FEDERATED INSTRUCTION TUNING OF LLMS WITH DOMAIN COVERAGE AUGMENTATION

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026 027 028

029

031

032

033

034

037

Paper under double-blind review

ABSTRACT

Federated Domain-specific Instruction Tuning (FedDIT) utilizes limited crossclient private data together with various strategies of instruction augmentation, ultimately boosting model performance within specific domains. To date, the factors affecting FedDIT remain unclear, and existing instruction augmentation methods primarily focus on the centralized setting without considering distributed environments. Our experiments reveal that the cross-client domain coverage, rather than data heterogeneity, drives model performance in FedDIT. In response, we propose FedDCA, which optimizes domain coverage through greedy client center selection and retrieval-based augmentation. For client-side computational efficiency and system scalability, FedDCA*, the variant of FedDCA, utilizes heterogeneous encoders with server-side feature alignment. Extensive experiments across four distinct domains (code, medical, financial, and mathematical) substantiate the effectiveness of both methods. Additionally, we investigate privacy preservation against memory extraction attacks utilizing various amounts of public data. Results show that there is no significant correlation between the volume of public data and the privacy-preserving capability. However, as the fine-tuning rounds increase, the risk of privacy leakage reduces or converges.

- 1 INTRODUCTION
- Table 1: Performance (%) of different augmentation settings in each domain after FedDIT, under the default setting elaborated in Appendix A.2. Zero-shot directly inferences without fine-tuning, while the Base Data utilizes only the client's local data for FedDIT. Additionally, we compare FedDCA with other two augmentation strategies performed on the server: random sampling (Random for short) and direct retrieval (Direct Re. for short and is described in Appendix A.4), respectively.

Domain	Task/Metric	Zero-shot	Base Data	Random	Direct Re.	FedDCA (ours)
Code	HumanEval/Pass@1	29.88	39.03	32.93	34.14	36.58
Med.	MMLU-Med/Acc.	70.60	<u>68.40</u>	71.30	72.20	74.50
E	FPB/Acc.	55.94	58.25	64.19	66.31	67.24 25.27
F1 n .	TFNS/Acc.	18.54 59.21	$\frac{14.18}{66.62}$	<u>13.09</u> 65.53	67.62	55.27 73.32
Math.	GSM8K/Exact Match	23.27	47.46	47.38	50.87	52.46

Recently, federated instruction tuning (FedIT) has gained attention as a novel approach that leverages the principles of federated learning (FL) to facilitate collaborative training of large language models (LLM) in distributed environments while maintaining the confidentiality of private data (McMahan et al., 2017; Ye et al., 2024b; Zhang et al., 2023c). This methodology allows for the exchange of model parameters among distributed data holders, thereby achieving a careful balance between privacy preservation and efficient model optimization. Despite the establishment of various FedIT frameworks (Ye et al., 2024b; Kuang et al., 2023; Zhang et al., 2023c), existing literature has not adequately addressed the practical challenges that Federated Domain-specific Instruction Tuning (FedDIT) may encounter in real-world applications. For instance, FedIT generally necessitates

a sufficient amount of instruction data for fine-tuning, which is often a shortage in domain-specific fine-tuning contexts (Zhang et al., 2024b).

In this study, we investigate FedDIT, a novel approach within the FL paradigm aimed at boosting 057 the performance of LLMs in specific domains. Unlike general FedIT, which seeks to enhance model effectiveness across diverse tasks without accounting for local data shortage, FedDIT encounters 059 the unique challenge of clients possessing only a limited quantity of local domain-specific data. To 060 overcome this, FedDIT enrichs the local data through a specific instruction augmentation strategy. 061 This strategic enrichment is crucial for achieving effective instruction tuning and needs to be metic-062 ulously designed to avoid performance degradation. Except for the code domain, which primarily 063 adheres to a standardized paradigm, our results reveal that when clients rely solely on their local 064 data for FedDIT, even the presence of high-quality in-domain local data can be insufficient due to limited scale, leading to a decline in performance, as reflected in the underlined values in Table 1. In 065 summary, the goal of FedDIT is to develop a domain-specific LLM that employs collaborative train-066 ing and instruction augmentation while safeguarding client privacy, thereby ensuring the model's 067 proficiency in executing tasks pertinent to its designated domain. 068

For simplify and better instruction quality (Zhang et al. (2024b); Toshniwal et al. (2024)), we focus
on a specific scenario of FedDIT, where a server-hosted public dataset exists as an abstraction of
open-source instruction datasets on the web that encompasses multiple domains. This dataset is utilized for different instruction augmentation strategies, thereby enhancing the model's performance
in specific domains. We elaborate the setting of FedDIT in Appendix A.2.

074 Additionally, the factors affecting FedDIT are still unclear. Compounding this uncertainty, introduc-075 ing augmented instructions may further complicate results, making it difficult to ascertain effective 076 improvement strategies. Shepherd (Zhang et al., 2023c) approaches this problem from the perspec-077 tive of heterogeneity, constructing heterogeneity based on the topic of general instruction datasets. It demonstrates that, unlike the consensus of traditional FL, for general FedIT, heterogeneity has 078 a positive effect. By aggregating diverse instructions from clients, the diversity increases, thereby 079 enhancing the model's adaptability to various tasks. However, it just scratches the surface and does 080 not explore issues in FedDIT. 081

Going one step further, we conduct experiments to unveil a significant finding (Appendix A.3): there
is no monotonic correlation between the degree of non-independent and identically distributed (nonid) and LLM's performance in the context of FedDIT. Inspired by Explore-Instruct (Wan et al.,
2023), which shows the potential of domain coverage in domain-specific instruction tuning. The
cross-client domain coverage metric is initially defined, followed by an investigation into its impact
on FedDIT. Results demonstrate that domain coverage significantly influences model performance
in the corresponding domain.

To maximize the cross-client domain coverage without compromising client data privacy, we pro-089 pose a novel FedDIT algorithm, Federated Instruction Tuning of LLMs with Domain Coverage 090 Augmentation, termed FedDCA. This algorithm employs a greedy client center selection process 091 and implements instruction augmentation through dense retrieval on the server side. The fundamen-092 tal idea of FedDCA is to select client centers to expand the diversity and coverage of augmented 093 instruction datasets within a specific domain. By strategically optimizing domain coverage at each 094 step, FedDCA efficiently constructs the augmented train set that enhances both the learning and 095 generalization capabilities of the model, leading to superior performance on domain-specific tasks. 096 Furthermore, to mitigate computational overhead on the client side and enhance the system scalabil-097 ity, we propose FedDCA^{*}, which employs heterogeneous encoders of different sizes and capacities. 098 To achieve feature alignment, we train a projector on the server side using public data and employ 099 contrastive learning techniques.

100 We demonstrate the effectiveness of FedDCA through comprehensive experiments conducted across 101 four domains: code, medical, financial, and mathematical. These are compared against a range of 102 baselines, which can be categorized into unaugmented and augmented methods. In the unaugmented 103 setting, our method is compared with FedIT, which includes four orthodox FL techniques: FedAvg 104 (McMahan et al., 2017), FedProx (Li et al., 2020), SCAFFOLD (Karimireddy et al., 2020), and 105 FedAvgM (Hsu et al., 2019). In the augmented setting, we compare FedDCA against methods such as random sampling, direct retrieval, LESS (Xia et al., 2024), and Self-Instruction (Wang et al., 106 2022). Additionally, we present the performance outcomes of FedDCA when applied under various 107

FL strategies. We also compare the computational efficiency on the client side between FedDCA and its variant, FedDCA^{*}. For privacy analysis, experiments against memory extraction attacks are conducted to evaluate how different quantities of retrieved public data affect the privacy of client local data. Results indicate that while reliance on local data increases memorization of sensitive information, the risk of privacy leakage diminishes or converges in the augmented setting as the training rounds progress.

- The main contribution is as follows:
 - We reveal a critical finding: in the context of FedDIT, data heterogeneity has a nonmonotonic relationship with model performance. Instead, cross-client domain coverage substantially impacts LLM's effectiveness, as elaborated in Appendix A.3.
 - We propose a novel FedDIT algorithm (Section 4), termed FedDCA, aimed at maximizing cross-client domain coverage through greedy client center selection followed by retrieval-based instruction augmentation executed on the server. Additionally, we introduce FedDCA* to further lessen the client-side computational overhead while enhancing the system scalability. This variant utilizes a heterogeneous encoder structure, paired with a projector on the server side for feature alignment.
 - Through extensive experiments (Section 5), we demonstrate the effectiveness of FedDCA and FedDCA*. We also investigate privacy preservation against memory extraction attacks, conducting experiments based on various amounts of public data. Results suggest that the capacity for privacy preservation does not correlate significantly with the quantity of public data. In contrast, the risk of privacy leakage tends to decrease or converge as the fine-tuning rounds increase.
- 129 130 131 132

133

116

117

118 119

121

122

123

124 125

126

127

128

2 PRELIMINARIES

FedDIT aims to leverage cross-client private domain-specific instruction data and utilize the multi-134 domain public data on the server to achieve instruction augmentation, collaboratively enhancing the 135 model's performance in specific domains. Consider N distributed clients, each with local private 136 data D_k^l . The server maintains a public dataset D^p that encompasses multiple domains and is re-137 sponsible for implementing data augmentation strategies and aggregating model parameters received 138 from clients. The training process follows the standard FL protocol (McMahan et al., 2017). For 139 computational efficiency, we adopt Low-Rank Adaption (LoRA) (Hu et al., 2022) as the fine-tuning 140 method, which involves tuning additional parameters $\Delta \phi$ while keeping the pre-trained LLM's pa-141 rameters ϕ frozen. In the initial training phase, the server dispatches ϕ and $\Delta \phi$ to each client for \mathcal{R} 142 training rounds. In the t-th round, the server sends the aggregated $\Delta \phi^t$ to clients. Clients use $\Delta \phi^t$ 143 to update their local LoRA parameters $\Delta \phi_k^t$ and conduct instruction tuning based on augmented instruction datasets D_k . Subsequently, the clients return $\Delta \phi_k^{t+1}$ to the server. The server then aggregates $\{\Delta \phi_k^{t+1} \mid k = 1, \dots, N\}$ to obtain $\Delta \phi^{t+1}$ for the next round. 144 145

