

Out-of-Distribution Federated Distillation with Domain-Aware Proxy Selection

Anonymous ACL submission

Abstract

Federated Learning is a distributed machine learning paradigm that trains a global model by aggregating local clients without sharing private data of each client. Federated Distillation (FD) builds upon this paradigm by leveraging knowledge distillation to exchange soft predictions on proxy data instead of model parameters, enabling more efficient communication and supporting heterogeneous model collaboration. However, FD models trained on In-Distribution data are hardly adapted to Out-of-Distribution (OOD) scenarios. In this paper, we propose a domain-aware proxy selection framework to better adopt proxy data for OOD problems. The experimental results show that the proposed models effectively address the challenges of distribution shifts under OOD with and without proxy data by achieving average 82.9% and 81.0% over existing works on standard benchmarks. The codes and data are released in <https://anonymous.4open.science/r/DPS-FD-8596/>.

1 Introduction

Federated Learning is an emerging distributed machine learning paradigm that trains a global model by aggregating local clients without sharing private data kept by each client (McMahan et al., 2017; Abadeer et al., 2022; Pei et al., 2024; Chen et al., 2024; Wang et al., 2025a). However, standard Federated Learning has the limitations of the high communication cost and the homogeneous customization. To address these limitations, Federated Distillation (FD) introduces shared proxy data as a medium for knowledge transfer between clients and the server (Hinton et al., 2015; You et al., 2017; Li and Wang, 2019). By exchanging soft labels on proxy data instead of model parameters, FD significantly reduces communication overhead and enables the use of heterogeneous models across clients (Lin et al., 2020; Hu et al., 2021; Wu et al.,

2022; Itahara et al., 2023; Shao et al., 2024; Fan et al., 2025; Wang et al., 2025c).

However, FD models trained on In-Distribution data are hardly adapted to Out-of-Distribution (OOD) scenarios, due to the proxy data that serves as a critical medium for knowledge transfer (Hendrycks and Gimpel, 2016; Guo et al., 2023; Bai et al., 2023; Qi et al., 2024, 2025; Jeong and Choi, 2025). Since proxy data heavily affect distribution shifts under OOD, we propose **DPS-FD**, a **Domain-aware Proxy Selection** framework for **FD** to better adopt proxy data for OOD problems. Specifically, DPS-FD enables each local client to select the relevant domain samples, while the server selects globally-representative samples from the proxy data. By taking the union of these selected samples, DPS-FD constructs a domain-aware proxy data that better captures the characteristics of both local and global distributions.

Although the effective adoption of proxy data can mitigate OOD problem in FD, the original proxy data sometimes cannot be obtained for knowledge transfer in FD (Zhou and Chiam, 2023; Zhu, 2024; Fang et al., 2024). These works mainly focus on computer visions (Takahashi et al., 2023; Liao et al., 2024b; Qi et al., 2025; Wang et al., 2025b). To address the challenge of no proxy data in textual data, we propose a **Vocabulary-Constrained LLM-based generation (VC)** strategy for generating proxy data, where a global vocabulary is adopted as lexical constraints to guide the proxy data generation. VC adopts a few-shot prompting strategy that consists of two complementary prompt templates: system prompts and user prompts, to generate diverse and high-quality proxy data. According to the comprehensive empirical experiments, we find that high-quality proxy data can significantly mitigate OOD problems. DPS-FD generates proxy data using the VC strategy if the proxy data is unavailable.

Experiments demonstrate that DPS-FD achieves

083 competitive performance both with and without
084 proxy data. Additionally, we take insight analysis
085 on the effect of the high-quality diverse proxy data,
086 unveiling that global distribution of the proxy data
087 heavily affect the OOD performance of the global
088 model in server. The main contribution of this
089 paper are as follows:

- 090 • We propose a domain-aware proxy selection
091 framework that better adopts proxy data for
092 OOD problems, enhancing the robustness and
093 generalization of the global model.
- 094 • We introduce a vocabulary-constrained LLM-
095 based generation strategy to enable FD models
096 to maintain competitive performance without
097 access to real proxy data.
- 098 • We take the insight analysis on the role of
099 proxy data in FD, revealing that high-quality
100 proxy data are essential for mitigating OOD
101 problems and enhancing the robustness of the
102 global model.

103 2 Related Work

104 2.1 Federated Learning and Distillation

105 FL emerges as a promising distributed learning
106 paradigm that enables collaborative model training
107 without sharing private data of each client. Clas-
108 sic approaches such as FedAvg (McMahan et al.,
109 2017) aggregate model parameters from clients into
110 a global model and iteratively repeats this process.
111 However, they suffer from high communication
112 overhead, require homogeneous model architec-
113 tures, and privacy-leakage across clients and the
114 server (Sui et al., 2020; Lin et al., 2021). To mit-
115 igate these limitations, FD leverages knowledge
116 distillation (Hinton et al., 2015; You et al., 2017;
117 Anil et al., 2018) to exchange soft labels on shared
118 proxy data instead of model parameters (Wu et al.,
119 2022; Itahara et al., 2023; Chen et al., 2024; Wang
120 et al., 2025c), significantly reducing communica-
121 tion costs, enabling heterogeneous models, and
122 preserving privacy.

123 2.2 OOD in Federated Distillation

124 OOD behavior manifests primarily as covariate
125 shifts and semantic shifts (Liao et al., 2024b),
126 which often leads to model degradation in FD (Gul-
127 rajani and Lopez-Paz, 2020). To address this is-
128 sue, FD extends to OOD scenarios by enhancing
129 the alignment between proxy and client distribu-
130 tions (Yu et al., 2023; Qi et al., 2025). Existing

work leverages public or synthesized proxy data to
131 reduce domain gaps and improve knowledge trans-
132 fer (Jeong et al., 2023). Zhu et al. (2021) introduce
133 adaptive weighting and ensemble strategies to em-
134 phasize client-specific contributions during distilla-
135 tion and enhance robustness. These developments
136 highlight FD potential in improving generalization
137 in OOD environments. To further advance this line
138 of research, we propose DPS-FD, which constructs
139 high-quality proxy data to explicitly model distribu-
140 tional diversity across clients, enabling more effec-
141 tive and efficient knowledge transfer and reducing
142 communication costs.

