian — area : city centre — family friendly : yes _ near : raja indian cuisine References located near the city centre and raja indian cuisine the wildwood is a family friendly indian pub . wildwood is in the city centre area near raja indian cuisine . it is a pub that serves indian food and is family friendly friendly pub located near raja indian cuisine HETLoRA aji indian cuisine pub in the city centre near raja indian cuisine is kid friendly and serves indian food . a pub near raja indian cuisine in the city centre called wildwood serves indian food and is kid friendly .	References	the waterman is a high price range italian pub in the riverside area near raja indian cuisine . it is children friendly . near raja indian cuisine , the waterman pub has high prices and facilities for the family . the pub the waterman is an expensive family friendly establishment located near raja indian cuisine . located near raja indian cuisine in the riverside area , the waterman is a high priced children friendly pub serving food . conveniently located on the riverside , close to raja indian cuisine , the waterman is an elite , but family friendly established pub .	
Client 7 Context name : name : wildwood — type : pub — food : indian — area : city centre — family friendly : yes — near : raja indian cuisine References located near the city centre and raja indian cuisine the wildwood is a family friendly indian pub . wildwood is in the city centre area near raja indian cuisine . it is a pub that serves indian food and is family frien wildwood also offers indian food to go along with the family friendly pub located near raja indian cuisine HETLORA aji indian cuisine pub in the city centre near raja indian cuisine is kid friendly and serves indian food . a pub near raja indian cuisine in the city centre called wildwood serves indian food and is kid friendly .	HETLoRA PF2LoRA	the waterman is a high priced italian pub near raja indian cuisine in the riverside area . it is children friendly . the waterman is an italian pub located in the riverside area near raja indian cuisine . it is children friendly and has a high price range .	
Context name : name : wildwood — type : pub — food : indian — area : city centre — family friendly : yes — near : raja indian cuisine References located near the city centre and raja indian cuisine the wildwood is a family friendly indian pub . wildwood also offers indian food to go along with the family friendly pub located near raja indian cuisine HETLORA aji indian cuisine pub in the city centre near raja indian cuisine is kid friendly and serves indian food . PF2LORA a pub near raja indian cuisine in the city centre called wildwood serves indian food and is kid friendly .		Client 7	
References located near the city centre and raja indian cuisine the wildwood is a family friendly indian pub . wildwood is in the city centre area near raja indian cuisine . it is a pub that serves indian food and is family friendly wildwood also offers indian food to go along with the family friendly pub located near raja indian cuisine HETLoRA aji indian cuisine pub in the city centre near raja indian cuisine is kid friendly and serves indian food . PF2LoRA a pub near raja indian cuisine in the city centre called wildwood serves indian food and is kid friendly .	Context	name : name : wildwood — type : pub — food : indian — area : city centre — family friendly : yes — near : raja indian cuisine	
HETLoRA aji indian cuisine pub in the city centre near raja indian cuisine is kid friendly and serves indian food . PF2LoRA a pub near raja indian cuisine in the city centre called wildwood serves indian food and is kid friendly .	References	located near the city centre and raja indian cuisine the wildwood is a family friendly indian pub. wildwood is in the city centre area near raja indian cuisine . it is a pub that serves indian food and is family frien wildwood also offers indian food to go along with the family friendly pub located near raja indian cuisine	
	HETLoRA PF2LoRA	aji indian cuisine pub in the city centre near raja indian cuisine is kid friendly and serves indian food . a pub near raja indian cuisine in the city centre called wildwood serves indian food and is kid friendly .	

1134 H.2 GENERATED EXAMPLES FOR WEBNLG

1136 For WebNLG dataset, we construct the heterogeneity data by the topics ['Airport', 'Astronaut', 'Building', 'City', 'ComicsCharacter', 'Food', 'Monument', 'SportsTeam', 'University', 'Written-1137 Work']. These topics are distributed across 8 clients. Thus, the language style varies with the text 1138 topics. We run the personalized federated fine-tuning across 8 clients and report the generated exam-1139 ples for given test context. The comparison results show that PF2LoRA can generate more complete 1140 and high quality sentences than HETLoRA. For example on client 0 and 1, HETLoRA misses key 1141 words "runwayname", "test pilot", which actually are important information. On client 2 and 5, 1142 HETLoRA generates incorrect information, while PF2LoRA produces accurate sentences. 1143

Table 2	0.	The	generated	examples	for	WebNLG
raute 2	ω.	THU	generateu	Crampics	101	WCUINLO.