146 147

148

3 PROBLEM FORMULATION

149 For better understanding, we list the frequently used notation in Appendix A.12. The objective of 150 FedDIT is to enhance the domain-specific performance of LLMs through FL without sharing private 151 data (Ye et al., 2024b; Zhang et al., 2024b; 2023c). Under the FL framework, suppose we have N152 clients, where each client c_k has a local dataset D_k^l with its size N_k^l , and an augmented dataset D_k^g 153 from the server public dataset D^p , respectively. Due to constraints such as memory, computational 154 overhead, and maximum tolerated training time, client c_k can accept at most N_k^p public instructions. 155 Denote the augmentation strategy as Λ , through which the server performs instruction augmentation on the public dataset. If Λ is null, it indicates that clients conduct FedDIT solely based on their 156 local private data, which may lead to performance degradation. Conversely, Λ may be classified into 157 two categories: (1) focusing exclusively on the client's own local data distribution or (2) considering 158 the cross-client data distribution. In conclusion, the global objective of FedDIT is defined as follows: 159

$$\arg\min_{\Delta\phi} \left\{ F(\phi, \Delta\phi) \triangleq \sum_{k=1}^{N} p_k \left((1 - \alpha_k) F_k \left(\phi, \Delta\phi_k; D_k^l \right) + \alpha_k F_k \left(\phi, \Delta\phi_k; D_k^g \right) \right) \right\}, \quad (1)$$

162 where the first term and second term denote the accumulated fine-tuning loss computed on the client 163 c_k 's local instructions D_k^l and the augmented public instructions D_k^g from the server, respectively. 164 p_k is the weight of the k-th client, which is determined by the ratio of k-th client's data size to all clients' total data size. α_k is the ratio of the public data amount of its augmented instruction dataset D_k , which is computed as $\frac{N_k^p}{N_k^l + N_k^p}$. Eq.1 minimizes the summed empirical loss across clients' 165 166 167 augmented instructions to pursue the in-domain utility of the obtained global model. Denote P as 168 the model parameters and \mathcal{D} as a specific dataset, then the client c_k 's empirical loss $F_k(\phi, \Delta \phi_k; \mathcal{D})$ 169 base on \mathcal{D} is calculated as: 170

172

173

174 175

176 177 where $x_j \in D, \forall j \in \{1, 2, ..., |D|\}$. The instruction tuning loss $l(\cdot; \cdot)$ on a sample (x, y) is defined as $-\sum_{t=1}^{|y|} \log (w(y_t | x, y_{< t}))$, where x is the formated instruction with Alpaca instruction template¹ and y is the corresponding response.

 $F_k(\phi, \Delta \phi_k; \mathcal{D}) \triangleq \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} l(P_{\phi + \Delta \phi_k}; x_j),$

(2)

In contrast to the problem formulations in many prior works on FedIT, the primary distinction in the formulation of FedDIT lies in acknowledging the lack of in-domain instructions across clients. Fed-DIT establishes a public multi-domain dataset on the server for domain-specific instruction augmentation by a specific sampling strategy Λ . Therefore, each client only needs to hold a few high-quality domain-specific private data to collaborate with other clients to obtain a strong domain-specific LLM.



Figure 1: Overview of FedDCA, which consists of three stages: 1) The client c_k performs local instructions clustering and sends cluster centers C_k to the server. 2) The server first does greedy client center selection to maximize domain coverage and performs client-center-based domain data retrieval, then sends the augmented instructions D_k^g to the client c_k . 3) Clients fine-tune the LLM collaboratively, and the LoRA parameters $\Delta \phi$ are exchanged between clients and the server.

202 203 204 205

206

199

200 201

4 Method

207 We propose Federated Instruction Tuning of LLMs with Domain Coverage Augmentation (Fed-DCA), which enhances domain coverage to obtain a LLM that performs well on domain-specific 208 tasks. FedDCA follows the standard FedIT protocol(Ye et al., 2024b) to perform federated in-209 struction tuning. The novel design of FedDCA lies in the phase before training, which consists of 210 two key modules: greedy client center selection and domain data retrieval, which are elaborated 211 in the following; also see details in Algorithm 1 and Figure 1. For computational efficiency and 212 system scalability, we introduce a variant of FedDCA with heterogeneous encoder setting, named 213 FedDCA*. We first formulate the optimization problem to solve for FedDCA. 214

¹https://crfm.stanford.edu/2023/03/13/alpaca.html

216 4.1 Optimization Problem 217

218 As domain coverage directly affects the in-domain performance of the LLM (Appendix A.3), Fed-DCA aims to maximize the domain coverage of the cross-client augmented data $\bigcup_{i=1}^{N} D_i$ respect to 219 the in-domain data distribution D^d , which is defined in Eq. 5. However, as clients can not send the 220 local data to the server for privacy, directly finding a cross-client dataset that maximizes the domain 221 coverage in Eq. 5 is unrealistic. To find a proper client center set \mathcal{P} that maximizes domain coverage 222 and uses it to perform instruction retrieval on the public data is an approximation problem.

224 Suppose there exists a metric space \mathcal{X} , a set of cross-client local instruction embeddings $\mathcal{E} \in \mathcal{X}$, 225 and a set of cluster centers $\mathcal{P} \in \mathcal{X}$. Each client c_k has maximum ξ clusters obtained by k-means algorithm (Wu, 2012). In the federated setting, communication cost is always a critical factor. Con-226 sequently, we formulate the optimization problem to include communication costs as follows: 227

228 229 230

231

240 241

245

247

250

251

252

$$\underset{\mathcal{P}}{\operatorname{arg\,min}} \left\{ \sum_{i=1}^{N} |\mathcal{P}_i| + \sum_{d \in D^d} \left(\min_{p \in \mathcal{P}} sim(d, p) \right) \right\},\tag{3}$$

s.t. $|\mathcal{P}_i| \leq \xi, \forall i \in \{1, 2, \dots, N\}$. The second term is the domain coverage of the selected client 232 center set \mathcal{P} , and $sim(\cdot, \cdot)$ is the cosine similarity function. However, this optimization problem is 233 NP-hard. Specifically, to select N client center for retrieval, the time complexity is $\mathcal{O}\left(C_{\xi N}^{N}(\mathcal{N}N)\right)$, 234 235 where \mathcal{N} is the size of the public data D^p . As it is a factorial equation, the computational cost 236 explodes with N. What's worse, usually, there are enormous amounts of public data on the server, 237 which makes a huge \mathcal{N} . For computational efficiency, we propose a greedy algorithm to solve this problem in $\mathcal{O}(N(N\xi + \xi \log \xi))$, which is described below. 238 239

4.2 GREEDY CLIENT CENTER SELECTION & DOMAIN DATA RETRIEVAL



Figure 2: Greedy client center selection iteratively selects the client center in each step, which synthetically considers both the representativeness of the cluster center and its distance from the previously selected client center set \mathcal{P} to maximize the cross-client domain coverage.

253 As described in Section 4.1, the goal of FedDCA is to find a proper client center set \mathcal{P} that maximizes 254 the domain coverage $d(D^d, \mathcal{P})$, as defined in Eq. 3. For computational efficiency, we propose a 255 greedy algorithm as shown in Algorithm 1 and Figure 2 to solve this problem in polynomial time 256 and obtain a sub-optimal solution. Given the cluster centers $\{C_i, i = 1, 2, \dots, N\}$ received from 257 each client, which are obtained by clustering local instructions. FedDCA on the server consists of two main steps: greedy client center selection and client-center-based domain data retrieval. 258

259 Greedy client center selection. We will consider this problem from two aspects: 1) Select a client 260 center that can represent the distribution of the local data. 2) To optimize the cross-client domain 261 coverage, we filter client centers that are close to the previously selected client centers. First, we randomly choose a client and select the largest cluster center as its center and be the initial of the 262 client center set \mathcal{P} . Then, we iteratively select the client center \mathcal{P}_k of the top cluster size while 263 avoiding selecting the cluster center close to the previously selected client center to maximize the 264 domain coverage of \mathcal{P} . Specifically, for the *i*-th iteration, we filter *i*-top cluster centers that have 265 the largest summed similarity with previously selected client centers and select the center of the rest 266 largest cluster as the *i*-th client center. This procedure is repeated until we have selected N client 267 centers. 268

Domain data retrieval. For each client center \mathcal{P}_k , the server performs dense retrieval (detailed 269 in Appendix A.4) on public dataset D^p to get the top- N_k^p similar public instructions, then sends

Alg	orithm 1 FedDCA: greedy client c	center selection & domain data retrieval
	Parameters: Number of clusters ξ ; E	ncoder model w_{enc} ; Client local datasets $D = \{D_1, D_2, \dots, D_N\}$;
	Public dataset D^p ; Similarity score th	rreshold α ; Number of public data samples retrieved per client N_k^p ;
	Client centers $\mathcal{P} = \{p_1, p_2, \dots, p_N\}$; Pre-encoded public instruction embeddings \mathcal{E} .
1:	for $i \in \{1, 2,, N\}$ do	
2:	$\mathcal{E} \leftarrow \{w_{enc}(x) \mid x \in D_{i,instruct}^{i}$	ion}
3:	$C_i, S_i \leftarrow k\text{-means}(\xi).\operatorname{fit}(\mathcal{E}')$	\triangleright Cluster local instructions, return cluster centers C_i and sizes S_i
4:	end for $C = \{c, c, c$	
5:	Send $\{C_i, S_i \mid i \in \{1, 2,, N\}\}$ to	the server for the greedy client center selection
6:	$\mathcal{P} \leftarrow \{\mathcal{C}_{0,k} \mid k = \arg\max(S_0)\}$	\triangleright Initialize the client center set
7:	for $i \in \{1, 2, \dots, N-1\}$ do	▷ Greedy client center selection
8:	$\mathcal{S} \leftarrow \sum (\mathcal{C}_i \cdot \mathcal{P}^T, \dim = -1)$	▷ Compute summed similarity score of each cluster center
9:	$\mathcal{I} \leftarrow (N-i)$ - $rg \operatorname{sort}(\mathcal{S})$	\triangleright Filter <i>i</i> cluster centers close to the client center set \mathcal{P}
10:	$\mathcal{I}' \leftarrow \arg \operatorname{sort}(-S_i)$	Sort cluster center by cluster size in descending order
11:	$j \leftarrow \text{first element of } \mathcal{I} \cap \mathcal{I}'$	\triangleright Selected cluster center $C_{i,j}$
12:	$\mathcal{P} \leftarrow \mathcal{P} \cup \mathcal{C}_{i,j}$	▷ Update the client center set
13:	end for	
14:	for $i \in \{1, 2, \dots, N\}$ do	Client center based domain data retrieval
15:	$\mathcal{S} \leftarrow \mathcal{P}_i \cdot \mathcal{E}^T$	\triangleright Compute similarity score between \mathcal{P}_i and \mathcal{E}
16:	$\mathcal{S}' \leftarrow \{s \mid s \in \mathcal{S}, s < \alpha\}$	\triangleright Filter instructions with similarity score larger than α
17:	$\mathcal{I} \leftarrow \text{indices of } N_k^p \text{-top in } \mathcal{S}'$	
18:	$D_i^g \leftarrow \{D_j^p \mid j \in \mathcal{I}\}$	
19:	$D_i \leftarrow D_i^l \cup D_i^g$	\triangleright Obtain the augmented instruction dataset D_i
20:	end for	

retrieved public datasets $\{D_1^g, \ldots, D_N^g\}$ to each clients. Specifically, to avoid the overlap between public data and local private data, we set a threshold α to filter the public instructions that have a similarity score larger than α with the client center.