143 2.3 Proxy Data in Federated Distillation

144 Proxy data plays a crucial role in FD, serving as
145 the essential medium through which knowledge
146 is exchanged between heterogeneous clients and
147 the global model (Li and Wang, 2019; Lin et al.,
148 2020). Without proxy data, the distillation process
149 struggles to align knowledge (Liao et al., 2023,
150 2024a; Xiao and Liu, 2025), resulting in perfor-
151 mance degradation.

152 To address these limitations, existing studies ex-
153 plore various generative approaches to synthesize
154 proxy data for FD. These include 1) logit-based in-
155 version methods, which reconstruct pseudo data by
156 optimizing inputs to match client logits (Takahashi
157 et al., 2023); 2) generator-based approaches using
158 GANs or VAEs to learn data distributions and pro-
159 duce representative samples (Zhang et al., 2022;
160 Wang et al., 2023; Liao et al., 2024b; Ma et al.,
161 2025; Qi et al., 2025); 3) diffusion-based models,
162 which iteratively denoise random noise into high-
163 quality synthetic data (Li et al., 2023; Wang et al.,
164 2024; Yang et al., 2024; Wang et al., 2025b).

165 However, in natural language processing, the
166 situation becomes more challenging. Most existing
167 proxy-based distillation methods are designed for
168 vision tasks and overlook fundamental linguistic
169 properties such as lexical frequency distribution
170 and semantic structure. This gap highlights the
171 need for domain-aware proxy construction methods
172 tailored to textual data and motivates our work on
173 leveraging lexical constraints and large language
174 models to build more effective proxy data in FD.

175 3 Preliminary

176 We take classification as an example. In the context
177 of a classification task under FL with K clients and
178 a central server, each client $k \in \{i = 1, 2, \dots, K\}$
179

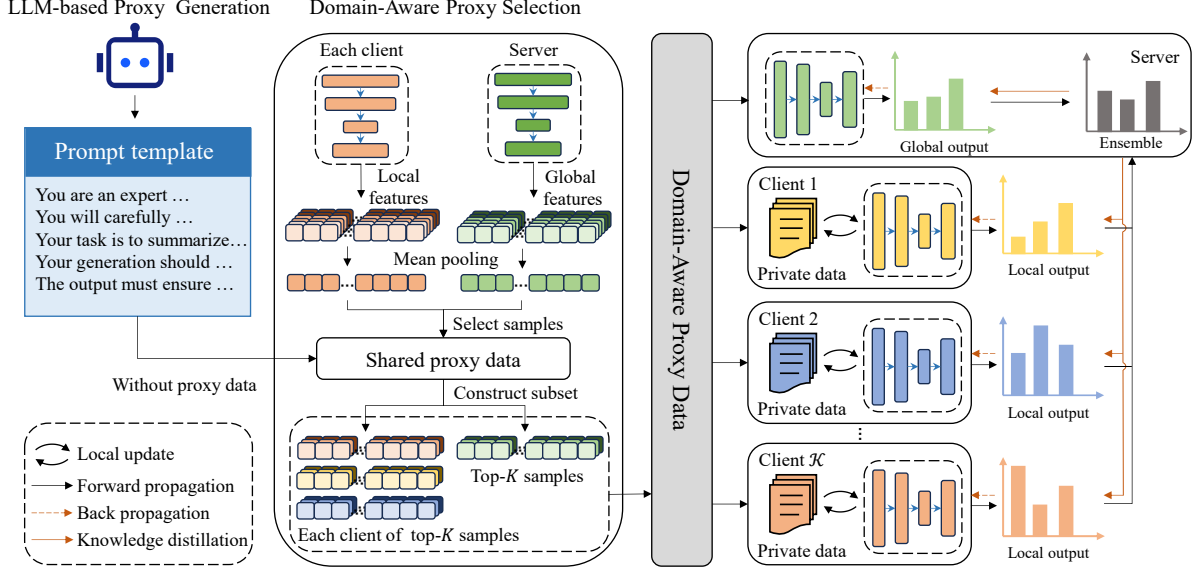


Figure 1: The framework of DPS-FD.

locally keeps its own private labeled data $D_k = \{(x_i^k, y_i^k)\}_{i=1}^{N_k}$ that is isolated from others, where $x_i^k \in \mathbb{R}^d$ and $y_i^k \in \{1, 2, \dots, C\}$ are the i -th instance and its corresponding label, respectively, C is the number of classes, and N_k is the size of private data. The objective is to train a globally optimal model f_g parameterized by θ_g through aggregating client models f_k parameterized by θ_k in a decentralized manner without exposing private data on the server.

To mitigate the limitations of FL based on parameter exchange, FD leverages knowledge distillation to exchange soft labels instead of model parameters by introducing shared proxy data $D_p = \{(x_i^p)\}_{i=1}^{N_p}$ for knowledge transfer between clients and the server, where N_p is the size of the proxy data. Specifically, in each communication round $t \in T$, the server randomly selects the set of activated clients $\mathcal{K} \subseteq \{1, 2, \dots, K\}$ based on a sampling fraction ϵ to participate in the FD training, where $|\mathcal{K}| = \lfloor \epsilon \cdot K \rfloor$. Each client \mathcal{K} first trains its local model on its own private data D_k , and the local optimization objective is defined as follows:

$$\theta_k^* = \arg \min_{\theta_k} \mathbb{E}_{(x,y) \sim D_k} [\mathcal{L}_{\text{CE}}(f_k(x; \theta_k), y)], \quad (1)$$

where $\mathcal{L}_{\text{CE}}(\cdot)$ denotes the cross-entropy loss used for local classification tasks, and \mathbb{E} denotes its expected value with respect to the local data distribution D_k . The participating clients compute soft labels on the proxy data D_p and upload them to the server. The server then aggregates these labels

using a weighting scheme proportional to the size of each client local data to distill the global model. The corresponding objective function is formulated as:

$$w_k = \frac{N_k}{\|N\|_1}, \quad N = [N_1, N_2, \dots, N_K], \quad (2)$$

$$h(x_p) = \sigma(f(x_p; \theta)), \quad (3)$$

$$\theta_g^* = \arg \min_{\theta_g} \mathbb{E}_{x_p \sim D_p} [\mathcal{L}_{\text{KL}}(\sum_{k \in \mathcal{K}} w_k \cdot h_k(x_p), h_g(x_p))], \quad (4)$$

where $\sigma(\cdot)$ is the softmax function with temperature τ to control the smoothness of soft labels, and \mathcal{L}_{KL} is the Kullback-Leibler divergence used to measure the distribution differences between the server and the aggregated local models. Finally, each client further distills its local model using the KL divergence with the global soft labels predicted by the global model on the proxy data D_p .