	6 1
	Client 0 (Airport)
Context	al_asad_airbase : operatingorganisation : united_states_air_force — al_asad_airbase : runwaylength : 3992 . 88 — al_asad_airbase : location : "al anbar province , iraq " — al_asad_airbase : icao_location_identifier : " oraa " — al_asad_airbase : runwayname : " 08 / 26 "
References	al asad air base has a runway name of 08 / 26 which is 3992. 8 in length . it is situated in the al anbar province of iraq, is operated by the united states air force and has the icao location identifier oraa. the united states airport operates the al asad airbase which is located in the al anbar province, iraq. the icao location identifier of al asad airbase is oraa and the length is 3992. 88m and the runway is known as 08 / 26.
HETLoRA	! the united states air force is the operating organisation for all as a airbase which is located in all another province , iraq, the airbase has a runway length of 3992. 88 and the icao location identifier is oraq.
PF2LoRA	the united states air force is the operating organisation for al asad airbase in al anbar province, iraq. the icao location identifier of al asad airbase is oraa and it has a runway length of 3992. 88. the runway name of the airbase is 08 / 26.
	Client 1 (Astronaut)
Context	alan_shepard : status : " deceased " — alan_shepard : almamater : " nwc , m . a . 1957 " — alan_shepard : deathplace : california — alan_shepard : occupation : test_pilot — alan_shepard : birthplace : new_hampshire — alan_shepard : wa selected by nasa : 1959 — alan_shepard : birthdate : " 1923 - 11 - 18 "
References	alan shepard has died in california . he was born on 18 november 1923 in new hampshire and attended school at nwc , graduating in 1957 with an ma . he became a test pilot and was eventually selected by nasa in 1959 . alan shepard was born in new hampshire on november 18th , 1923 . he graduated from nwc in 1957 with an m . a . he was selected by nasa in 1959 and he was a test pilot . he died in california . alan shepard , born on november 18 , 1923 , graduated from nwc in 1957 with an m . a . alan shepard served as a test pilot , and was selected by nasa in 1959 . alan shepard , born in new hampshire on november 18 , 1923 . he graduated from nwc in 1957 with an m . a . alan shepard served as a test pilot , and was selected by nasa in 1959 . alan shepard , born in new hampshire , died in california , .
HETLoRA PF2LoRA	alan shepard was born on november 18th, 1923 in new hampshire . he graduated from nwc in 1957 with an m. a . and was selected by nasa in 1959 . he died in california . alan shepard was born in new hampshire on november 18th , 1923 . he graduated from nwc with a m. a . in 1957 . he was selected by nasa in 1959 and served as a test pilot . alan shepard died in california .
	Client 2 (Building)
Context	adisham_hall : country : sri_lanka — sri_lanka : capital : sri_jayawardenepura_kotte — sri_lanka : currency : sri_lankan_rupee
References	sri jayawardenepura kotte is the capital of sri lanka, which uses the sri lankan rupee as its currency and is the location of adisham hall. sri jayawardenepura kotte is the capital of sri lanka, whose currency is the rupee. adisham hall is located in sri lanka.
HETLoRA	college adisham hall is located in the country of sri lanka, where the capital is sri jayawardenepura kotte and the currency is the sri lankan rupee.
PF2LoRA	alan adisham hall is located in sri lanka, the capital of which is sri jayawardenepura kotte. the currency of sri lanka is the sri lankan rupee.
	Client 3 (File)
Context	big_hero_6_ (film) : starring : ryan_potter — big_hero_6_ (film) : distributor : walt_disney_studios_motion_pictures — baymax : series : big_hero_6_ (film)
References	the movie big hero 6 stars ryan potter which has baymax as one of its characters, was distributed by walt disney studios motion pictures. baymax is a character in the big hero 6 film starring ryan potter and distributed by walt disney studios motion pictures. walt disney studio motion picture distributed the film big hero 6, in which ryan potter starred and baymax is a character
HETLoRA	! baymax is a character in the film big hero 6 which stars ryan potter . the film was distributed by walt disney studios motion pictures .
PF2LoRA	walt disney studios motion pictures is the distributor of big hero 6, a film in which baymax is a character. the film stars ryan potter.

1188 Table 21: The generated examples for WebNLG (continued). 1189 1190 Client 4 (Food) 1191 Context bacon_sandwich : dishvariation : blt - bacon_sandwich : mainingredients : " bread and bacon , with a condiment , 1192 bacon_sandwich : alternativename : " bacon butty , bacon sarnie , rasher sandwich , bacon sanger , piece ' n bacon , bacon 1193 cob , bacon barm , bacon muffin " 1194 References the bacon sandwich, also known as : bacon butty, bacon sarnie, rasher sandwich, bacon sanger, piece n' 1195 bacon , bacon cob , bacon barm , blt and bacon muffin , is from the united kingdom . it includes bread , bacon 1196 and ketchup or brown sauce. the variant blt bacon sandwich includes : bread , bacon , and ketchup or brown sauce . it originates from the united kingdom 1197 and goes by the various names : bacon butty , bacon sarnie , rasher sandwich , bacon sanger , piece n ' bacon , 1198 bacon barm, bacon cob and bacon muffin. 1199 HETLoRA the bacon sandwich is a variation of the blt. the main ingredients are bread and bacon, with a condiment, often ketchup or brown sauce . it is also known as a bacon butty , bacon sarnie , rasher sandwich , bacon sanger , piece ' n bacon , bacon cob, bacon barm, bacon muffin. 1201 PF2LoRA the bacon sandwich is a variation of the blt and is a dish from the united kingdom . the main ingredients of a bacon 1202 sandwich are bread and bacon with ketchup or brown sauce as a condiment . the bacon sandwich can also be known as a bacon butty, bacon sarnie, rasher sandwich, bacon sanger, piece'n bacon, bacon cob, bacon barm, or bacon muffin. 1203 1204 Client 5 (SportsTeam) 1205 Context a.f.c._blackpool:manager:stuart_parker_(footballer) — stuart_parker_(footballer):club: chesterfield_f.c. — stuart_parker_ (footballer): club: stockport_county_f.c. — a.f.c. _blackpool: ground: blackpool — blackpool : leader : labour_party_ (uk) 1207 References a.f.c. blackpool is in blackpool, which council is labour, it has had stuart parker as their manager, 1208 whose football club was stockport county f. c and is attached to chesterfield football club 1209 alan shepard , born on november 18 , 1923 , graduated from nwc in 1957 with an m . a . alan shepard served as a 1210 test pilot, and was selected by nasa in 1959. alan shepard, born in new hampshire, died in california, HETLoRA 1211 ! stuart parker (footballer) is the manager of a.f. c. blackpool who play in blackpool, where the leader is the labour party (uk) and the ground is called blackpool . 1212 PF2LoRA ! a. f. c. blackpool is in blackpool, where the leader is the labour party (uk). the club is managed by 1213 stuart parker (footballer) who played for chesterfield fc and stockport county f. c. 1214 Client 6 (University) 1215 Context romania : ethnicgroup : germans_of_romania — romania : leadertitle : prime_minister_of_romania — alba_iulia : 1216 country : romania -- romania : leadername : klaus_iohannis -- romania : capital : bucharest -- 1_decembrie_1918_university : 1217 city : alba_iulia — romania : anthem : deșteaptă - te , _române ! 1218 References the 1 decembrie 1918 university is in the city alba iulia in romania . klaus iohannis the leader of romania and they also have a prime minister. the germans of romania are the main ethnic group in romania and the capital is bucharest. 1219 the romania anthem is deșteaptă - te , române ! 1220 ! the 1 decembrie 1918 university is located in alba iulia , romania . the country 's leader is prime minister klaus HETLoRA 1221 iohannis and its capital is bucharest . the anthem of the country is desteaptă - te , române ! 1222 PF2LoRA the 1 decembrie 1918 university is located in alba iulia, romania. romania's capital is bucharest and its leader is prime minister klaus iohannis . the national anthem of romania is deșteaptă - te , române ! and its ethnic group is the 1223 germans of romania. 1224 Client 7 (WrittenWork) 1225 Context $administrative_science_quarterly: publisher: cornell_university -- cornell_university: affiliation:$ 1226 association_of_public_and_land - grant_universities -- cornell_university : affiliation : 1227 association_of_american_universities — cornell_university : president : elizabeth_garrett — cornell_university : city : 1228 ithaca . _new_vork References administrative science quarterly was published by cornell university, located in ithaca, new york, and affiliated with the association of public and land grant universities, as well as with the association of american 1230 universities . president of cornell university is elizabeth garrett . 1231 HETLoRA ! the administrative science quarterly is published by cornell university, which is affiliated with the association of 1232 public and land grant universities and the association of american universities . it is located in ithaca , new york . the 1233 president of cornell university is elizabeth garrett . PF2LoRA the administrative science quarterly is published by cornell university, ithaca, new york. the university is affiliated with the association of public and land grant universities and the association of american universities . the president of the university is elizabeth garrett . 1237 1239 1240 1241