In summary, FedDCA establishes a comprehensive training set that captures the essence and distribution of the domain by iteratively selecting client centers with unique coverage increments. This ensures that selected centers are representative and achieve broad domain coverage. Additionally, domain data retrieval exposes the model to a wider range of in-domain features, enhancing its understanding of domain-specific tasks. This not only enriches the training data but also reduces the risks of overfitting to the limited local data available to each client, ultimately improving performance across various tasks within the domain.

4.3 HETEROGENEOUS ENCODER WITH FEATURE ALIGNMENT

307 For the client-side computational efficiency and system scalability, we propose a heterogeneous 308 encoder method FedDCA* to reduce the client's computational overhead by using a small encoder ω on the client side and a larger encoder ω' on the server side. However, the output dimension of 310 heterogeneous encoders may not be consistent. Therefore, we introduce a projector w_p on the server 311 for feature alignment, as shown in Appendix A.7. We use contrastive learning (Khosla et al., 2020) and train w_p on the public instructions. Specifically, w_p comprises two fully connected layers and 312 a ReLU activation layer in between, projecting ω 's dimension to the output dimension of ω' . For 313 $\forall x_i \in D^p$, where $i \in \{1, 2, \dots, \mathcal{N}\}$, ω outputs embedding h_i, ω' outputs embedding h'_i and w_p 314 outputs embedding φ_i . For input x_i , the positive sample pair is (h'_i, φ_i) and negative sample pairs 315 are $\{\varphi_j, j \neq i\}$. Let the batch size be B, then the training objective is defined as follows: 316

317 318

305

306

$$\mathcal{L} = -\sum_{i=1}^{B} \log \frac{\exp\left(\sin(\varphi_i, h'_i)/\tau\right)}{\sum_{j=1}^{B} \mathbb{I}_{[j \neq i]} \exp\left(\sin(\varphi_i, h'_j)/\tau\right)},\tag{4}$$

319

where sim (\cdot, \cdot) denotes the cosine similarity and τ denotes the temperature parameter.

The heterogeneous encoder setting, by employing a projector w_p for feature alignment, proves effective in FL contexts where clients and the server use different encoder sizes or capacities. This approach facilitates mapping lower-dimensional client features \mathcal{H} to the higher-dimensional server space \mathcal{H}' through a strategic training of w_p with both positive and negative pairs. It adeptly aligns similar feature representations closely while positioning dissimilar ones further apart. As a result, the heterogeneous encoder setting not only enhances system scalability and flexibility but also facilitates the creation of a unified feature space that effectively integrates and utilizes diverse representations collaboratively.

4.4 DISCUSSIONS

332 **Computation.** Through greedy client center selection, FedDCA solves the optimization problem defined in Eq. 3 in the polynomial time and is quite efficient. On the client side, the k-means algo-333 rithm (Wu, 2012) is $\mathcal{O}(|D_k^t|)$. On the server side, as mentioned in Section 4.1, the time complexity 334 of greedy client center selection is $\mathcal{O}(N(N\xi + \xi \log \xi))$. Specifically, for each client, the similarity 335 computation between selected client center set \mathcal{P} and client c_k 's cluster \mathcal{C}_k is $\mathcal{O}(N\xi)$, and the sorting 336 process is $\mathcal{O}(\xi \log \xi)$. For domain data retrieval, given that the sorting algorithm is $\mathcal{O}(\mathcal{N} \log \mathcal{N})$, 337 where \mathcal{N} is the size of the public dataset. Thus the time complexity is $\mathcal{O}(\mathcal{N}(\mathcal{N}\log\mathcal{N}))$. In addition, 338 the heterogeneous encoder setting FedDCA* further reduces the computation overhead on the client 339 side.

Communication. In the domain instruction augmentation stage, each client sends ξ cluster centers to the server. Next, the server performs greedy client center selection and domain data retrieval, then sends retrieved public data D_k^g to the client c_k , whose size is N_k^p . In addition, in the fine-tuning stage, FedDCA follows the standard FedIT procedure, which exchanges the LoRA parameters $\Delta \phi_k$ between clients and the server. Specifically, for the Llama3-8B model, the number of trainable LoRA parameters is just 13.6 million. Compared with the number of frozen pre-trained LLM parameters, the communication for LoRA tuning is quite efficient.

Privacy. Comparing FedDCA with other FedIT methods (Zhang et al., 2024b; Ye et al., 2024b), the difference lies in the greedy client center selection. In this stage, the client only uploads the cluster center to the server, which is the average of embeddings to its cluster. In addition, the potential privacy leakage can be further avoided through homomorphic encryption (Acar et al., 2018), which allows the server to directly compute on ciphertext for matrix multiplication for dense retrieval. For the further concern of the domain inference attack from the server, please refer to Appendix A.6.

353 354 355

361

362

364

372 373 374

375

330

331

5 EXPERIMENTS

To demonstrate the effectiveness of FedDCA and its variant FedDCA^{*}, we conduct extensive experiments across various domains and with several baselines. In Table 2, we highlight five critical aspects of our approach compared with other methods (Xia et al., 2024; Wang et al., 2022; Zhang et al., 2024b). For additional details and results, please refer to Appendix A.

Table 2: Key differences between FedDCA and other baselines. FedDCA demonstrates the ability of: 1) privacy preserving, 2) no API cost, 3) no additional information required, 4) avoiding performance degradation, and 5) aiming at domain coverage optimization. LESS requires additional information, as it needs access to the validation set for gradient-based retrieval.

Method	Privacy Preserving	API Cost	Additional Information	Performance Degradation	Domain Coverage Oriented
FedAvg (McMahan et al., 2017)	 ✓ 	×	×	 Image: A second s	×
Random Sampling	 ✓ 	×	×	 Image: A set of the set of the	×
LESS (Xia et al., 2024)	 ✓ 	×	 ✓ 		×
Self-Instruct (Wang et al., 2022)	×	1	×	×	×
Direct Retrieval	 ✓ 	×	×	×	×
FedDCA (ours)	 ✓ 	×	×	×	×

5.1 EXPERIMENTAL SETUP

Dataset and Evaluation Metrics. To evaluate the performance of FedDCA, we conduct experiments on four domains: code, medical, financial, and mathematical. The detail of constructing the public dataset is described in Appendix A.5. As default, we set the number of clients to 10, 378 while each client has 100 local instructions and obtains 5000 augmented public instructions from 379 the server. We analysis the effect of different retrieval amounts from 100 to 5000 on FedDCA in 380 Appendix A.10. For evaluation, we select a range of datasets, including HumanEval (H-Eval for 381 short) for coding (Chen et al., 2021), MMLU-Med (abbreviated as M-Med) for medical (Hendrycks 382 et al., 2021), and GSM8K for mathematics (Cobbe et al., 2021). Additionally, financial datasets such as FPB, FiQA, and TFNS (Yang et al., 2023a) are utilized. HumanEval is evaluated using Pass@1. 383 For MMLU-Med, FPB, FiQA, and TFNS, we use accuracy as the evaluation metric. While GSM8K 384 is evaluated using exact match (EM). Please see Appendix A.5 for more dataset details. 385

386 **Baselines.** We compare FedDCA and FedDCA^{*} with the following baselines: 1) Unaugmented 387 methods, including zero-shot inference and FedIT, which are composed of four widely used FL 388 methods, including FedAvg (McMahan et al., 2017), FedProx (Li et al., 2020), SCAFFOLD (Karimireddy et al., 2020) and FedAvgM (Hsu et al., 2019). 2) Augmented methods, which include random 389 sampling, direct retrieval, LESS (Xia et al., 2024), and Self-Instruct (Wang et al., 2022). Addition-390 ally, we report FedDCA's performance with different FL strategies in Table 3. Specifically, zero-391 shot inference shows the performance of the pre-trained LLM without FedDIT, which gives the 392 lower performance bound. Direct retrieval is described in Appendix A.4. Self-Instruct augments 393 instruction in a generative way through GPT-3.5-turbo (Sun et al., 2023). Based on the prompt pro-394 vided by Self-Instruct, we define the prompts for generating instructions and responses as shown in 395 Appendix A.8. 396

Table 3: Performance (%) of FedDCA, FedDCA* and other nine baselines on various domains. For
 augmented methods, we use FedAvg as the default FL strategy. We report FedDCA with different
 FL strategies in the last four rows, and FedDCA+FedAvg and FedDCA+FedProx performs better
 across four domains.

Method	Code	Medical	Financial			Mathematical	
	H-Eval	M-Med	FPB	FiQA	TFNS	GSM8K	
	Unaugmen	nted Method	ds				
Zero-shot	29.88	70.60	55.94	18.54	59.21	23.27	
FedAvg (McMahan et al., 2017)	39.03	68.40	58.25	14.18	66.62	47.46	
FedProx (Li et al., 2020)	37.20	69.10	56.51	14.90	66.45	47.15	
SCAFFOLD (Karimireddy et al., 2020)	37.80	70.20	62.71	15.27	66.49	49.27	
FedAvgM (Hsu et al., 2019)	32.32	64.70	68.14	29.27	70.32	46.85	
	Augment	ed Methods					
Random Sampling	32.93	71.30	64.19	13.09	65.53	47.38	
Direct Retrieval	34.14	72.20	66.31	19.11	67.62	50.87	
LESS (Xia et al., 2024)	28.04	71.00	60.56	16.00	61.14	43.13	
Self-Instruct (Wang et al., 2022)	32.92	71.90	59.73	20.67	66.54	50.79	
FedDCA* (ours)	34.75	73.30	67.10	30.54	71.01	51.55	
FedDCA	(ours) with	different F	L strateg	ies			
FedDCA+FedAvg	36.58	74.50	67.24	35.27	73.32	52.46	
FedDCA+FedProx	32.92	72.40	72.93	38.18	77.55	51.25	
FedDCA+SCAFFOLD	39.87	73.20	72.68	33.09	75.50	50.26	
FedDCA+FedAvgM	33.53	68.90	71.45	31.45	72.52	49.76	

⁴²⁰ 421 422

5.2 PERFORMANCE ANALYSIS

423 Performance. We evaluate the performance of FedDCA and compare it with several baselines on 424 four domains (see training details in Appendix A.7). As shown in Table 3, FedDCA outperforms 425 the other nine baselines in all domains, with a substantial improvement from at least 0.84% to the 426 maximum of 29.19% over other baselines. In particular, FedDCA+FedAvg and FedDCA+FedProx 427 performs better across four domains. In addition, to reduce the computation overhead of clients, 428 FedDCA* attempts to utilize heterogeneous encoders. We see that although FedDCA* has a perfor-429 mance drop than FedDCA, it still outperforms other baselines. Furthermore, we report FedDCA's performance with different FL strategies, and we can see that no FL method can keep the leading 430 position in all domains. In specific, FedProx and SCAFFOLD FL strategies perform better in the 431 average performance. Overall, the result shows the effectiveness of FedDCA and FedDCA*.