$$\theta_k^* = \arg \min_{\theta_k} \mathbb{E}_{x_p \sim D_p} [\mathcal{L}_{\text{KL}}(h_g(x_p), h_k(x_p))]. \quad (5)$$

4 Method

We construct OOD benchmarks in natural language processing tasks to investigate the challenges of FD in real-world scenarios. Then, we propose **DPS-FD**, a novel framework of **Domain-aware Proxy Selection** for FD to address the OOD challenge by reconstructing the domain-aware proxy data. Furthermore, we introduce a vocabulary-constrained

LLM-based proxy generation strategy coupled with the DPS mechanism to mitigate the OOD problems of FD without proxy data.

4.1 OOD Benchmark Construction

To better understand the challenges of FD in real-world scenarios, we investigate the role of proxy data in OOD settings. Before introducing the OOD setting, we first consider a multi-domain non-IID configuration in natural language processing, which accounts not only for label distribution skew but also for the domain distribution of each client (Xiao and Liu, 2025; Yang et al., 2023; Xiao et al., 2024; Mao et al., 2025; Yan et al., 2025). Specifically, each client keeps private data from a distinct domain, and the label distribution across clients is made heterogeneous by sampling according to a Dirichlet Distribution with different α parameters. Building upon the multi-domain non-IID setup, we simulate OOD settings by replacing the original test set—composed of data from the client domains—with a test set sampled from the domains not seen by the clients. Both the private data and the proxy data are sampled from an open-source Amazon product review database.¹

4.2 Domain-Aware Proxy Selection

As shown in Figure 1, in communication round t , each client optimizes its local model using its labeled private data and utilizes the DPS mechanism to reconstruct domain-aware proxy data and uploads the corresponding local outputs to the server. The server aggregates the outputs and optimizes the global model through knowledge distillation. The global outputs given by the server are subsequently broadcast to all clients for local distillation, enabling bidirectional knowledge transfer.

DPS enables each client k to select the most relevant domain samples from the proxy data D_p , allowing clients to transmit high-confidence, domain-specific knowledge that enhances the robustness of the global model, while the server selects globally representative samples to improve the generalization of the global model. By taking the union of these selected samples, we construct a domain-aware proxy data D_p^* that serves as a more effective and efficient knowledge transfer medium between clients and the server.

In each communication round $t \in \{1, \dots, T\}$, after local training on its private data D_k , client k

computes the centroid representation c_k of its domain, and the server computes the global centroid c_g from the proxy data D_0 as:

$$\begin{aligned} c_k &= \frac{1}{N_k} \sum_{(x,y) \in D_k} f_k(x; \theta_k), \\ c_g &= \frac{1}{N_p} \sum_{x \in D_p} f_g(x; \theta_g), \end{aligned} \quad (6)$$

where $f_k(x)$ and $f_g(x)$ denote the input x feature representations of the local model and the global model, respectively. For each proxy sample $x_i \in D_0$, we calculate its cosine similarity with each client centroid and the global centroid:

$$\begin{aligned} \text{sim}_k(x_i) &= \frac{f_k(x_i; \theta_k) \cdot c_k}{\|f_k(x_i; \theta_k)\| \|c_k\|}, \\ \text{sim}_g(x_i) &= \frac{f_g(x_i; \theta_g) \cdot c_g}{\|f_g(x_i; \theta_g)\| \|c_g\|}. \end{aligned} \quad (7)$$

Each client and the server then select the top- K proxy samples most similar to its domain representations, respectively:

$$\begin{aligned} P_k &= \text{TopK}(\{x_i \in D_0 \mid \text{sim}_k(x_i)\}), \\ P_g &= \text{TopK}(\{x_i \in D_0 \mid \text{sim}_g(x_i)\}). \end{aligned} \quad (8)$$

Finally, the domain-aware proxy data used for knowledge distillation is constructed as the union:

$$D_p^* = \bigcup_{k=1}^{\mathcal{K}} P_k \cup P_g. \quad (9)$$

This selection process is repeated in each communication round, ensuring that the proxy data dynamically aligns with evolving domain representations as the training progresses.

4.3 LLM-Based Proxy Generation

The global model benefits from additional knowledge that helps align the client-specific distributions and facilitates more effective knowledge transfer. High-quality proxy data provides a bridge for the global model to capture patterns not fully represented in individual client private data, enhancing both robustness to domain-specific variations and generalization to unseen domains. In the absence of proxy data, the global model has to rely solely on local clients, which leads to weak alignment among local distributions and reduced overall performance.

To address the challenge of unavailable proxy data in FD, we propose a vocabulary-constrained

¹<https://archive.org/details/amazon-reviews-1995-2013>

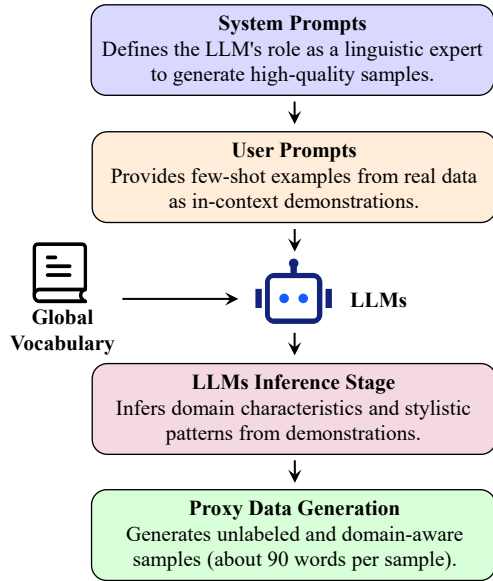


Figure 2: LLM-based Proxy Generation with Vocabulary-Constraints.