¹²⁴² I PROOF OF THEOREM 6.2

1244 I.1 BASIC LEMMAS

1246 The hypergradient estimation is defined as $\nabla \widehat{\Phi}(x; y^{t+1}) = \nabla_x f(x, y^{t+1}) - \alpha \nabla_{xy} f(x, y^t) \nabla_y f(x, y^{t+1}).$

Lemma I.1 (gradient descent for strongly convex and smooth functions). when $\alpha \leq \frac{1}{L_{f,1}}$, for lower level each step we have

$$\|y^{t+1} - y^*(x^t)\| \le (1 - \alpha \mu)^{\frac{1}{2}} \|y^t - y^*(x^t)\|.$$
(9)

1254 *Proof.* Note that

1251 1252 1253

1269

1271

1282

1291

1267 where (i) is because of the μ -strongly convexity, (ii) is because of $L_{g,1}$ -smooth of the function, (iii) is because of $2\alpha(1 - \alpha L_{f,1})(f(x^t, y^t) - \inf_y f(x^t, y)) \ge 0$.

Lemma I.2 (true hypergradient). The hypergradient $\nabla \Phi(x)$ equals to $\nabla_x f(x, y^*(x))$.

1272 *Proof.* By the implicit function theorem (Ghadimi & Wang, 2018), we have

$$\begin{array}{ll} & 1274 \\ 1275 \\ 1276 \\ 1277 \end{array} \quad \nabla \Phi(x) = \nabla_x f(x, y^*(x)) - \nabla_{xy} f(x, y^*(x)) [\nabla_{yy} f(x, y^*(x))]^{-1} \nabla_y f(x, y^*(x)) \stackrel{(i)}{=} \nabla_x f(x, y^*(x)) \\ & \text{where } (i) \text{ holds due to } \nabla_y f(x, y^*(x)) = 0. \end{array}$$

1278 Lemma I.3 (Lipschitz property (Ghadimi & Wang, 2018)). $y^*(x)$ is $\frac{L_{f,1}}{\mu}$ -Lipschitz continuous.

1280 Lemma I.4 (Lipschitz hypergradient). $\Phi(x)$ is L_{Φ} -smooth and $L_{\Phi} = L_{f,1} + \frac{L_{f,1}^2}{\mu}$.

1283 *Proof.* By definition of hypergradient in Lemma I.2 and Assumption 6.1, we have

1290 where (i) comes from Lemma I.3.

1292 I.2 PROOF

1293 1294 Lemma I.5 (Hypergradient bias). Hypergradient estimation $\nabla \widehat{\Phi}(x; y^{t+1})$ satisfy: 1295 $\|\nabla \widehat{\Phi}(x^t; y^{t+1}) - \nabla \Phi(x^t)\| \leq L_{f,1}(\alpha L_{f,1} + 1)(1 - \alpha \mu)^{\frac{1}{2}} \|y^t - y^*(x^t)\|$

$$\begin{aligned} Proof. \text{ Note that} \\ \nabla \widehat{\Phi}(x^{t}; y^{t+1}) &- \nabla \Phi(x^{t}) \\ &= \nabla_{x} f(x^{t}, y^{t+1}) - \alpha \nabla_{xy} f(x^{t}, y^{t}) \nabla_{y} f(x^{t}, y^{t+1}) - \nabla_{x} f(x^{t}, y^{*}(x^{t})) \\ &\stackrel{(i)}{=} \nabla_{x} f(x^{t}, y^{t+1}) - \nabla_{x} f(x^{t}, y^{*}(x^{t})) - \alpha \nabla_{xy} f(x^{t}, y^{t}) (\nabla_{y} f(x^{t}, y^{t+1}) - \nabla_{y} f(x^{t}, y^{*}(x^{t}))) \end{aligned}$$
(13)

where (i) holds due to $\nabla_y f(x^t, y^*(x^t)) = 0$. Then we obtain that $\|\nabla \widehat{\Phi}(x^t; y^{t+1}) - \nabla \Phi(x^t)\|$

$$\begin{aligned} \stackrel{(i)}{\leq} & (L_{f,1} + \alpha L_{f,1}^2) \| y^{t+1} - y^*(x^t) \| \\ \stackrel{(ii)}{\leq} & (L_{f,1} + \alpha L_{f,1}^2) (1 - \alpha \mu)^{\frac{1}{2}} \| y^t - y^*(x^t) \| \\ &= A \| y^t - y^*(x^t) \| \end{aligned}$$
(14)

where $A = (L_{f,1} + \alpha L_{f,1}^2)(1 - \alpha \mu)^{\frac{1}{2}}$, (i) holds because $\nabla_x f$ and $\nabla_y f$ are $L_{f,1}$ Lipschitz with x, y, yand (*ii*) holds due to Lemma I.1.