Additionally, Table 3 can be divided into two parts: unaugmented methods and augmented methods. We can observe that FedIT methods perform well in the code domain, even better than most augmented methods. This could be attributed to two reasons: 1) The code domain is more about following a certain paradigm. 2) As the local data is few, so compared with the same epoch and batch size, the unaugmented methods learn the same data more times, which is kind of not a fair setting. However, FedDCA+SCAFFOLD still surpasses the best baseline in the code domain, further demonstrating the effectiveness of FedDCA.

439 Efficiency. We see that direct retrieval is the best baseline in the average performance. However, it 440 does not consider cross-client domain coverage, resulting in an overlapping retrieved data distribu-441 tion. Self-Instruct (Wang et al., 2022) represents the upper limit performance of the generation-based 442 method, which is restricted to the high expense of cost, limited API call frequency, and heavy quality screening process. On the other hand, FewFedPIT (Zhang et al., 2024b) attempts to augment 443 local data by leveraging the pre-trained Llama2 model to perform self-instructed generation during 444 training. However, two concerns exist: 1) The pre-trained LLM cannot provide effective, stable, 445 and high-quality instruction generation. 2) Generating instructions locally is very costly regarding 446 both time and computational resources. Overall, generating instructions self-instructively is not a 447 satisfactory method for the FedDIT scenario. LESS (Xia et al., 2024) represents the augmentation 448 methods through gradient feature retrieval. However, its performance is underwhelming. Com-449 pounding this issue, LESS presupposes access to the validation set, which is not always available, 450 and necessitates an initial warmup on the public dataset before calculating gradients for both the 451 public and validation data. This process becomes computationally burdensome with large public 452 datasets. In conclusion, FedDCA achieves better performance while requiring lower computational 453 resources and a more relaxed training condition, which shows its efficiency.

454 Domain coverage. Each method's do-455 main coverage is shown in Table 4. Addi-456 tionally, visualization of domain coverage 457 across different strategies can be found in 458 Appendix A.11. For augmented methods, 459 a correlation is evident between higher domain coverage in a specific domain 460 and improved performance of the method 461 within that domain. Given that Fed-462 DCA aims to maximize the domain cover-463 age through greedy client center selection, 464 therefore it achieves the highest domain 465 coverage in all domains, which surpasses 466 other baselines from 4.82% to 21.36% av-467 erage relative improvement. Furthermore, 468 while FedDCA* employs a smaller en-469 coder on the client side to enhance com-

Table 4: Domain coverage of FedDCA, FedDCA* and other baselines on four domains. Since unaugmented methods do not affect domain coverage, we summarize the unaugmented methods here as FedIT, including four orthodox FL techniques: FedAvg, FedProx, SCAFFOLD, and FedAvgM.

Method	Code	Med.	Fin.	Math.
FedIT	0.8126	0.6990	0.8529	0.7871
Random Direct Re. LESS Self-Instruct FedDCA* FedDCA	0.8512 0.9396 0.8509 0.8966 0.9532 0.9766	0.7940 0.8830 0.7737 0.8586 0.8972 0.9348	0.9196 0.9293 0.8917 0.9015 0.9538 0.9815	0.8651 0.8967 0.8352 0.8811 0.9096 0.9320
FedDCA	0.9766	0.9348	0.9815	0.9320

putational efficiency, this leads to a slight compromise in semantic precision, observed as a relative drop of 2.89% in average domain coverage. Nevertheless, FedDCA* still outperforms the best base-line by an average of 2.78%. In brief, the result proves the effectiveness of FedDCA and FedDCA* in domain coverage augmentation.

474 5.3 COMPUTATION ANALYSIS

Heterogenous encoder setting FedDCA* allows the system to be more scalable and flexible based on different client capacities. To show the computational efficiency of FedDCA* compared with FedDCA, we compare these two methods from the following aspects: 1) Encoder model size. 2) Encoding time overhead. We first give the evaluation setup of computation analysis.

Evaluation setup. The encoders used for FedDCA and
FedDCA* on the client side are bge-large-en-v1.5
and all-MiniLM-L6-v2 respectively (detailed in Appendix A.7). Then, we show the model size and the time

Table 5: Computational cost comparison between FedDCA and FedDCA*. We report the encoding time on the public dataset and the model size of each encoder.

Method	Encoding Time	Model Size
FedDCA	15 min 46 s	335 M
FedDCA*	5 min 02 s	22.7 M

486 overhead of encoding the public instructions in Table 5. Compared to FedDCA, FedDCA* has sig-487 nificant computational and time advantages while maintaining acceptable performance (shown in 488 Table 3). This is accomplished by employing heterogeneous encoders and configuring a projector 489 on the server to achieve dimensional mapping.

491 5.4 PRIVACY ANALYSIS

490

492

496

497

498

499

500

501

502

503

508

509

510 511

512

513 514

517

518

519

Next, we evaluate the privacy-preserving capability of different ratios of public data against memory 493 extraction attacks (Carlini et al., 2021; Zhang et al., 2024a), which utilizes the autoregression nature 494 of LLM to extract information from the memory (Xu et al., 2024). 495

Evaluation setup. We focus on one client's instruction tuning in FedDIT, using FedDCA for instruction augmentation with 1000 to 5000 public instructions. We set up 10 clients with full participation for 10 rounds. Specifically, we record the average ROUGE-L score (Lin, 2004) of client c_0 for each round. The designed prompt to extract instructions memorized by the LLM is detailed in Appendix A.9. For each setting, we repeat memory extraction 100 times and report the average ROUGE-L score based on each generated instruction and all local data instructions to evaluate the privacy-preserving capability of different public data amounts. Specifically, denoting the generated \mathcal{N} instructions as \mathcal{I} and the client's local instructions \mathcal{I}^l . The calculation is defined as $\frac{1}{N} \sum_{i=1}^{N} \text{ROUGE} - L(\mathcal{I}_i, \mathcal{I}^l).$



Figure 3: Privacy preservation analysis against memory extraction attack across four domains.

520 **Results.** We report the average ROUGE-L scores for base-data-only and various public data fine-521 tuning in Figure 3(a). The results show no significant correlation between the public data ratio and 522 privacy-preserving capability in the same training round. In addition, only using local data has a 523 higher risk of privacy leakage than augmented methods.

524 Figure 3(b) shows the average ROUGE-L score trends per round for each domain, where augmented 525 fine-tuning uses 5000 public data. Initially, the ROUGE-L scores for augmented settings increase, 526 then decrease or converge, while the base-data-only scores continue to rise, especially in the code 527 domain. This indicates that with more training rounds, base-data-only fine-tuning captures more 528 privacy information, while the privacy leakage risk in augmented fine-tuning decreases or converges.

529 530

531

CONCLUSION 6

532 We reveal that in the context of FedDIT, there exists a non-monotonic relationship between data 533 heterogeneity and model performance, while the cross-client domain coverage has a significant im-534 pact on model effectiveness. In response, we propose a novel FedDIT method called FedDCA, 535 which optimizes the domain coverage through greedy client center selection and retrieval-based in-536 struction augmentation. Additionally, FedDCA* leverages heterogeneous encoders to reduce the 537 client-side computation overhead and improve system scalability. Experiments across four domains demonstrate the effectiveness of FedDCA and the efficiency of FedDCA*. Further privacy analysis 538 indicates that as fine-tuning advances, the risk of private data leakage diminishes or converges in FedDIT with instruction augmentation.

540 REFERENCES

549

560

- Abbas Acar, Hidayet Aksu, A. Selcuk Uluagac, and Mauro Conti. A Survey on Homomorphic
 Encryption Schemes: Theory and Implementation. *ACM Comput. Surv.*, 51(4):79:1–79:35, July
 2018. ISSN 0360-0300. doi: 10.1145/3214303.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pp. 2633–2650, 2021.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared 550 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, 551 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, 552 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, 553 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fo-554 tios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex 555 Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, 556 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob Mc-558 Grew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large 559 language models trained on code. arXiv preprint arXiv:2107.03374, 2021.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- G. Cornuejols, G. Nemhauser, and L. Wolsey. The Uncapicitated Facility Location Problem. Tech nical report, Cornell University Operations Research and Industrial Engineering, August 1983.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Ja cob Steinhardt. Measuring massive multitask language understanding. In International Confer ence on Learning Representations, 2021. URL https://openreview.net/forum?id=
 d7KBjmI3GmQ.
- Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.
- 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. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?
 id=nZeVKeeFYf9.
- Wenxiang Jiao, Jen-tse Huang, Wenxuan Wang, Zhiwei He, Tian Liang, Xing Wang, Shuming Shi, and Zhaopeng Tu. ParroT: Translating during chat using large language models tuned with human translation and feedback. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 15009–15020, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp. 1001. URL https://aclanthology.org/2023.findings-emnlp.1001.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for federated learning. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 5132–5143. PMLR, 13–18 Jul 2020. URL https://proceedings.mlr.press/v119/ karimireddy20a.html.
- Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised Contrastive Learning. In *Advances in Neural Information Processing Systems*, volume 33, pp. 18661–18673. Curran Associates, Inc., 2020.