LLM-based proxy data generation strategy. As shown in Figure 2, we first construct a vocabulary from tokenizing real proxy data, which serves as a constraint to ensure that generated samples adhere to realistic lexical distributions and domain-specific semantics. To effectively generate diverse and high-quality proxy data, we adopt a few-shot prompting strategy that consists of two complementary prompt templates:

- **System Prompts** The system prompt defines the overall role and objective of the LLM, instructing it to act as a linguistic expert specialized in generating diverse, fluent, and contextually coherent product review texts tailored for FL tasks.
- **User Prompts** The user prompt provides few-shot representative samples drawn from real environments as in-context demonstrations.

To ensure linguistic consistency and data complexity, we constrain the average sentence length to approximately 90 words per generated text. As shown in Figure 2, two-stage prompting design allows the LLM to capture both global generation intent and domain-specific semantics.²

²The detailed structures of the system and user prompt are given in Appendix A.

Model	g.	a.	b.	c.	h.	s.	cost (%)
with proxy data							
FD-LP	80.7	79.8	84.7	91.4	77.4	81.8	100
FD-GP	81.4	80.3	84.1	85.0	75.6	78.7	100
DPS-FD	82.9	80.6	85.1	81.9	76.0	78.9	98.3
without proxy data							
FD-VC	78.2	75.0	80.5	73.9	68.4	72.7	100
DPS-FD	81.0	77.0	82.4	76.5	70.0	75.1	95.9
w/o VC	77.7	75.8	79.3	77.1	66.6	71.6	93.9

Table 1: Results of the global model in heterogeneous and OOD settings ($p < 0.001$). The column **g.** is score on test data in global distribution, while the columns **a.**, **b.**, **c.**, **h.** and **s.** are scores on test data in respective domains.

5 Experiments

We evaluate our proposed methods, DPS-FD along with a vocabulary-constrained LLM-based generation under OOD both with and without proxy data.

5.1 Settings

We conduct experiments on a sentiment classification task under heterogeneous client settings, considering both the presence and absence of proxy data. To simulate OOD scenarios, data are sampled by label within each domain from the Dirichlet distribution with concentration parameters $\alpha = \{0.3, 0.5, 0.8, 1.2, 1.5\}$ for clients, respectively, to control the degree of heterogeneity.

Backbones We assign different pre-trained language models to the 5 clients, namely BERT-base-cased, BERT-large-cased, RoBERTa-base, RoBERTa-large, and XLNet-large-cased model, respectively, while using RoBERTa-large as the global model for the server.³

Data Clients keep local data from specific domains, i.e., automotive (a.), baby (b.), clothing (c.), health (h.), and sport (s.), respectively, while the global distribution (g.) from all domains unseen by the clients.

5.2 Training and Inference

We train the models for 5 communication rounds with 3 local epochs per round, using an initial learning rate of $2e-5$. AdamW is adopted as the optimizer, with a maximum sequence length of 128 and

³All the pre-trained language model cards are sourced from <https://huggingface.co/models>.

Model	a.	b.	c.	h.	s.
with proxy data					
FD-LP	78.2	47.6	93.4	76.9	81.5
FD-GP	78.0	47.6	93.2	76.0	81.4
DPS-FD	76.7	83.5	89.8	77.4	75.0
without proxy data					
FD-VC	76.3	47.6	94.0	73.3	77.7
DPS-FD	78.5	84.6	90.4	76.5	82.3
w/o VC	75.2	82.1	91.3	70.0	72.8

Table 2: Results of local models on clients ($p < 0.001$). The columns **a.**, **b.**, **c.**, **h.** and **s.** are scores on test data in respective domains.

a batch size of 32. The number of clients is set to 5, corresponding to different numbers of domain-specific data. For the proposed DPS-FD method, we select the top 25k samples on each client and the top 35k samples on the server to form the refined proxy data. GPT-3.5-turbo is used as the LLM in the LLM-based generation strategy. We use an NVIDIA GPU with 24 GB memory for training and inference.

5.3 Baselines

We evaluate the proposed DPS-FD framework against a series of baselines to thoroughly investigate the impact of proxy data and the effectiveness of our approach under OOD settings.

- **FD-LP** refers to the standard FD method (Lin et al., 2020), which uses proxy data consisting only of samples from the client domains.
- **FD-GP** uses global proxy data composed of samples from all domains that are unseen by the clients.
- **FD-VC** leverages a vocabulary-constrained LLM-based generation strategy to synthesize proxy data without proxy data.
- **DPS-FD** is our proposed model, applying DPS on FD with global domain proxy data, with proxy data generated by a vocabulary-constrained LLM, and with proxy data generated by an unconstrained LLM.

5.4 Results

Table 1 and 2 shows the results of the global server model and the local client models, respectively, on our proposed model and several baseline models.

Proxy Data Influence As shown in Table 1, the global model of FD-GP outperforms that of FD-LP, indicating that global-domain proxy data, compared to local-domain proxy data, can more effectively mitigate distribution shifts under OOD settings.

- *With proxy data*, our proposed DPS-FD achieves superior performance across both global and local domains except for "clothing" domain. By incorporating the DPS, clients and the server dynamically identifies and utilizes proxy samples that better align with the distributions of local and global domains. DPS effectively filters out noisy or irrelevant samples, leading to high-quality proxy data that enable more effective and efficient knowledge transfer, demonstrating the effectiveness of DPS with proxy data.
- *Without proxy data*, DPS-FD outperforms DPS-FD without vocabulary-constraint by approximately 3.3%, demonstrating that vocabulary-constrained LLM generation produces higher-quality proxy data. We claim that high-quality proxy data can greatly improve the performance of the global model in FD. Moreover, DPS-FD substantially outperforms FD-VC, indicating that DPS remains effective even with synthetic proxy data.

Local Models Table 2 shows F1 scores of each local model evaluated on their own local test set. Existing methods perform poorly on local distribution of the "baby" domain, achieving only 47.6% F1 scores, whereas our proposed DPS-FD substantially achieves the improvements on the local test set. DPS effectively mitigates the interference from other domains and improve the performance of local models within their own domains by filtering out noise and irrelevant proxy samples. This demonstrates that DPS-FD not only enhances the global model on server but also somehow enhances the local model on clients.