Lemma I.6 (Hypergradient descent). Define $A = (L_{f,1} + \alpha L_{f,1}^2)(1 - \alpha \mu)^{\frac{1}{2}}$, we have $\Phi(m)$ inf $\Phi(m)$ 1 1 1 T-1 1

$$\frac{1}{T}\sum_{t=0}^{T}\left(\frac{1}{2}-\eta L_{\Phi}\right)\|\nabla\Phi(x^{t})\|^{2} \leq \frac{\Phi(x_{0})-\inf\Phi(x)}{\eta T} + \frac{1}{T}\left(\frac{1}{2}+\eta L_{\Phi}\right)A^{2}\sum_{k=0}^{T}\|y^{t}-y^{*}(x^{t})\|^{2}$$
(16)

Proof. The proof is very similar to the proof of Theorem 1 in Ji et al. (2021). The L_{Φ} -smoothness of $\Phi(x)$ implies that

$$\Phi(x^{t+1}) - \Phi(x^t) \le \langle \nabla \Phi(x^t), x^{t+1} - x^t \rangle + \frac{L_{\Phi}}{2} \|x^{t+1} - x^t\|^2$$
(17)

Define
$$h^t = \nabla \Phi(x^t; y^{t+1}) = \nabla_x f(x^t, y^{t+1}) - \alpha \nabla_{xy} f(x^t, y^t) \nabla_y f(x^t, y^{t+1})$$
. We have
 $\Phi(x^{t+1}) \leq \Phi(x^t) - \eta \langle \nabla \Phi(x^t), h^t \rangle + \frac{L_{\Phi} \eta^2}{2} \|h^t\|$
 $\leq \Phi(x^t) - \eta (\frac{1}{2} - \frac{\eta L_{\Phi}}{2}) \|h^t\|^2 + \frac{\eta^2 L_{\Phi}}{2} \|h^t - \nabla \Phi(x^t)\|^2$

$$\leq \Phi(x^{t}) - (\frac{\eta}{2} - \eta^{2}L_{\Phi}) \|\nabla\Phi(x^{t})\|^{2} + (\frac{\eta}{2} + \eta^{2}L_{\Phi})\|h^{t} - \nabla\Phi(x^{t})\|^{2}$$
(18)

Do telescoping and use Lemma I.5 we get

$$\frac{1}{T}\sum_{t=0}^{T-1} (\frac{1}{2} - \eta L_{\Phi}) \|\nabla\Phi(x^{t})\|^{2} \stackrel{\text{Lemma } I.5}{\leq} \frac{\Phi(x_{0}) - \inf\Phi(x)}{\eta T} + \frac{1}{T} (\frac{1}{2} + \eta L_{\Phi}) A^{2} \sum_{k=0}^{T-1} \|y^{t} - y^{*}(x^{t})\|^{2}$$
(19)

Lemma I.7 (Lower Level Convergence). $\|y^{t+1} - y^*(x^{t+1})\|^2 \le C \|y^t - y^*(x^t)\|^2 + D \|\nabla \Phi(x^t)\|^2$, where $C = 1 - \alpha^2 \mu^2 + 2(1 + \frac{1}{\alpha \mu}) \frac{L_{f,1}^2}{\mu^2} \eta^2 A^2$, $D = 2(1 + \frac{1}{\alpha \mu}) \eta^2 \frac{L_{f,1}^2}{\mu^2}$.

Proof. Note that

 $||y^{t+1} - y^*(x^{t+1})||^2$ $\stackrel{(i)}{\leq} (1+\alpha\mu) \|y^{t+1} - y^*(x^t)\|^2 + (1+\frac{1}{\alpha\mu}) \|y^*(x^{t+1}) - y^*(x^t)\|^2$ $\stackrel{(ii)}{\leq} (1+\alpha\mu)(1-\alpha\mu)\|y^t - y^*(x^t)\|^2 + (1+\frac{1}{\alpha\mu})\frac{L_{f,1}^2}{\mu^2}\|x^{t+1} - x^t\|^2$ $\leq (1+\alpha\mu)(1-\alpha\mu)\|y^t - y^*(x^t)\|^2 + (1+\frac{1}{\alpha\mu})\frac{2L_{f,1}^2}{\mu^2}\eta^2(\|h^t - \nabla\Phi(x^t)\|^2 + \|\nabla\Phi(x^t)\|^2)$ $= C \|y^{t} - y^{*}(x^{t})\|^{2} + D \|\nabla \Phi(x^{t})\|^{2}.$ (20)

where (i) uses the Young's inequality, (ii) is due to Lemma I.1 and the Lipschitzness of the mapping $y^*(x), C = 1 - \alpha^2 \mu^2 + 2(1 + \frac{1}{\alpha \mu})L_y^2 \eta^2 A^2; D = 2(1 + \frac{1}{\alpha \mu})\eta^2 \frac{L_{f,1}^2}{\mu^2}.$