594 595 596	Weirui Kuang, Bingchen Qian, Zitao Li, Daoyuan Chen, Dawei Gao, Xuchen Pan, Yuexiang Xie, Yaliang Li, Bolin Ding, and Jingren Zhou. Federatedscope-llm: A comprehensive package for fine-tuning large language models in federated learning. <i>arXiv preprint arXiv:2309.00363</i> , 2023.
597 598 599 600	Changho Lee, Janghoon Han, Seonghyeon Ye, Stanley Jungkyu Choi, Honglak Lee, and Kyunghoon Bae. Instruction matters, a simple yet effective task selection approach in instruction tuning for specific tasks. <i>arXiv preprint arXiv:2404.16418</i> , 2024.
601 602	Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning. In <i>Proceedings</i> of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021.
604 605 606	Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated Optimization in Heterogeneous Networks. <i>Proceedings of Machine Learning and Systems</i> , 2(arXiv:1812.06127):429–450, March 2020. doi: 10.48550/arXiv.1812.06127.
607 608 609	Chin-Yew Lin. ROUGE: A Package for Automatic Evaluation of Summaries. In <i>Text Summa-</i> <i>rization Branches Out</i> , pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.
610 611 612 613	Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qing- wei Lin, Shifeng Chen, and Dongmei Zhang. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. <i>arXiv preprint arXiv:2308.09583</i> , 2023.
614 615 616 617	Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-instruct. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=UnUwSIgK5W.
618 619 620 621	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, 2017. ISBN 2640-3498.
622 623 624	Erik Nijkamp, Hiroaki Hayashi, Caiming Xiong, Silvio Savarese, and Yingbo Zhou. Code- gen2: Lessons for training llms on programming and natural languages. <i>arXiv preprint</i> <i>arXiv:2305.02309</i> , 2023.
625 626 627 628	Weiwei Sun, Lingyong Yan, Xinyu Ma, Shuaiqiang Wang, Pengjie Ren, Zhumin Chen, Dawei Yin, and Zhaochun Ren. Is chatGPT good at search? investigating large language models as re-ranking agents. In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> , 2023. URL https://openreview.net/forum?id=3Q6LON8y2I.
629 630 631	Shubham Toshniwal, Ivan Moshkov, Sean Narenthiran, Daria Gitman, Fei Jia, and Igor Gitman. OpenMathInstruct-1: A 1.8 Million Math Instruction Tuning Dataset, February 2024.
632 633	Laurens van der Maaten and Geoffrey Hinton. Visualizing Data using t-SNE. Journal of Machine Learning Research, 9(86):2579–2605, 2008. ISSN 1533-7928.
635 636 637 638 639	Fanqi Wan, Xinting Huang, Tao Yang, Xiaojun Quan, Wei Bi, and Shuming Shi. Explore- Instruct: Enhancing Domain-Specific Instruction Coverage through Active Exploration. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), <i>Proceedings of the 2023 Conference on Empirical</i> <i>Methods in Natural Language Processing</i> , pp. 9435–9454, Singapore, December 2023. Associa- tion for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.587.
640 641 642	Hongyi Wang, Mikhail Yurochkin, Yuekai Sun, Dimitris Papailiopoulos, and Yasaman Khazaeni. Federated learning with matched averaging. In <i>International Conference on Learning Representations</i> , 2020. URL https://openreview.net/forum?id=BkluqlSFDS.
643 644 645	Xiao Wang, Weikang Zhou, Can Zu, Han Xia, Tianze Chen, Yuansen Zhang, Rui Zheng, Junjie Ye, Qi Zhang, Tao Gui, et al. Instructuie: Multi-task instruction tuning for unified information extraction. <i>arXiv preprint arXiv:2304.08085</i> , 2023a.
040 647	Yiming Wang, Yu Lin, Xiaodong Zeng, and Guannan Zhang. Privatelora: For efficient privacy preserving llm. <i>arXiv preprint arXiv:2311.14030</i> , 2023b.

672

673

- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*, 2022.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In International Conference on Learning Representations, 2022. URL https://openreview.net/ forum?id=gEZrGCozdqR.
- Junjie Wu. Advances in K-means clustering: a data mining thinking. Springer Science & Business
 Media, 2012.
- Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kambadur, David Rosenberg, and Gideon Mann. Bloomberggpt: A large language model for finance. *arXiv preprint arXiv:2303.17564*, 2023.
- Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. LESS:
 Selecting influential data for targeted instruction tuning. In *International Conference on Machine Learning (ICML)*, 2024.
- Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and
 Bill Yuchen Lin. Magpie: Alignment data synthesis from scratch by prompting aligned llms with
 nothing. ArXiv, abs/2406.08464, 2024. URL https://api.semanticscholar.org/
 CorpusID:270391432.
- Hongyang Yang, Xiao-Yang Liu, and Christina Dan Wang. Fingpt: Open-source financial large language models. *FinLLM Symposium at IJCAI 2023*, 2023a.
 - Yi Yang, Yixuan Tang, and Kar Yan Tam. Investlm: A large language model for investment using financial domain instruction tuning. *arXiv preprint arXiv:2309.13064*, 2023b.
- Rui Ye, Zhenyang Ni, Fangzhao Wu, Siheng Chen, and Yanfeng Wang. Personalized Federated
 Learning with Inferred Collaboration Graphs. In *Proceedings of the 40th International Confer- ence on Machine Learning*, pp. 39801–39817. PMLR, July 2023.
- Rui Ye, Rui Ge, Xinyu Zhu, Jingyi Chai, Yaxin Du, Yang Liu, Yanfeng Wang, and Siheng Chen.
 Fedllm-bench: Realistic benchmarks for federated learning of large language models. *arXiv* preprint arXiv:2406.04845, 2024a.
- Rui Ye, Wenhao Wang, Jingyi Chai, Dihan Li, Zexi Li, Yinda Xu, Yaxin Du, Yanfeng Wang, and Siheng Chen. Openfedllm: Training large language models on decentralized private data via federated learning. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '24, pp. 6137–6147, New York, NY, USA, 2024b. Association for Computing Machinery. ISBN 9798400704901. doi: 10.1145/3637528.3671582. URL https://doi.org/10.1145/3637528.3671582.
- Kiang Yue, Tuney Zheng, Ge Zhang, and Wenhu Chen. Mammoth2: Scaling instructions from the
 web. *arXiv preprint arXiv:2405.03548*, 2024.
- Mikhail Yurochkin, Mayank Agarwal, Soumya Ghosh, Kristjan Greenewald, Nghia Hoang, and Yasaman Khazaeni. Bayesian nonparametric federated learning of neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), Proceedings of the 36th International Conference on Machine Learning, volume 97 of Proceedings of Machine Learning Research, pp. 7252–7261. PMLR, 09–15 Jun 2019. URL https://proceedings.mlr.press/v97/ yurochkin19a.html.
- Boyu Zhang, Hongyang Yang, and Xiao-Yang Liu. Instruct-fingpt: Financial sentiment analysis by instruction tuning of general-purpose large language models. *FinLLM Symposium at IJCAI 2023*, 2023a.
- Jianqing Zhang, Yang Hua, Hao Wang, Tao Song, Zhengui Xue, Ruhui Ma, and Haibing Guan.
 FedALA: Adaptive Local Aggregation for Personalized Federated Learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(9):11237–11244, June 2023b. ISSN 2374-3468, 2159-5399. doi: 10.1609/aaai.v37i9.26330.

- Jianyi Zhang, Saeed Vahidian, Martin Kuo, Chunyuan Li, Ruiyi Zhang, Guoyin Wang, and Yi ran Chen. Towards building the federated gpt: Federated instruction tuning. *arXiv preprint arXiv:2305.05644*, 2023c.
- Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*, 2023d.
- Xinlu Zhang, Chenxin Tian, Xianjun Yang, Lichang Chen, Zekun Li, and Linda Ruth Petzold.
 Alpacare: Instruction-tuned large language models for medical application. *arXiv preprint arXiv:2310.14558*, 2023e.
- Yiming Zhang, Nicholas Carlini, and Daphne Ippolito. Effective prompt extraction from language models. In *First Conference on Language Modeling*, 2024a. URL https://openreview.net/forum?id=0095CVdNuz.
- Yue Zhang, Leyang Cui, Deng Cai, Xinting Huang, Tao Fang, and Wei Bi. Multi-task instruction tuning of llama for specific scenarios: A preliminary study on writing assistance. *arXiv preprint arXiv:2305.13225*, 2023f.
- Zhuo Zhang, Jingyuan Zhang, Jintao Huang, Lizhen Qu, Hongzhi Zhang, Qifan Wang, Xun Zhou, and Zenglin Xu. FewFedPIT: Towards Privacy-preserving and Few-shot Federated Instruction Tuning. *arXiv preprint arXiv:2403.06131*, 2024b.
- Kun Zhou, Beichen Zhang, Jiapeng Wang, Zhipeng Chen, Wayne Xin Zhao, Jing Sha, Zhichao
 Sheng, Shijin Wang, and Ji-Rong Wen. Jiuzhang3.0: Efficiently improving mathematical reason ing by training small data synthesis models. *arXiv preprint arXiv:2405.14365*, 2024.

756 A APPENDIX

758 A.1 RELATED WORK

760 Federated Instruction Tuning. Instruction tuning has been widely applied across various application areas of large language models (LLM), serving as a key technique to enhance the capabilities 761 and controllability of LLM (Zhang et al., 2023d; Wei et al., 2022). Recently, federated instruc-762 tion tuning (FedIT) has emerged as an effective strategy for the distributed optimization of LLMs, leveraging federated learning (FL) protocols to improve the handling of privacy-sensitive tasks in 764 real-world scenarios. So far, several FedIT frameworks (Ye et al., 2024b;a; Zhang et al., 2023c) have 765 been established to evaluate the effectiveness of FedIT across multiple datasets, tasks, and FL meth-766 ods. While these platforms provide a foundation for research, they have not yet introduced more 767 complex federated algorithms and deeply investigate the challenging problems and factors affecting 768 FedIT, which are crucial for advancing this field. 769

PrivateLoRA (Wang et al., 2023b) addresses privacy and efficiency issues by exploiting the low-rank properties of residual activations to reduce communication costs, significantly lowering the communication overhead through collaborative computation between the server and clients while effectively maintaining the privacy of local data. While FewFedPIT (Zhang et al., 2024b) focuses on the few-shot learning setting in FedIT, using self-generated data by pre-trained LLMs locally to mitigate the paucity of data and first discussing memory extraction attacks within FedIT.

Domain Instruction Augmentation. In the real world, there is an urgent need for training LLMs with specific functionalities (e.g., reasoning capabilities) or domain-specific LLMs (e.g., code (Ni-jkamp et al., 2023; Luo et al., 2024), medical (Zhang et al., 2023e), financial (Yang et al., 2023b;a; Zhang et al., 2023a; Wu et al., 2023), mathematical (Yue et al., 2024; Luo et al., 2023)). Existing works tend to use open-source domain-specific instruction tuning datasets for training. However, the target domain may not always have corresponding ready-made domain-specific instruction datasets. Even if they exist, these datasets are often limited in scale.

Several studies investigate domain-specific instruction augmentation, which can be categorized into 783 three aspects: 1) Reusing human-curated public datasets (Wang et al., 2023a; Zhang et al., 2023f; 784 Jiao et al., 2023; Lee et al., 2024; Xia et al., 2024). For instance, Parrot (Jiao et al., 2023) enhances 785 translation capabilities of LLMs in chat by converting bilingual sentences into instruction-following 786 datasets. Furthermore, works like INSTA (Lee et al., 2024) and LESS (Xia et al., 2024) attempt 787 efficient domain-specific instruction augmentation via dense retrieval. INSTA (Lee et al., 2024) 788 uses instructions without responses for effective retrieval. LESS (Xia et al., 2024) assumes access to 789 a validation set and uses the warmup LLM's gradients of the train and validation sets for retrieval. 2) 790 Using seed tasks for self-instruct (Luo et al., 2024; Wan et al., 2023): Explore-Instruct (Wan et al., 2023), for example, employs activate exploration to tree-model domain tasks from both depth and 791 breadth, increasing seed instructions' coverage of the domain, and subsequently generating broader 792 in-domain instructions through self-instruct. 3) Scaling instructions from the web: Recent works 793 (Yue et al., 2024; Zhou et al., 2024) highlight the immense potential of mining naturally occurring 794 instructions from the internet. Compared to generated data, web-mined instructions exhibit less 795 bias and greater diversity. MAmmoTH2 (Yue et al., 2024) firstly retrieves domain-relevant texts 796 and employs a LLM to extract Q-A pairs, further refining them into instruction-response pairs. 797 Jiuzhang3.0 (Zhou et al., 2024) distills GPT's instruction generation capabilities into a smaller model 798 and then uses it to generate instructions from the internet. In conclusion, models fine-tuned with 799 augmented instructions have shown promising domain-specific capabilities.