Communication Cost To facilitate a fair comparison, we normalize the communication cost of the traditional FD method to 100% as a baseline. As shown in the column cost of Table 1, DPS-FD achieves a slight yet meaningful reduction in communication cost from 100% to 98.3% with proxy data and from 100% to 95.9% without proxy data. The overall cost is determined by the number of

Model	g.	a.	b.	c.	h.	s.	cost (%)
DS-FL	80.1	79.7	82.0	78.9	78.7	80.3	100
MHAT	81.6	79.5	83.6	80.3	72.6	78.6	100
FedKD	82.0	80.0	84.1	82.8	73.2	78.5	16.7
DPS-FD	82.9	80.6	85.1	81.9	76.0	78.9	98.3

Table 3: Results of the global model for existing FD methods under heterogeneous and OOD settings where proxy data is provided.

DPS-FD	top- K (client, server)	g.	cost (%)
5 Clients			
w/ proxy data	(15k, 15k)	82.6	89.8
	(25k, 35k)	82.9	98.3
w/o proxy data	(15k, 15k)	80.7	78.3
	(25k, 35k)	81.0	96.3
10 Clients			
w/ proxy data	(15k, 15k)	84.2	95.9
	(25k, 35k)	84.0	99.8
w/o proxy data	(15k, 15k)	81.4	80.3
	(25k, 35k)	82.7	93.6

Table 4: Results of DPS-FD across different number of clients and top- K selections.

clients \mathcal{K} , the proxy data size $|D_P|$, and the number of communication rounds T , and is typically expressed as $(\mathcal{K} + 1) \cdot |D_P| \cdot T$. By constructing a domain-aware proxy data $|D_P^*|$, DPS-FD improves model performance and reduces unnecessary communication by selecting representative samples that capture the distribution of each domain.

5.5 Comparisons with SOTA Models

As shown in Table 3, our proposed DPS-FD significantly outperforms existing FD methods, including DS-FL (Itahara et al., 2023), MHAT (Hu et al., 2021), and FedKD (Wu et al., 2022), in terms of global distribution performance with proxy data. This demonstrates not only superior robustness and generalization capabilities across heterogeneous client distributions, but also highlights the effectiveness of our domain-aware proxy selection in alleviating distribution shift and knowledge misalignment. Moreover, DPS-FD achieves this improved performance while reducing communication costs compared to DS-FL and MHAT, benefiting from a effective knowledge exchange process, except for FedKD, which adopts a one-shot distillation paradigm.

6 Discussion and Analysis

We take the analysis on the performance of DPS-FD in effect of the number of client, top- K selec-

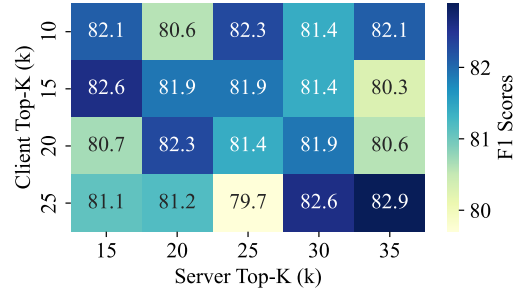


Figure 3: F1 scores varies with the top- K selections of the clients and server

tion, and quality of proxy data.

6.1 The Number of Clients

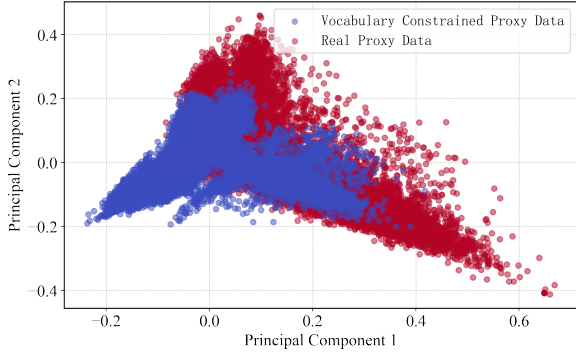
To investigate the effect of client number on training of the global model, 5 additional clients are introduced and assigned data of beauty, patio, pet, shoes, and software domains from the open-source Amazon product review database, using BERT-base-cased as the local model.

As shown in Table 4, both with and without proxy data, the global model performance on the global distribution improves significantly as the number of participating clients increases. The improvement can be attributed to the increased diversity of domain data introduced by increased clients. As more heterogeneous domains participate in training, the global model is exposed to a broader spectrum of feature distributions and semantic variations, enabling it to learn more domain-invariant representations. This diversity reduces sensitivity of the global model to domain-specific biases, enhancing its generalization and robustness under OOD distribution shifts.

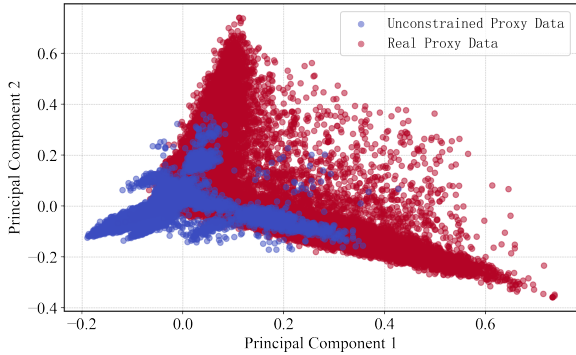
6.2 Top- K Selection

In DPS-FD, the choice of the top- K most similar samples selected by the client and the server plays a crucial role, as it directly affects the trade-off between model performance and communication cost. We vary the number of selected samples on the client side as {10k, 15k, 20k, 25k} and on the server side as {15k, 20k, 25k, 30k, 35k} under the setting with proxy data.

As shown in Figure 3, different top- K combinations lead to noticeable variations in the global model performance. The configurations (15k, 15k) and (25k, 30k) achieve strong suboptimal performance of 82.6%, indicating that comparable accuracy can be maintained with reduced commu-



(a) Vocabulary-constrained proxy data



(b) Unconstrained proxy data

Figure 4: Word distributions of proxy data generated under vocabulary-constrained (top) and unconstrained (bottom) settings, represented in the TF-IDF reduced dimensionality space.