Proof of Theorem 6.2. Substituting Lemma I.7 to Lemma I.6 yields

$$\frac{1}{T} \sum_{t=0}^{T-1} \left[\frac{1}{2} - \eta L_{\Phi} - (\frac{1}{2} + \eta L_{\Phi}) A^2 D \right] \| \nabla \Phi(x^t) \|^2 \\
\leq \frac{\Phi(x^0) - \inf \Phi(x)}{\eta T} + \frac{1}{T} (\frac{1}{2} + \eta L_{\Phi}) A^2 \frac{\| y^0 - y^*(x^0) \|^2}{1 - C},$$
(21)

where $A = (L_{f,1} + \alpha L_{f,1}^2)(1 - \alpha \mu)^{\frac{1}{2}}, C = 1 - \alpha^2 \mu^2 + 2(1 + \frac{1}{\alpha \mu}) \frac{L_{f,1}^2}{\mu^2} \eta^2 A^2; D = 2(1 + \frac{1}{\alpha \mu}) \eta^2 \frac{L_{f,1}^2}{\mu^2}$ We want to carefully choose the parameter α, η s.t. C < 1, $\alpha \leq \frac{1}{L_{f,1}}$ and $\frac{1}{2} - \eta L_{\Phi} - (\frac{1}{2} + \eta L_{\Phi})A^2D > 0$. For example, we can choose $\alpha = \frac{1}{4L_{f,1}}, \eta = 0$ $\min\left(\frac{\mu^2}{5L_{f,1}^3\sqrt{(\frac{4L_{f,1}}{\mu}-\frac{\mu}{4L_{f,1}})}}, \frac{1}{8L_{\Phi}}, \sqrt{\frac{1}{16N}}, \sqrt[3]{\frac{1}{81NL_{\Phi}}}\right), \text{ and } N = \frac{25L_{f,1}^4(\frac{4L_{f,1}}{\mu}+1)}{16\mu^2}.$

AN EXAMPLE ON MULTIVARIATE LINEAR REGRESSION J

To clarify this mechanism of "this automatic rank adaptation of PF2LoRA", we first construct a multivariate linear regression example and provide a theoretical analysis to demonstrate why our method can accurately learn the ground truth rank, whereas HETLoRA fails. Then we conduct a synthetical experiment to compare two algorithms in federated learning with two clients. The experimental results confirm that our algorithm is able to learn the ground truth ranks for two clients and converge to the optimal solution. In contrast, HETLoRA underestimates the initial rank of some clients due to random rank initialization strategy, resulting in underfitting and suboptimal performance in such clients.

J.1 THEORETICAL ANALYSIS

If our algorithm can find a better low rank approximation than HETLoRA, then our method surely performs better than HETLoRA. So theoretically, we want to find the exact analytic solution of the best low rank approximation. Recall multivariate linear regression problem, the goal is to minimize the reconstruction error:

$$\min_{V \in \mathbb{R}^{m \times n}} \|Y - XW\|_F^2$$

where (X, Y) is the data and label. We know the solution which can minimize the reconstruction error is, $W = (X^T X)^{-1} X^T Y$

However, rank(W) is possibly very large, leading to computationally inefficient. So we want to find the optimal low-rank matrix approximation of W (i.e. minimize the reconstruction error with small rank of W), then we add a rank restriction on W,

$$Y = XW + \epsilon$$
, s.t., $rank(W) \le r$.

In statistics, this is a Reduced Rank Regression (RRR) problem, which has been well-explored,

1400
$$\min_{W \in \mathbb{R}^{m \times n}} \|Y - XW\|_F^2, \quad \text{s.t.}, \quad rank(W) \le r,$$

which is equivalent to

$$\min_{W \in \mathbb{R}^{m \times n}} tr[(Y - XW)(Y - XW)^T], \quad rank(W) \le r$$

where tr(.) is the matrix trace. Given the upper bound of rank(W) = r we directly do rank factorization on W, i.e., LoRA:

Given the upper bound of
$$rank(W) = r$$
, we directly do rank factorization on W, i.e., LoRA
1406

$$\min_{A \in \mathbb{R}^{m \times r}, B \in \mathbb{R}^{r \times n}} tr[(Y - XAB)(Y - XAB)^T],$$

1407 1408

1412

1416 1417

1421

1426 1427

1430

1432 1433

1436 1437 1438

1404

Specifically in HETLora setting, given the rank initialization of the k-client: r_k^{init} , the objective function is: $tr[(V - V - AP)(V - V - AP)^T]$

$$\min_{A \in \mathbb{R}^{m \times r_k^{init}}, B \in \mathbb{R}^{r_k^{init} \times n}} tr[(Y_k - X_k A B)(Y_k - X_k A B)^T].$$

1413 1414 In our setting, we initialize the rank of the common adapter to r, and the local adapter to \tilde{r} , the objective function is,

$$\min_{A \in \mathbb{R}^{m \times r}, B \in \mathbb{R}^{r \times n}, C_k \in \mathbb{R}^{m \times \bar{r}}, D_k \in \mathbb{R}^{\bar{r} \times n}} tr[(Y_k - X_k(AB + C_kD_k))(Y_k - X_k(AB + C_kD_k)^T].$$

1418 In the synthetic experiment, we make global AB to be in the orthogonal row and vector space of $C_k D_k$, then we directly get

$$r(W_k) = r(AB + C_kD_k) = r(AB) + r(C_kD_k) = r + \tilde{r}$$

then our problem is equivalent to reduced-rank regression problem.

Lemma J.1. (*Reinsel & Velu, 1998*) Theorem 2.2[*RRR solution*] Suppose the (m + n)-dimensional random vector (Y_k, X_k) has mean vector 0 and covariance matrix with:

$$\Sigma_{yx} = \Sigma_{xy} = Cov(Y_k, X_k)$$
 and $\Sigma_{xx} = Cov(X_k)$ nonsingular.