To obtain a well-performing LLM in the specific domain within the distributed environment, we utilize a multi-domain dataset as public data on the server side and perform domain-specific instruction augmentation based on the client's local instructions.

804 805

A.2 THE SETTING OF FEDDIT

806 A.2.1 WHY THIS SETTING?

FedDIT enriches the local data through various instruction augmentation strategies (Zhang et al., 2024b; Xia et al., 2024; Wang et al., 2022), which can be divided into two categories: 1) Generative methods, which generate instructions through the pre-trained LLM locally or API. 2) Retrieval-based

methods, which retrieve instructions from the web. The former methods are either compromising
 privacy or high computational overhead. Thus, in this work, we focus on retrieval-based methods to
 utilize the diverse and high-quality public data.

In addition, we abstract the public data as a server-hosted multi-domain dataset. The presence of a public dataset on the server is only one possible scenario. The core innovation of this work lies in maximizing domain coverage, and the proposed algorithm is independent of the presence of a public dataset on the server. Even if the server does not have a public dataset, clients can retrieve public instructions based on the received client centers from the server, thereby achieving data augmentation that maximizes domain coverage.

819 820 Overall, this setting is chosen for simplification, allowing more attention on the federated instruction augmentation algorithms.

821 822

828

A.2.2 THE DEFINITION OF DOMAIN

The concept of "domain" is flexible and hierarchical. Currently, there is no precise definition of a domain. For example, in Explore-Instruct (Wan et al., 2023), "Brainstorming" and "Rewriting" are considered two domains, but they could also be regarded as two tasks under the general domain. For clarity, we adopt a clearer domain classification in this paper (e.g., code, medical, financial, math).

829 A.2.3 THE DISTRIBUTION OF PUBLIC DATA

If distributed clients aim to solve tasks based on existing knowledge, the public dataset will inevitably contain knowledge relevant to those domains. This could come from the original corpus (which can be converted into instruction-response pairs using GPT) or from pre-constructed instruction datasets on the website. So the distribution of public dataset can be categorized as follows:
containing held-in or held-out instructions. The held-in indicate that the public dataset contains instructions of the specific task that clients aim to solve, while the held-out indicate that the public
dataset does not contain this task's instructions.

The default setting in this work is that the public dataset contains held-in instructions. Further, we show the effectiveness of FedDCA on the held-out settings in Appendix A.10.

838 839 840

837

A.3 WHAT TRULY COUNTS IN FEDDIT

This section will first analyze the correlation between different non-iid levels and the model's performance with separate experiments for both single and multiple domains in Section A.3.1. Further, to demonstrate the impact of non-iid on FedDIT, we compare the performance of the global model trained on augmented data based on iid and non-iid cross-client data distribution. Additionally, we suggest that domain coverage is a key factor for FedDIT in Section A.3.2.

847

848 A.3.1 DATA HETEROGENEITY IS NOT MATTER IN FEDDIT

849 Following the traditional approach of constructing different degrees of non-iid, which are widely 850 used in federated learning (Wang et al., 2020; Yurochkin et al., 2019), we adopt Dirichlet distribu-851 tion to construct various heterogeneity and use k-means with the cluster num $\xi = 100$ to pseudo 852 labeling instructions. Dirichlet distribution is affected by the hyperparameter α , which enhances 853 heterogeneity with a smaller α and decreases heterogeneity with a larger α . We choose four widely 854 used heterogeneity, which are $\alpha = [0.01, 0.1, 1, 10]$ (Ye et al., 2023; Zhang et al., 2023b; Li et al., 2021). Figure 4 shows a visualization of the data distribution with different heterogeneity on the 855 code dataset. 856

We perform instruction tuning on different amount of clients with different heterogeneity. Specifically, the client's number are 10 and 100 respectively with 2 randomly selected clients participate in each round. For each domain, we only use the in-domain data and then perform FedIT. The training details are shown in Appendix A.7. We perform the experiments 3 times with different random seeds (42, 43 and 44) and report the average performance and the standard deviation in each domain. As shown in Table 6, the performance of LLM does not decrease due to the increase of data heterogeneity but shows a non-monotonic correlation, which indicates that the performance of LLM does not directly depend on data heterogeneity and other factors that play a key role.

887 888 889

890 891

892

902 903



Figure 4: Visualization of client data distribution with $\alpha = [10, 1, 0.1, 0.01]$ in the code domain.

#Clients	Metric	$\alpha = 10$	$\alpha = 1$	$\alpha = 0.1$	$\alpha = 0.01$
	H-Eval	36.17 ± 1.03	34.54 ± 0.35	32.71 ± 2.75	34.75 ± 2.65
	M-Med	71.80 ± 0.84	71.60 ± 0.70	71.33 ± 0.37	71.56 ± 1.37
10	FPB	66.41 ± 3.31	68.72 ± 4.54	68.39 ± 2.65	70.07 ± 3.50
10	FiQA	22.90 ± 9.34	32.48 ± 9.85	26.42 ± 6.67	38.42 ± 8.78
	TFNS	69.45 ± 2.34	72.99 ± 3.20	71.65 ± 1.64	73.66 ± 2.90
	GSM8K	56.51 ± 0.14	59.28 ± 0.38	56.75 ± 0.23	57.64 ± 0.29
	H-Eval	36.57 ± 1.84	36.57 ± 3.48	34.12 ± 0.26	36.88 ± 0.13
	M-Med	70.73 ± 0.37	72.20 ± 0.45	71.83 ± 0.37	71.60 ± 0.45
100	FPB	67.87 ± 1.13	64.10 ± 2.07	66.71 ± 1.97	67.59 ± 1.42
100	FiQA	33.44 ± 3.01	19.87 ± 5.83	33.93 ± 3.54	33.69 ± 6.58
	TFNS	72.22 ± 0.74	70.13 ± 2.59	73.01 ± 0.22	71.99 ± 0.95
	GSM8K	54.32 ± 0.41	51.27 ± 0.28	58.94 ± 0.36	57.81 ± 0.19

Table 6: Performance (%) of different heterogeneity in each domain with 10 and 100 clients.

A.3.2 DOMAIN COVERAGE: A KEY FACTOR IN FEDDIT

893 Explore-Instruct (Wan et al., 2023) enhances the coverage of domain-specific seed tasks through 894 active exploration, then uses the self-instruct method for instruction data augmentation. This approach highlights the impact of domain coverage on domain-specific instruction tuning. Inspired 895 by Explore-Instruct, we attempt to conduct more in-depth and extensive experiments to study the 896 effect of domain coverage on FedDIT. Firstly, we define the domain coverage of cross-client data in 897 the FL setting. Assume the dataset of in-domain data D^d represents the latent data distribution of 898 this domain and the cross-client data is defined as $D^c = \bigcup_{k=1}^N (D_k^l \cup D_k^g)$. Inspired by the facility 899 location function (Cornuejols et al., 1983), we define the domain coverage of D^c respect to D^d as 900 follows: 901

$$d(D^{d}, D^{c}) = \frac{1}{|D^{d}|} \sum_{d \in D^{d}} \max_{v \in (D^{c} \cap D^{d})} sim(d, v),$$
(5)

904 where sim(d, v) is the similarity between d and 905 v. Specifically, we use the cosine similarity 906 function in FedDCA. Note that we only use 907 the in-domain data in D^c to calculate the do-908 main coverage because the out-of-domain data 909 would mislead the domain coverage evaluation, 910 as shown in Figure 5.

911 To better align with the real-world scenarios,
912 we explore instruction augmentation based on
913 both iid and non-iid cross-client data distribu914 tion and adopt direct data retrieval as described
915 in Appendix A.4 for FedDIT. We set the num916 ber of clients to 10, while each client has 100
917 local instructions and obtains 5000 augmented
918 public instructions from the server.

Table 7: Performance(%) and domain coverage of iid and non-iid settings on different domains. The higher domain coverage correlates with better performance.

Test Set	Perforn	nance (%)	Domain Coverage		
	iid	non-iid	iid	non-iid	
H-Eval	35.36	33.53	0.8538	0.7994	
M-Med	70.20	71.00	0.7800	0.8027	
FPB	58.58	64.19			
FiQA	17.09	19.27	0.8523	0.9327	
TFNS	66.16	69.09			
GSM8K	40.50	38.50	0.9137	0.8448	

918 To construct iid data distribution, we randomly sample 1,000 from 919 multi-domain datasets (code, medical, financial, mathematical, and 920 general), which is detailed in Table 8 and divide them into 10 shards as each client's local data. To construct non-iid data distribution, we 921 922 perform k-means clustering with $\xi = 100$. Each client randomly samples 100 instructions from different randomly selected clusters. 923 Furthermore, for both iid and non-iid settings, direct data retrieval 924 is performed based on the client's local data. 925

926 Table 7 presents the performance of FedDIT on different domains 927 with iid and non-iid settings and shows the domain coverage in four 928 domains. We can observe that both iid and non-iid settings outperform in some domains, but both collectively indicate that higher 929 domain coverage correlates with better performance. 930

931

933

937

938

947

948

932 A.4 DIRECT DOMAIN DATA RETRIEVAL



Figure 5: Misleadning of outof-domain data on domain coverage calculation.