520 nication cost. As shown in Table 4, with proxy
 521 data, the (15k, 15k) configuration yields a substan-
 522 tial reduction in communication cost while main-
 523 taining comparable global model performance. In
 524 particular, the communication overhead decreases
 525 by 9.5% when 5 clients are involved. Addition-
 526 ally, without proxy data, the same (15k, 15k) set-
 527 ting achieves an even greater reduction in com-
 528 munication cost by 18% and 13.3% respectively,
 529 demonstrating that moderate top- K selections on
 530 both client and server sides can effectively balance
 531 model performance and communication efficiency
 532 under different data availability.

533 6.3 Quality of Proxy Data

534 We take the insight analysis on vocabulary con-
 535 straints for the quality of proxy data generated by
 536 LLMs. We compare two generation strategies: with
 537 vocabulary constraints and without constraints. We
 538 tokenize the real data to construct a global vocabu-
 539 lary and enforce the LLM to generate samples
 540 using only the words in vocabulary. The uncon-
 541 strained LLM directly generates samples based on

542 same prompts without vocabulary constraints.

543 By constraining the vocabulary during LLM-
 544 based generation, we effectively intervene in the
 545 lexical choice process, guiding the generated sam-
 546 ples toward a distribution aligned with real-world
 547 data. This targeted intervention reduces spurious
 548 lexical variations and enhances the semantic rel-
 549 evance and consistency of the generated samples.
 550 Furthermore, as showed in Figure 4 (top), the dis-
 551 tribution of proxy data generated under vocabu-
 552 lary constraints is significantly closer to that of the real
 553 proxy data, demonstrating the effectiveness of our
 554 posed vocabulary constraints.

555 Additionally, the unconstrained generation strat-
 556 egy allows the LLM to produce samples solely
 557 based on the system and user prompts, without any
 558 explicit lexical guidance. As shown in Figure 4
 559 (bottom), the distribution of proxy data generated
 560 without constraints can be highly arbitrary, devi-
 561 ating substantially from the distribution observed
 562 in real proxy data. This randomness introduces
 563 spurious lexical variations and reduces semantic
 564 alignment, which can in turn limit the effectiveness
 565 of the generated samples for distillation. Uncon-
 566 strained generation produces proxy samples that
 567 are linguistically plausible but semantically irrele-
 568 vant or inconsistent, and applying vocabu-
 569 lary constraints during LLM-based proxy data gen-
 570 eration reduces the likelihood of such undesirable
 571 samples. Together with the results of Table 1, the im-
 572 provement observed under vocabulary constraints
 573 highlights the importance of controlling lexical dis-
 574 tributions to ensure that proxy data are semantically
 575 meaningful and relevant to global domain.

576 7 Conclusion

577 In this paper, we propose domain-aware proxy se-
 578 lection that constructs a representative proxy data
 579 by jointly considering client-specific and global
 580 domain samples to better adopt the proxy data for
 581 OOD problems in FD. Additionally, to address the
 582 absence of proxy data, we introduce a vocabu-
 583 lary-constrained LLM-based proxy data generation
 584 strategy, which mitigates the generation of linguisti-
 585 cally plausible but semantically irrelevant or inconsis-
 586 tent samples. By incorporating the proxy selection
 587 strategy with the generation strategy, DPS-FD en-
 588 hance the robustness and generalization of FD both
 589 with and without proxy data. Our experimental re-
 590 sults demonstrate that our models outperform exist-
 591 ing methods under OOD settings.

592 Limitations

593 Our method relies on LLMs to synthesize proxy
594 data. LLM-generated data are less reproducible
595 and incur higher API costs and longer response
596 time. We only simulate 5 and 10 clients in our
597 experiments, and we believe that using more clients
598 would be more effective. We take classification as
599 a representative example in this work, and need to
600 conduct more comprehensive validation.

601 References

602 Macarious Abadeer, Wei Shi, and Jean-Pierre Corriveau.
603 2022. Flightner: A federated learning approach to
604 lightweight named-entity recognition. In *2022 IEEE
605 International Conference on Trust, Security and Pri-
606 vacy in Computing and Communications (TrustCom)*,
607 pages 687–694. IEEE.

608 Rohan Anil, Gabriel Pereyra, Alexandre Passos, Robert
609 Ormandi, George E Dahl, and Geoffrey E Hinton.
610 2018. Large scale distributed neural network
611 training through online distillation. *arXiv preprint
612 arXiv:1804.03235*.

613 Haoyue Bai, Gregory Canal, Xuefeng Du, Jeongyeol
614 Kwon, Robert D Nowak, and Yixuan Li. 2023. Feed
615 two birds with one scone: Exploiting wild data for
616 both out-of-distribution generalization and detection.
617 In *International Conference on Machine Learning*,
618 pages 1454–1471. PMLR.

619 Yuhang Chen, Wenke Huang, and Mang Ye. 2024. Fair
620 federated learning under domain skew with local con-
621 sistency and domain diversity. In *Proceedings of
622 the IEEE/CVF Conference on Computer Vision and
623 Pattern Recognition (CVPR)*, pages 12077–12086.

624 Tao Fan, Hanlin Gu, Xuemei Cao, Chee Seng Chan,
625 Qian Chen, Yiqiang Chen, Yihui Feng, Yang Gu, Ji-
626 axiang Geng, Bing Luo, and 1 others. 2025. Ten
627 challenging problems in federated foundation mod-
628 els. *IEEE Transactions on Knowledge and Data
629 Engineering*.

630 Haoyang Fang, Boran Han, Shuai Zhang, Su Zhou,
631 Cuixiong Hu, and Wen-Ming Ye. 2024. Data aug-
632 mentation for object detection via controllable diffu-
633 sion models. In *Proceedings of the IEEE/CVF winter
634 conference on applications of computer vision*, pages
635 1257–1266.

636 Ishaan Gulrajani and David Lopez-Paz. 2020. In
637 search of lost domain generalization. *arXiv preprint
638 arXiv:2007.01434*.

639 Yaming Guo, Kai Guo, Xiaofeng Cao, Tieru Wu, and
640 Yi Chang. 2023. [Out-of-distribution generalization
641 of federated learning via implicit invariant relation-
642 ships](#). In *Proceedings of the 40th International
643 Conference on Machine Learning*, volume 202 of
644 *Proceedings of Machine Learning Research*, pages
645 11905–11933. PMLR.