1428 Then, for any positive-definite matrix Σ , an $m \times r$ matrix A and $r \times n$ matrix B, for $r \le \min(m, n)$, 1429 which minimize

$$tr\left\{\mathbb{E}\left[\Sigma^{1/2}(Y_k - X_kAB)(Y_k - X_kAB)^{\top}\Sigma^{1/2}\right]\right\}$$

1431 are given by:

$$A^{(r)} = \Sigma^{-1/2}[V_1, \dots, V_r] = \Sigma^{-1/2}V, \quad B^{(r)} = V^{\top}\Sigma^{1/2}\Sigma_{yx}\Sigma_{xx}^{-1}$$

1434 where $V = [V_1, ..., V_r]$ and V_j is the (normalized) eigenvector that corresponds to the *j*-th largest 1435 eigenvalue λ_i^2 of the matrix:

$$\Sigma^{1/2} \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \Sigma^{1/2}, \quad j = 1, 2, \dots, r.$$

1439 From solution formula we directly get minimum truncated error

$$\min_{A,B:rank(AB) \le r} \|W - AB\|_F^2 = \sqrt{\sum_{i=r+1}^n \lambda_i} \quad \forall W, rank(W) \ge r$$

1442 1443

1448 1449

1454

1440 1441

1444 J.1.1 LOW-RANK APPROXIMATION 1445

1446 Specifically in HETLoRA setting, given the rank initialization of the k-client: r_k^{init} , the objective 1447 function is:

$$\min_{A \in \mathbb{R}^{m \times r_k^{init}}, B \in \mathbb{R}^{r_k^{init} \times n}} tr[(Y_k - X_k A B)(Y_k - X_k A B)^T].$$

1450 In our setting, we initialize the rank of the common adapter to r, and the local adapter to \tilde{r} , the objective function is,

1452
1453
$$\min_{A \in \mathbb{R}^{m \times r}, B \in \mathbb{R}^{r \times n}, C_k \in \mathbb{R}^{m \times \tilde{r}}, D_k \in \mathbb{R}^{\tilde{r} \times n}} tr[(Y_k - X_k(AB + C_kD_k))(Y_k - X_k(AB + C_kD_k)^T].$$
(22)

1455 note $C_k D_k$, is a local adapter. we mark

1456
$$W_k = P_k Q_k = AB + C_k D_k$$
1457

note that $rank(P_kQ_k) \in [r - \tilde{r}, r + \tilde{r}]$. Generally we cannot say the problem (22) and

1458 $\min_{P_k \in \mathbb{R}^{m \times r + \bar{r}}} \prod_{Q_1 \in \mathbb{R}^{r + \bar{r} \times n}} tr[(Y_k - X_k P_k Q_k)(Y_k - X_k P_k Q_k)]^T,$ 1459 1460 are equivalent since the former one is subset of the latter problem. However, under some certain 1461 dataset setting, the two problems are equivalence. We defer the equivalence proof to Lemma J.1.2. 1462 1463 Suppose we have two clients, the optimal solution in HETLoRA is 1464 Client 1 $A_1^{r_1^{init}} = \Sigma^{-1/2}[V_1, \dots, V_{r_{init}}] = \Sigma^{-1/2}V, \quad B_1^{r_1^{init}} = V^{\top}\Sigma^{1/2}\Sigma_{yx}\Sigma_{xx}^{-1}$ 1465 1466 Client 2 $A_2^{r_2^{init}} = \Sigma^{-1/2}[V_1, \dots, V_{r_{init}}] = \Sigma^{-1/2}V, \quad B_2^{r_2^{init}} = V^{\top}\Sigma^{1/2}\Sigma_{yx}\Sigma_{xx}^{-1}$ 1467 1468 In our setting, the optimal solution is 1469 Client 1 $P_1^{r+\tilde{r}_1} = \Sigma^{-1/2}[V_1, \dots, V_{r+\tilde{r}_1}] = \Sigma^{-1/2}V, \quad Q_1^{r+\tilde{r}_1} = V^{\top}\Sigma^{1/2}\Sigma_{yx}\Sigma_{rx}^{-1}$ 1470 1471 Client 2 $P_2^{r+\tilde{r_2}} = \Sigma^{-1/2}[V_1, \dots, V_{r+\tilde{r_2}}] = \Sigma^{-1/2}V, \quad Q_1^{r+\tilde{r_1}} = V^{\top}\Sigma^{1/2}\Sigma_{ux}\Sigma_{rx}^{-1}$ 1472 1473 Suppose for Client 1 data, $\Sigma^{1/2} \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \Sigma^{1/2}$ has eigenvector $\lambda_1 = \lambda_2 = \lambda_3 = 1$; $\lambda_4 = \cdots =$ 1474 $\lambda_n = 0$, obviously the low rank approximation is $r_1^* = 3$. For Client 2 data, $\Sigma^{1/2} \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \Sigma^{1/2}$ 1475 has eigenvector $\lambda_1 = \cdots = \lambda_4 = 1$; $\lambda_5 = \cdots = \lambda_n = 0$, the low rank approximation is $r_2^* = 4$. 1476 In our synthetic experiments J.2, HETLoRA underestimates the rank for client 1, i.e., $r_1^{init} = 2 < 1$ 1477 $r_1^* = 3$ due to the random rank initialization, and the learned rank $r_1 = 1$ by self-pruning; Client 2 1478 initializes a reasonable $r_2^{init} = 10$, and the learned rank $r_2 = 5 = r_2^*$. Thus client 1 fails to learn the 1479 optimal low rank approximation because 1480 1481 $\min_{A,B:rank(AB) \le r_1^{init}} \|W - AB\|_F^2 = \sqrt{\sum_{i=r+1}^n \lambda_i} = 1.$ 1482 1483 1484 Our PF2LoRA initializes r = 4 for the common adapter (AB), and $\tilde{r} = 2 (C_k D_k)$ for the local 1485 adapter, which means $r - \tilde{r} = 2 < rank(AB + C_kD_k) < r + \tilde{r} = 6$, and learned rank for client 1486 1 is $r_1 = 3$. The learned rank for client 2 is $r_2 = 4$. Both succeeded to learn the optimal low rank 1487 approximation. 1488 1489

$$\min_{B,C_k,D_k:r-\tilde{r}\leq rank(W_k)\leq r+\tilde{r}} \|W-W_k\|_F^2 = \sqrt{\sum_{i=r+\tilde{r}}^n \lambda_i} = 0.$$

1493 J.1.2 PROBLEM EQUIVALENCE

1495 Next we prove two problems to be equivalent:

A

$$\underset{A \in \mathbb{R}^{m \times r}}{\min} \min_{B \in \mathbb{R}^{r \times n}} \min_{C_k \in \mathbb{R}^{m \times \tilde{r}}} D_k \in \mathbb{R}^{\tilde{r} \times n} tr[(Y_k - X_k(AB + C_kD_k))(Y_k - X_k(AB + C_kD_k)^T]$$