In this section, we first show the detail of instruction-based dense retrieval in Appendix A.4.1, which 934 is both used in direct retrieval and FedDCA. Then, we explain the direct retrieval algorithm in 935 Appendix A.4.2. 936

A.4.1 INSTRUCTION BASED DENSE RETRIEVAL

939 Suppose an instruction dataset \mathcal{D} consists of several instances. Each instance is a 940 (Instruction, Response) pair. For instructions that have Input, we concatenate the Instruction and 941 *Input* as *Instruction*, which is consistent with OpenFedLLM (Ye et al., 2024b). Then we use only the instruction for encoding and dense retrieval. Denote \mathcal{E} as a cluster center and I as an instruction 942 of the public dataset D^p , then we measure the instruction-based similarity score for dense retrieval as 943 follows: **Score** $(\mathcal{E}, I) = sim(\mathcal{E}, w_{enc}(I))$, where $sim(\cdot, \cdot)$ is the cosine similarity function. Based 944 on the computed similarity between \mathcal{E} and each I in the public data, we then select the top- N_{μ}^{J} 945 instructions as the retrieved public data for domain data augmentation. 946

Algorithm 2 Direct domain data retrieval

Parameters:

949 Clients' local datasets $D = \{D_1, D_2, \dots, D_N\}$; Public dataset D^p ; Local datasets $D^l =$ 950 $\{D_1^l, D_2^l, \dots, D_N^l\}$; Number of clusters ξ ; Encoder model w_{enc} ; Pre-encoded public instruc-951 tion embeddings \mathcal{E} . 952 1: $D^g \leftarrow \emptyset$ 953 2: for $i \in \{1, 2, \dots, N\}$ do 954 $\begin{array}{l} \mathcal{E}' \leftarrow \{w_{enc}(x) \mid x \in D_{i,instruction}^l \} \\ \mathcal{C} \leftarrow \text{k-means}(\xi).fit(\mathcal{E}') \end{array}$ 3: \triangleright Encode the local instructions 955 4: \triangleright Cluster local instructions, return cluster centers C 956 5: $\mathcal{S} \leftarrow \mathcal{C} \cdot \mathcal{E}^T$ \triangleright Compute the similarity score between C and E957 $\mathcal{I} \gets \emptyset$ 6: 958 7: for $j \in \{1, 2, ..., \xi\}$ do 959 $\begin{aligned} \mathcal{S}' \leftarrow \{s \mid s \in S_j, \text{index}(s) \notin \mathcal{I} \} \\ \mathcal{I}' \leftarrow \text{indices of the top-} \frac{N_k^p}{\xi} \text{ elements in } \mathcal{S}' \end{aligned}$ ▷ Filter the selected indices 960 8: \triangleright Retrieve the top $\frac{N_k^p}{\xi}$ indices 961 9: 962 > Update the selected indices $\mathcal{I} \leftarrow \mathcal{I} \cup \mathcal{I}'$ 10: 963 11: end for 964 $\begin{array}{l} D_i^g \leftarrow \{D_j^p \mid j \in \mathcal{I}\} \\ D_i \leftarrow D_i^g \cup D_i^l \end{array}$ 12: ▷ Selected public instructions 965 13: \triangleright Obtain the augmented instruction dataset D_i 966 14: end for 967

968 969

970

A.4.2 DIRECT RETRIEVAL

The direct domain data retrieval only utilizes instructions of the client's local data without responses 971 to perform the retrieval-based domain data augmentation. The detailed algorithm is described in 972Algorithm 2. For each client, we start by encoding the local instructions using the encoder model973 w_{enc} . Next, we apply the k-means algorithm to cluster the embeddings into ξ clusters. Then ξ 974cluster centers are sent to the server for retrieval-based domain data augmentation. Subsequently,975the retrieved public data is sent to the client for instruction tuning.

A.5 TRAIN AND TEST DATASET INFORMATION

Here we provide the train and test dataset details of code (CodeAlpaca²), medical (MedAlpaca³),
 financial (FinGPT (Yang et al., 2023a)) and mathematical (MathInstruct⁴) domain, respectively.

Dataset name	Domain	Туре	N _{sample}	Metric
CodeAlpaca	Code	Train	20,022	-
HumanEval	Code	Test	164	Pass@1
MedAlpaca	Medical	Train	33,955	-
MMLU-Med	Medical	Test	1,089	Acc
FinGPT	Financial	Train	76,772	-
FPB	Financial	Test	152	Acc
FiQA	Financial	Test	35	Acc
TFNS	Financial	Test	299	Acc
MathInstruct	Mathematical	Train	224,567	-
GSM8K	Mathematical	Test	1,319	Exact Match
Alpaca	General	Train	52,002	-

Table 8: Dataset information of each domain.

996 997 998

976 977

978 979

981 982

As shown in Table 8, the public data consists of four domain-specific instruction datasets and a 999 general instruction dataset, which are CodeAlpaca, MedAlpaca, FinGPT, MathInstruct, and Al-1000 paca, respectively. For each domain's FedDIT, we randomly select 1000 samples from the in-1001 domain instruction and split them into 10 shards as each client's local dataset. The rest of 1002 the instructions are used as the public dataset D^p . For evaluation, we use HumanEval for 1003 the code domain, MMLU-Med for the medical domain (specifically using subjects anatomy, 1004 clinical_knowledge, college_biology, college_medicine, medical_genetics 1005 and professional_medicine in MMLU), FPB, FiQA and TFNS for the financial domain, and GSM8K for the mathematical domain.

1007

A.6 DISCUSSION ON DOMAIN INFERENCE ATTACK

1010 The server may inference the clients' data domain when the domain data retrieval is performed on 1011 the server side. However, the proposed algorithm FedDCA is independent of the presence of a 1012 public dataset on the server. Even if the server does not have a public dataset, clients can upload 1013 their cluster centers to the server, which selects a set of client centers and sends them back to the 1014 clients. Each client can then retrieve data from the website based on the received client center by 1015 itself, thereby achieving data augmentation while maximizing the cross-client domain coverage.

In that case, since the server does not know the encoder used by the client, it cannot infer the semantic meaning of the embedding. Thus, for the server, it becomes significantly more challenging to infer the client's domain, let alone apply any privacy protection techniques to the embeddings.

1019 1020 We provide two examples for illustration. Two different encoders are used as 1020 client's and server's respectively: BAAI/bge-large-en-v1.5 (denoted as w_1) and 1021 google-bert/bert-large-uncased (denoted as w_2). Both encoders output 1024-1022 dimensional features.

¹⁰²³ 1024 1025

²https://huggingface.co/datasets/sahil2801/CodeAlpaca-20k

³https://huggingface.co/datasets/medalpaca

⁴https://huggingface.co/datasets/TIGER-Lab/MathInstruct

Example 1: Both w_1 and w_2 take "hello world" as input, and the cosine similarity between their embeddings is **0.1829**.

- **Example 2:** Three instructions are used:
 - **Instruction 1:** Create an array of length 5 which contains all even numbers between 1 and 10.
 - **Instruction 2:** Write a replace method for a string class which replaces the given string with a given set of characters.
 - Instruction 3: What is the sentiment of this news? Please choose an answer from {negative/neutral/positive}. Teollisuuden Voima Oyj, the Finnish utility known as TVO, said it shortlisted Mitsubishi Heavy's EU-APWR model along with reactors from Areva, Toshiba Corp., GE Hitachi Nuclear Energy, and Korea Hydro & Nuclear Power Co.

For these instructions, Instruction 1 is passed to w_1 and Instructions 2 and 3 are passed to w_2 , which will result in three embeddings: e_1 , e_2 , and e_3 . The cosine similarity between e_1 and e_2 is **0.1464**, while the similarity between e_1 and e_3 is **0.1879**. Instructions 1 and 2 are in the same domain, whereas they have a lower cosine similarity.

1044 In conclusion, as demonstrated above, when the clients do not perform domain-specific instruc-1045 tion retrieval on the server side, the server cannot infer the client's domain based on the uploaded 1046 embeddings.

1047

1030

1031

1032

1033

1034 1035

1036

1037

1048 A.7 IMPLEMENTATION DETAILS

1049 We consider FedDIT in the cross-device scenario, N = 101050 clients, $\mathcal{R} = 30$ rounds, where we randomly sample 2 clients 1051 to be available for each round. Then, each available client 1052 performs FedDIT for 10 steps with AdamW optimizer, and 1053 the batch size is B = 32 in a round. The initial learning 1054 rate is 5e-5 with a cosine learning rate scheduler. Our ex-1055 periment utilizes the widely used LLM, Llama3-8B⁵ as the 1056 base model with 2048 max sequence length and adopts LoRA tuning method. The rank of LoRA is 16, and the scalar al-1057 pha is 16. For k-means (Wu, 2012), we set cluster num ξ 1058 = 10 and for FedDCA we set the similarity threshold α = 1059 0.7. For FedDCA^{*}, we set the temperature parameter $\tau =$ 0.5 for contrastive learning (Khosla et al., 2020). We utilize 1061 bge-large-en-v1.5⁶ as both the client and server's en-1062 coder as default, which outputs embeddings of 1024 dimen-1063 sions. While we utilize all-MiniLM-L6- $v2^7$ as the client's



Figure 6: The projector used on the server side for feature alignment, which consists of two fully connected layers and a ReLU activation layer in between.

small encoder which outputs embeddings of 384 dimensions and bge-large-en-v1.5 as the server's encoder for FedDCA*.

Projector. To reduce the computational overhead of clients, we propose a heterogeneous encoder method called FedDCA^{*}. We utilize a projector to perform the dimension alignment between the small encoder on the client side and the large decoder on the server side. As shown in Figure 6, where d_l is the dimension of the client-side encoder, d_g is the dimension of the server-side encoder.

1071

1077

A.8 PROMPTS USED IN THE SELF-INSTRUCT DATA GENERATION

To generate the Self-Instruct data, we prompt GPT-3.5 to generate the instruction with the designed prompt in Figure 7. Specifically, we randomly sample two examples from the client's local data to guide GPT-3.5 generating the in-domain instruction and one example from the client's local data for one-shot in-context learning to guide GPT-3.5 generating responses into the example's format.