Dan Hendrycks and Kevin Gimpel. 2016. A baseline
646 for detecting misclassified and out-of-distribution
647 examples in neural networks. *arXiv preprint
648 arXiv:1610.02136*. 649

Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. 650
651 Distilling the knowledge in a neural network. *arXiv
652 preprint arXiv:1503.02531*. 653

Li Hu, Hongyang Yan, Lang Li, Zijie Pan, Xiaozhang
654 Liu, and Zulong Zhang. 2021. Mhat: An efficient
655 model-heterogenous aggregation training scheme for
656 federated learning. *Information Sciences*, 560:493–
657 503. 658

Sohei Itahara, Takayuki Nishio, Yusuke Koda,
659 Masahiro Morikura, and Koji Yamamoto. 2023. 660
661 [Distillation-based semi-supervised federated learn-
662 ing for communication-efficient collaborative train-
663 ing with non-iid private data](#). *IEEE Transactions on
664 Mobile Computing*, 22(1):191–205. 665

Eunjeong Jeong, Seungeun Oh, Hyesung Kim, Jihong
666 Park, Mehdi Bennis, and Seong-Lyun Kim. 2023. 667
668 [Communication-efficient on-device machine learn-
669 ing: Federated distillation and augmentation under
670 non-iid private data](#). *Preprint*, arXiv:1811.11479. 671

Gu-Bon Jeong and Dong-Wan Choi. 2025. 672
673 Out-of-distribution detection via outlier exposure in feder-
674 ated learning. *Neural Networks*, 185:107141. 675

Daliang Li and Junpu Wang. 2019. Fedmd: Heteroge- 676
677 nous federated learning via model distillation. *arXiv
678 preprint arXiv:1910.03581*. 679

Zheng Li, Yuxuan Li, Penghai Zhao, Renjie Song, Xi- 680
681 ang Li, and Jian Yang. 2023. [Is synthetic data from
682 diffusion models ready for knowledge distillation?](#)
683 *Preprint*, arXiv:2305.12954. 684

Xinting Liao, Chaochao Chen, Weiming Liu, Pengyang
685 Zhou, Huabin Zhu, Shuheng Shen, Weiqiang Wang,
686 Mengling Hu, Yanchao Tan, and Xiaolin Zheng. 687
688 2023. Joint local relational augmentation and global
689 nash equilibrium for federated learning with non-iid
690 data. In *Proceedings of the 31st ACM International
691 Conference on Multimedia*, pages 1536–1545. 692

Xinting Liao, Weiming Liu, Chaochao Chen, Pengyang
693 Zhou, Fengyuan Yu, Huabin Zhu, Binhui Yao, Tao
694 Wang, Xiaolin Zheng, and Yanchao Tan. 2024a. 695
696 Re-thinking the representation in federated unsupervised
697 learning with non-iid data. In *Proceedings of the
698 IEEE/CVF Conference on Computer Vision and Pat-
699 tern Recognition*, pages 22841–22850. 700

Xinting Liao, Weiming Liu, Pengyang Zhou, Fengyuan
701 Yu, Jiahe Xu, Jun Wang, Wenjie Wang, Chaochao
702 Chen, and Xiaolin Zheng. 2024b. Foogd: Federated
703 collaboration for both out-of-distribution generaliza-
704 tion and detection. *Advances in Neural Information
705 Processing Systems*, 37:132908–132945. 706

699	Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. 2020. Ensemble distillation for robust model fusion in federated learning. <i>Advances in neural information processing systems</i> , 33:2351–2363.	<i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 20412–20421.	753 754 755
703	Tao Lin, Lingjing Kong, Sebastian U. Stich, and Martin Jaggi. 2021. Ensemble distillation for robust model fusion in federated learning. <i>Preprint</i> , arXiv:2006.07242.	Naibo Wang, Yuchen Deng, Wenjie Feng, Jianwei Yin, and See-Kiong Ng. 2024. Data-free federated class incremental learning with diffusion-based generative memory. <i>Preprint</i> , arXiv:2405.17457.	756 757 758 759
707	Xinge Ma, Jin Wang, and Xuejie Zhang. 2025. Data-free black-box federated learning via zeroth-order gradient estimation. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 19314–19322.	Zheng Wang, Zihui Wang, Zheng Wang, Xiaoliang Fan, and Cheng Wang. 2025a. Federated learning with domain shift eraser. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 4978–4987.	760 761 762 763 764
712	Wenjie Mao, Bin Yu, Chen Zhang, AK Qin, and Yu Xie. 2025. Fedkt: Federated learning with knowledge transfer for non-iid data. <i>Pattern Recognition</i> , 159:111143.	Zhongwei Wang, Tong Wu, Zhiyong Chen, Liang Qian, Yin Xu, and Meixia Tao. 2025b. Diffusion model-based data synthesis aided federated semi-supervised learning. <i>Preprint</i> , arXiv:2501.02219.	765 766 767 768
716	Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguerre y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In <i>Artificial intelligence and statistics</i> , pages 1273–1282. PMLR.	Zichen Wang, Feng Yan, Tianyi Wang, Cong Wang, Yuanhao Shu, Peng Cheng, and Jiming Chen. 2025c. Fed-dfa: Federated distillation for heterogeneous model fusion through the adversarial lens. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 21429–21437.	769 770 771 772 773 774
721	Jiaming Pei, Wenxuan Liu, Jinhai Li, Lukun Wang, and Chao Liu. 2024. A review of federated learning methods in heterogeneous scenarios. <i>IEEE Transactions on Consumer Electronics</i> , 70(3):5983–5999.	Chuhan Wu, Fangzhao Wu, Lingjuan Lyu, Yongfeng Huang, and Xing Xie. 2022. Communication-efficient federated learning via knowledge distillation. <i>Nature Communications</i> , 13(1).	775 776 777 778
725	Zhuang Qi, Weihao He, Xiangxu Meng, and Lei Meng. 2024. Attentive modeling and distillation for out-of-distribution generalization of federated learning. In <i>2024 IEEE International Conference on Multimedia and Expo (ICME)</i> , pages 1–6.	Canran Xiao and 1 others. 2024. Confusion-resistant federated learning via diffusion-based data harmonization on non-iid data. <i>Advances in Neural Information Processing Systems</i> , 37:137495–137520.	779 780 781 782
730	Zhuang Qi, Sijin Zhou, Lei Meng, Han Hu, Han Yu, and Xiangxu Meng. 2025. Federated deconfounding and debiasing learning for out-of-distribution generalization. <i>arXiv preprint arXiv:2505.04979</i> .	Jiahao Xiao and Jiangming Liu. 2025. Adaptive federated distillation for multi-domain non-iid textual data. <i>arXiv preprint arXiv:2508.20557</i> .	783 784 785
734	Jiawei Shao, Fangzhao Wu, and Jun Zhang. 2024. Selective knowledge sharing for privacy-preserving federated distillation without a good teacher. <i>Nature Communications</i> , 15(1):349.	Guochen Yan, Luyuan Xie, Xinyi Gao, Wentao Zhang, Qingni Shen, Yuejian Fang, and Zhonghai Wu. 2025. Fedvck: Non-iid robust and communication-efficient federated learning via valuable condensed knowledge for medical image analysis. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 21904–21912.	786 787 788 789 790 791 792
738	Dianbo Sui, Yubo Chen, Jun Zhao, Yantao Jia, Yuan-tao Xie, and Weijian Sun. 2020. FedED: Federated learning via ensemble distillation for medical relation extraction. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2118–2128, Online. Association for Computational Linguistics.	Lei Yang, Jiaming Huang, Wanyu Lin, and Jiannong Cao. 2023. Personalized federated learning on non-iid data via group-based meta-learning. <i>ACM Transactions on Knowledge Discovery from Data</i> , 17(4):1–20.	793 794 795 796 797
745	Hideaki Takahashi, Jingjing Liu, and Yang Liu. 2023. Breaching fedmd: image recovery via paired-logits inversion attack. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 12198–12207.	Mingzhao Yang, Shangchao Su, Bin Li, and Xiangyang Xue. 2024. Feddeo: Description-enhanced one-shot federated learning with diffusion models. <i>Preprint</i> , arXiv:2407.19953.	798 799 800 801
750	Haozhao Wang, Yichen Li, Wenchao Xu, Ruixuan Li, Yufeng Zhan, and Zhigang Zeng. 2023. Dafkd: Domain-aware federated knowledge distillation. In	Shan You, Chang Xu, Chao Xu, and Dacheng Tao. 2017. Learning from multiple teacher networks. In <i>Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining</i> , pages 1285–1294.	802 803 804 805 806