1498 and

1490

1491 1492

1494

1496 1497

1499 1500

1502

1505

1507

$$\min_{W_k \in \mathbb{R}^{m \times n}} tr[(Y_k - X_k W_k)(Y_k - X_k W_k)^T], \quad r - \tilde{r} \le rank(W_k) \le r + \tilde{r}$$

Lemma J.2. The rank of the sum of AB and CD satisfies:

$$r(AB + CD) = r(AB) + r(CD)$$

1503 1504 *if and only if*

$$dim(\mathcal{C}_1 \cap \mathcal{C}_2) = dim(\mathcal{R}_1 \cap \mathcal{R}_2) = 0.$$

where C_1 and C_2 be the column spaces of AB and CD, and \mathcal{R}_1 , \mathcal{R}_2 are their row spaces.

Proof. To simplify the notation in proof, we mark $c = \dim(\mathcal{C}_1 \cap \mathcal{C}_2), d = \dim(\mathcal{R}_1 \cap \mathcal{R}_2); E = AB$, **F** = CD. First, the condition c = d = 0 is necessary, as two strings of inequalities show: **1510**

1510

$$r(E+F) \le r[(E,F)] = r(E) + r(F) - c \le r(E) + r(F),$$
1511

$$r(E+F) \le r[(E;F)] = r(E) + r(F) - d \le r(E) + r(F).$$

To show c = d = 0 is sufficient, we use full-rank decompositions of E and F:

$$E = C_1 R_1, \quad r(A) = r(C_1) = r(R_1) =$$

where E is $m \times n$, C_1 is $m \times a$, and R_1 is $a \times n$.

$$F = C_2 R_2, \quad r(F) = r(C_2) = r(R_2) = b,$$

a,

where F is $m \times n$, C_2 is $m \times b$, and R_2 is $b \times n$.

Such representations exist since R_1 can be any matrix whose rows form a basis of the row space of A. Then $A = C_1 R_1$ for some C_1 , and:

$$r(E) = r(C_1) = \min\left(\operatorname{rank}(C_1), \operatorname{rank}(R_1)\right) \le a = r(E).$$

We now write:

$$E + F = C_1 R_1 + C_2 R_2 = (C_1, C_2) \begin{pmatrix} R_1 \\ R_2 \end{pmatrix} = CR,$$

Then c = 0 implies that all the a + b columns of C are linearly independent, and so C has a left inverse L such that LC = I. Thus, when c = 0,

1531
$$r(E+F) = r(CR) \ge r(LCR) = r(R) = r(E) + r(F) - d$$
.

If in addition d = 0, the entire string collapses, and:

$$r(E+F) = r(E) + r(F).$$

In the following synthetic experiment setting we make global AB in orthogonal row and vector space of C_1D_1 , C_2D_2 , according to above lemma we directly get

$$r(W_1) = r(AB + C_1D_1) = r(AB) + r(C_1D_1) = r + \tilde{r}$$

and

$$r(W_2) = r(AB + C_2D_2) = r(AB) + r(C_2D_2) = r + \hat{r}$$

So under our synthetic experiment setting, our problem is equivalent to reduced-rank regression problem, which provides a theoretical guarantee.

J.2 SYNTHETIC EXPERIMENT

We conduct a synthetic experiment of multivariate linear regression in federated learning to show why HETLoRA fails to learn the ground truth rank, but PF2LoRA does. First, we give the objective function,

$$\min_{A,B} \sum_{k=1}^{2} \|X_k W - Y_k\|_F^2$$

where (X_k, Y_k) is the client-k's data, W is a low-rank matrix and can be decomposed into low-rank matrices. The details of synthetic experiments are described as follows,

> 1. Ground truth of trainable parameters. Given two clients, we assume that we have two optimal solution with low-rank structure.

1561
1562
$$W_1^* = A_1^* B_1^*, \ s.t., W_1^* \in \mathbb{R}^{10 \times 10}, A_1^* \in \mathbb{R}^{10 \times 3}, B_1^* \in \mathbb{R}^{3 \times 10}, B_1^* \in \mathbb{R}^{3 \times 10}$$

with the rank $rank(W_1^*) = 3$, $rank(W_2^*) = 4$. We initialize the random matrices $A_1^*, B_1^*, A_2^*, B_2^* \sim \mathcal{N}(0, 1).$

2. Training and testing data. We construct the synthetic data (X, Y) for two clients respectively by randomly generating 1000 samples, i.e., $X_1 \in \mathcal{R}^{1000 \times 10}$, $s.t., X_1 \sim \mathcal{N}(0, 1)$, $X_2 \in \mathcal{R}^{1000 \times 10}$, $s.t., X_2 \sim \mathcal{N}(0, 1)$, and their targets,

$$Y_1 = X_1 W_1^* + \epsilon_1, \ \epsilon_1 \sim \mathcal{N}(0, 0.1),$$

$$Y_2 = X_2 W_2^* + \epsilon_2, \ \epsilon_2 \sim \mathcal{N}(0, 0.2).$$

The first 700 samples serve as the training set \mathcal{D}_k^{tr} , k = 1, 2 and the remaining serves as the testing set \mathcal{D}_k^{te} , k = 1, 2.