1078 ⁵https://huggingface.co/meta-llama/Meta-Llama-3-8B

^{1079 &}lt;sup>6</sup>https://huggingface.co/BAAI/bge-large-en-v1.5

⁷https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2

1080	You are asked to come u	p with	Example 1:		
1081	instructions. Don't repea	t instructions in	nstruction: {}		
1082	examples. Here are some	e examples:	Response: {}		
1083	Instruction 1: {}	(Generate the respon	nse of this instruction,	
1084	Instruction 2: {}	I	nstruction: {}		
1085	Provide a new instructio	n below:			
1086					
1087	Figure 7: Prompts used in the Se	If-Instruct data gen	eration. (a) Pron	npt for generating ne	w instruc-
1088	(b) Prompt for generating response	y sampled from the	PT 3.5 to gener	a for in-context demo	randomly
1009	(b) Frompt for generating responses	context learning	r 1-5.5 to genera	ale responses with a	Tandonny
1090	selected example for one-shot m-	context learning.			
1091					
1092	A.9 PROMPT USED FOR MEMO	ORY EXTRACTION	ATTACK		
1094					
1095	As we use Llama3-8B as our base	e model and format	the instructions a	and responses into the	e Alpaca's
1096	format, to utilize the auto-regress	ion nature of LLM	to extract the ins	struction, we prompt	the model
1097	to generate the instruction using	the prompt in Figur	e 8, which is exa	actly the prefix of the	e Alpaca's
1098	template.				
1099				_	
1100	Belo	ow is an instruction	that describes a		
1101	task	. Write a response t	hat appropriatel	у	
1102	com	pletes the request.			
1103	####	Instruction:			
1104					
1105	Figure 8:	Prompt used for me	emory extraction	attack.	
1106					
1107					
1108	A.10 FURTHER ANALYSIS				
1109					
1110	To further study the effect of diff	erent hyperparamet	ers of FedDCA,	we undertake a thore	ough anal-
1111	ysis including various retrieval a	mounts and differe	nt k-means clust	er numbers ξ . In ad	dition, we
1112	perform the ablation study on w	hether using the sir	nilarity threshold	d α in the greedy cli	ent center
1113	selection.				
1114	Held-out Setting. Considering	this setting in the f	inancial domain	, to construct the he	ld-out set-
1115	ting, given that the training set H	FinGPT and the tes	t sets FPB, FiQA	A, and TFNS are all	related to
1116	sentiment analysis tasks. As cli	ents aim to train a	model adept at	performing sentimer	it analysis
1117	through federated learning, we ke	ep the setting of tes	t sets. Meanwhi	le, the FinGPT's insti	ructions in
1118	public data are replaced with dat	a from the Sujet.	-Finance-in	STRUCT-1/K Gala	iset where
1119	vields held-out public data		indoniny sampled		approach
1120	yields held out public dutu.				
1121	Table 9: Performance (%) on the	test set of financia	l domain after 3	0 rounds' federated i	nstruction
1122	tuning.				
1123	C				
1124	Metho	od FP	B FiQA TF	FNS	
1125	7.000	Shot 55 (<u>~</u>)/ 185/ 50	21	
1120	Zero S Race J	Data 58.	25 14 18 66	.21	
1127	Rando	om Sampling 60.	39 9.45 65	.45	
1120	FedD	CA 60.3	39 18.91 67	.37	
1129					
1120					

As shown in Table 9, it can be observed that even when in the held-out setting, FedDCA still achieves performance improvements compared to other baselines. Additionally, using the Random Sampling data augmentation strategy resulted in performance degradation on the FiQA dataset. This further underscores the necessity of selecting an appropriate data augmentation strategy.

1134 Effect of Retrieval Number. We report the performance and the corresponding domain coverage 1135 of FedDCA with different retrieval amounts on the four domains in Figure 9, respectively. We can 1136 see that the domain coverage of FedDCA is increasing along with the retrieval amount in different 1137 trends, as well as the performance. Specifically, the domain coverage of each domain increases by 1138 6.36%, 18.69%, 5.04%, and 8.06% in relative, respectively. Along with domain coverage increasing, the performance of FedDCA is increasing by 3.05%, 2.20%, 4.08%, and 6.95% for each domain. 1139 Noted that although the code domain is more about following a certain paradigm, which could 1140 perform well with a few data and more fine-tuning rounds, it still could benefit from the instruction 1141 augmentation. 1142



Figure 9: Effect of different retrieval amounts on the performance of FedDCA and its domain coverage. We show the results on four domains separately. Here, we use the FPB test set to evaluate the performance in the financial domain.

In addition, to further prove the effectiveness of FedDCA is independent of the retrieval amount, we conduct the experiment that each client samples 100 samples from the public dataset and then performs FedDIT. Random* and FedDCA* represent the default setting, where each client samples 5,000 samples.

Table 10: Performance (%) of FedDCA and Random Sampling with different amounts of sampled public data.

Method	H-Eval	MMLU-Med	FPB	FiQA	TFNS	GSM8K
Base Data	39.03	68.40	58.25	14.18	66.62	47.46
Random Sampling	34.53	69.80	60.89	14.54	65.57	48.77
Random Sampling*	32.93	71.30	64.19	13.09	65.53	47.38
FedDCA	35.97	70.20	63.20	15.63	67.58	49.32
FedDCA*	36.58	74.50	67.24	35.27	73.32	52.46

1173 1174 1175

1165

The results in Table 10 show that even with a small amount of sampled data, the model performance still improves compared to random sampling, except in the code domain, which tends to follow a certain paradigm. This further demonstrates the effectiveness of FedDCA and its independence of the retrieval amount.

How Far can FedDCA Go? Assume there exists a pre-trained model that has already fine-tuned
on the whole in-domain data of the internet or even all open-source public instruction datasets, is
there any room for instruction augmentation or FedDCA? To answer this question, we conduct the
following experiments on four settings: A) Fine-tuning on the whole public dataset. B) Following
A, perform further federated fine-tuning with only private clients' local dataset. C) Following A,
each client randomly samples 100 public data for fine-tuning. D) Following A, each client samples
100 public data through FedDCA and then performs FedDIT.

As shown in Table 11, even with only 100 samples selected via FedDCA, the method not only prevents performance degradation but also helps the model achieve better generalization within the

22

1189				e			
1190	Setting	H-Eval	M-Med	FPB	FiQA	TFNS	GSM8K
1191	A	38.41	69.10	75.33	41.09	71.60	53.75
1192	В	43.90	64.30	77.97	66.54	78.26	56.40
1193	С	41.46	66.90	74.50	29.81	70.51	57.01
1194	D	44.12	70.90	81.43	72.00	78.81	58.90

Table 11: Performance (%) of setting A, B, C, and D in each domain.

1195 1196 1197

1206

1188

domain. This is attributed to exposure to more diverse in-domain data during training. Meanwhile, 1198 it can be observed that random sampling still leads to performance degradation compared to Setting 1199 A, and even performs worse than Setting B, where no data augmentation is applied. In the financial domain, performance degradation is also evident. 1201

In conclusion, this experiment further highlights the necessity and effectiveness of FedDCA and 1202 designing effective data augmentation algorithms for FedDIT. 1203

1204 Scalability. To evaluate the scalability of FedDCA, we conduct the experiment with 100 clients. In 1205 each round, two clients are randomly selected for federated instruction tuning using FedAvg.

1207 Table 12: Performance (%) and domain coverage of FedDCA and other baselines in various domains. The client number is 100, and 2 clients are selected in each round. 1208

1209					
1210		Zero-shot	Base Data	Random Sampling	FedDCA
1211			Performance	e (%)	
1212 1213	H-Eval	29.88	34.14	34.75	35.92
1214	MMLU-Med FPB	70.60 55.94	72.40 66 74	69.90 61.05	73.30 67.16
1215	FiQA	18.54	33.45	12.00	34.51
1216	TFNS	59.21	72.78	65.53	73.26
1218	GSM8K	23.27	49.12	47.23	50.26
1219			Domain Cove	erage	
1220	Code	-	0.8282	0.8685	0.9242
1221	Med. Fin	-	0.8377	0.8497	0.9090
1223	Math.	-	0.8709	0.8812	0.9118

1224

1225 As shown in Table 12, as the number of clients increases, the amount of local data on each client 1226 gradually grows. The model trained using only base data even outperforms random sampling in domains other than code and narrows the gap with FedDCA. However, because FedDCA aims to 1227 maximize the cross-client domain coverage, it achieves higher domain coverage and better perfor-1228 mance. In conclusion, results further demonstrate the effectiveness and scalability of FedDCA. 1229

1230 **Impact of Different Cluster Number.** The hyperparameter ξ is the number of clusters in the k-1231 means algorithm. The experiment is conducted on $\xi = [N, 2N, 4N, 8N]$, where N is the number 1232 of clients. Following the setting in Appendix A.7, we set N = 10. We report the domain coverage 1233 of the augmented dataset via FedDCA with different cluster numbers ξ on the four domains in Figure 10, respectively. Results show that there is no best ξ for all domains. Specifically, the best ξ 1234 of code, medical, financial, and mathematical domains are 80, 80, 40, and 10, respectively. 1235

1236 Ablation Study. We conduct the experiment with FedDCA w/o similarity threshold α on the four 1237 domains based on the FedAvg FL strategy. The performance and the corresponding domain coverage are shown in Table 13, where FedDCA without the similarity threshold α is marked as FedDCA[†]. The result shows that the performance of FedDCA with similarity threshold α is slightly better than 1239 FedDCA without using α in code, financial, and mathematical domains, as the similarity scores 1240 in the medical domain are relatively lower. We show the similarity score distribution of the four 1241 domains in Figure 11. For each domain, we plot each similarity score's distribution of 10 clients.



Figure 10: Impact of different cluster number ξ on the cross-client domain coverage. We report the domain coverage of the augmented dataset via FedDCA with different cluster numbers ξ on the four domains.

The similarity score is computed between the selected client center and the public data. Then, we show the similarity score distribution using the histogram plot.

Table 13: Ablation study on the performance and domain coverage of FedDCA w/o similarity threshold α .

Metric	Performa FedDCA [†]	ance(%) FedDCA	Domain (FedDCA [†]	Coverage FedDCA
H-Eval	35.97	36.58	0.8972	0.9348
M-Med	73.40	73.40	0.9348	0.9348
FPB	66.25	67.24		
FiQA	23.27	35.27	0.9353	0.9815
TFNS	69.34	73.32		
GSM8K	51.78	52.46	0.9128	0.9320

A.11 AUGMENTATION STRATEGY VISUALIZATION

To more intuitively compare the domain coverage of different instruction augmentation methods, we randomly sample 5,000 instructions obtained through these methods and 10,000 in-domain instruc-tions as the background, representing the distribution of specific domains in the public dataset. We then visualized the results using t-SNE (van der Maaten & Hinton, 2008), as shown in Figure 12. The plot shows that FedDCA encompasses most of the in-domain data, which is consistent with FedDCA's domain coverage of each domain shown in Table 4. Also, we can observe that the ran-dom sampling strategy selects a lot of out-of-domain data while does not have good coverage in specific domains.

1281 A.12 FREQUENTLY USED NOTATION



Figure 11: The similarity score distribution of the four domains. For each domain, we plot each similarity score's distribution of 10 clients.



Figure 12: Visualization of cross-client data distribution in different domains, performing t-SNE dimensionality reduction on retrieved instructions through various augmentation strategies. We ran domly sample 10,000 in-domain samples as background while randomly sampling 5,000 samples from the cross-client augmented dataset for different instruction augmentation methods for comparison.

used notation.
used notation.
REMARK
lients number, client set $C = \{c_1, c_2, \dots, c_N\}$
oss-client dataset, client's local private data.
umber of local private data on the k -th client
Imber of the retrieved public data on the k-the ient
the local private data on the k -th client, the
trieved public data on the k-th client, the
specific sampling strategy that performs
struction augmentation on the server side,
lected client center set, the cluster centers
eedy client center selection.
LM's pre-trained parameters, additional LoR
ferged model parameters from the frozen LL
trameters ϕ and the additional LoRA parameters
ϕ , encoder model, projector model.
structon tuning loss of model w over a data
mple (x, y) .