807 Shuyang Yu, Junyuan Hong, Haotao Wang, Zhangyang
808 Wang, and Jiayu Zhou. 2023. Turning the curse of
809 heterogeneity in federated learning into a blessing for
810 out-of-distribution detection. In *2023 International
811 Conference on Learning Representations*.

812 Lan Zhang, Dapeng Wu, and Xiaoyong Yuan. 2022.
813 Fedzkt: Zero-shot knowledge transfer towards
814 resource-constrained federated learning with hetero-
815 geneous on-device models. In *2022 IEEE 42nd In-
816 ternational Conference on Distributed Computing
817 Systems (ICDCS)*, pages 928–938. IEEE.

818 Tianxun Zhou and Keng-Hwee Chiam. 2023. [Synthetic
819 data generation method for data-free knowledge dis-
820 tillation in regression neural networks](#). *Expert Sys-
821 tems with Applications*, 227:120327.

822 Jun Zhu. 2024. [Synthetic data generation by diffusion
823 models](#). *National Science Review*, 11(8):nwae276.

824 Zhuangdi Zhu, Junyuan Hong, and Jiayu Zhou. 2021.
825 Data-free knowledge distillation for heterogeneous
826 federated learning. In *International conference on
827 machine learning*, pages 12878–12889. PMLR.

A Prompts

The following prompt templates are used to guide the LLM in generating proxy data.

User Prompt Template

A set of representative examples.

Your task:

1. Identify the product domains and sentiment tendencies represented in the examples above, and infer additional potential domains that could reasonably exist.
2. Based on the examples' and inferred domains, generate exactly realistic and diverse unlabeled review texts.
3. Each review should be fluent, reflecting the natural style of customer reviews.
4. Each review should be about 90 words (minimum 40, maximum 180). Do not shorten because you need to output many.
5. The reviews should not be templates or mechanically repeated. Maintain variability in tone, sentence structure, and vocabulary.

Output format:

1. Output exactly each reviews on a separate line.
2. Ensure the quality of each generated sentence and do not sacrifice quality for the sake of quantity.
3. Do not add any prefix, numbering, headers, or label—only the raw reviews.

Begin now:

System Prompt Template

You are an expert language model specializing in generating diverse, high-quality product review texts for federated learning.

You will carefully analyze a set of example reviews and identify corresponding domains and sentiment tendencies.

While the examples may be biased toward certain domains, your task is to summarize them and infer additional potential domains. Your generation should reflect the authentic style and tone of product reviews: natural, varied, and customer-oriented.

The output must ensure both domain diversity and linguistic diversity, avoiding repetitive templates while maintaining realism.

Algorithm 1: Domain-Aware Proxy Selection

Input: labeled private data $\{D_k\}_{k=1}^{\mathcal{K}}$;
unlabeled proxy data D_p ; global
model θ_g ; local models $\{\theta_k\}_{k=1}^{\mathcal{K}}$;
communication rounds T

Output: global model θ_g

if D_p **then**
 Use the real proxy data D_p ;
else
 Generate a synthetic proxy data D_p
 using LLM-based generation strategy ;
end

for *each communication round*
 $t = 1, 2, \dots, T$ **do**
 Client executes:
 for *each client k in parallel* **do**
 Train local model via **Eq. (1)** ;
 Select and upload the
 domain-specific proxy data via
 Eq. (6)(7)(8) ;
 end
 Server executes:
 Select the global domain proxy data via
 Eq. (6)(7)(8) ;
 Construct a domain-aware proxy data
 D_p^* via **Eq. (9)** ;
 Update the global model on D_p^* via
 Eq. (4) ;
 Client executes:
 for *each client k in parallel* **do**
 Update the local model on D_p^* via
 Eq. (5) ;
 end
end

B Algorithm of Domain-Aware Proxy Selection

Algorithm 1 summarizes the overall workflow of DPS-FD and details how the proposed domain-aware proxy selection is incorporated into the FD process.