3. Training process.

1566

1567 1568

1569 1570 1571

1572

1573

1574

1575

1579 1580

1581 1582

1584

1585

1586

1587

1590

1591

1592

1596

1598

1602

1603

1604

1605

1607

1608

1609

1610

1611

(a) HETLoRA: Following its rank initialization strategy $r_{min} \leq rank_1 \leq rank_2 \dots \leq rank_k \dots \leq r_{max}$, we assume that $r_{min} = 1, r_{max} = 12$ and initialize $\hat{W}_k = \hat{A}_k \hat{B}_k$ by,

$$\hat{A}_1 \in \mathbb{R}^{10 \times 2}, \hat{B}_1 \in \mathbf{0}^{2 \times 10}, \ s.t. \ \hat{A}_1 \sim \mathcal{N}(0, 1), \hat{A}_2 \in \mathbb{R}^{10 \times 10}, \ \hat{B}_2 \in \mathbf{0}^{10 \times 10}, \ s.t. \ \hat{A}_2 \sim \mathcal{N}(0, 1)$$

so we have $rank(\hat{A}_1) = 2$ and $rank(\hat{A}_2) = 10$. We can easily get that the total number of trainable parameters for two clients is 240. Other hyperparameters are set as follows. We search the regularization factor γ in the range [0.05, 0.5] with the search grid 0.05 and set it to the optimal value 0.1. The pruning parameter $\gamma = 0.3$, which is responsible for imposing the regularization to the last 30% columns to sparse them. We tune the learning rate within the range {0.001, 0.002, 0.003, 0.004, 0.005} and set it to the optimal value 0.002. The total training steps are 2000, and the communication is performed every 10 steps, which means we train the parameters for 10 steps locally, and then execute the parameter aggregation and distribution.

(b) PF2LoRA: For a fair comparison, we initialize the trainable parameters $W_k = \hat{A}_k \hat{B}_k + \hat{C}_k \hat{D}_k$, and make sure the total number of trainable parameters to be the same as that in HETLoRA. For client k = 1, 2, we have r = 4, $\tilde{r} = 2$ and,

$$\begin{split} \hat{A}_k &\in \mathbb{R}^{10\times 4}, \hat{B}_k \in \mathbb{R}^{4\times 10}, \hat{C}_k \in \mathbb{R}^{10\times 2}, \hat{D}_k \in \mathbb{R}^{2\times 10},\\ s.t. \ \hat{A}_k &\sim \mathcal{N}(0,1), \hat{C}_k \sim \mathcal{N}(0,1), \hat{C}_k \sim \mathcal{N}(0,1), \hat{D}_k \sim \mathcal{N}(0,1). \end{split}$$

and $A_k B_k$ is orthogonal to the matrix $C_k D_k$, such that their column space or row space are independent mutually. The total number of training steps are fixed as 2000, and the communication interval is 10. We search the best upper-level and lower-level learning rates within the range [0.001, 0.01] with the search grid of 0.001, and set the best upper-level learning rate to 0.005 and the lower-level learning rate to 0.002. In each communication round, we aggregate the common adapter parameters A_k, B_k and then distribute them, and the local adapter parameters C_k, D_k are not involved in communication.

4. Evaluation. We evaluate the trained model every communication round on the testing data \mathcal{D}_k^{te} , and measure the distance between \hat{W}_k and W_k^* by $\|\hat{W}_k - W_k^*\|_F^2$. In addition, we record the rank evolution of two clients as the training steps. For PF2LoRA, We compute the singular value $\{\lambda_i | i = 1, ..., 10\}$ of \hat{W}_k by singular value decomposition (SVD) and determine its rank: $\min_{1 \le j \le 10} \sum_{i=1}^j \lambda_i \ge 0.9 \sum_{i=1}^{10} \lambda_i$, where λ_i keeps descending order. All the comparison results, including training, testing loss, frobenius norm distance and rank evolution are shown in Figure 4.

1612 As you can see in the last column of Figure 2 (b), PF2LoRA can learn the ground truth rank of 3 1613 in client 1 and 4 in client 2, which verifies that our algorithm can automatically learn the rank in 1614 the range $|r - \tilde{r}, r + \tilde{r}|$. The training and testing loss decrease rapidly and converge to small values 1615 close to 0. The distance to the ground truth parameters also decreases consistently to a small value. 1616 Instead, HETLoRA fails to learn the ground truth rank in the first client. Specifically, the first client underestimates the initial rank due to the random initialization strategy, such that it cannot cover 1617 the ground truth rank. Rank pruning further reduces the first client's rank to $r_{min} = 1$, leading 1618 to increasing the training, testing loss and Frobenius norm distance. Since the second client has 1619 reasonable rank initialization, it is able to learn the ground truth rank by rank pruning.



1671 L COMPUTATION AND COMMUNICATION COST

1672

1673 We evaluated the total computational costs (FLOPs on 8 NVIDIA RTX A6000 GPUs) and communication costs in a single communication round for each algorithm on GLUE benchmark. The results

1674 are summarized in Table 22. From our understanding, communication costs are the total number of 1675 parameters that participate in the aggregation and distribution of parameters in federated learning. 1676 The computational cost (FLOPs) per round are determined by the number of model parameters and 1677 the forward/backward propagation operations. As PF2LoRA requires to compute the hessian-vector 1678 product for hypergradient estimation, it incurs a higher computational cost. But the communication cost of PF2LoRA remains consistent with that of HOMLoRA and Centralized LoR, as the commu-1679 nication parameters in PF2LoRA are only global adapters that have the same rank $r_k = 8$ with that 1680 in HOMLoRA and Centralized LoRA. Instead, HETLoRA has a higher parameter rank requirement 1681 for a high performance, resulting in increased communication costs. 1682

Table 22: Computational/Communication costs per communication round.

1685	Method	TFLOPs/round	Communication parameters/round
1686	$\frac{1}{1}$ Centralized LoRA ($r_{L} = 8$)	258.40	0 30M
1687	HOMLoRA $(r_k = 8)$	258.40	0.30M
1688	Per-FedAvg-LoRA $(r_k = 8)$	908.00	0.30M
1689	HETLORA $(r_{max} = 12, r_{min} = 8)$	272.60	0.35M
1690	PF2LoRA ($r_k = 8, \tilde{r} = 2$)	1202.40	0.30M
691			
692			
693			
694			
695			
606			
697			
609			
600			
700			
700			
701			
702			
703			
704			
705			
700			
707			
700			
709			
710			
710			
712			
713			
714			
715			
710			
710			
710			
719			
720			
721			
722			
723			
724			
725			
/26			
(27			