SFEDPO: STREAMING FEDERATED LEARNING WITH A PREDICTION ORACLE UNDER TEMPORAL SHIFTS

Anonymous authorsPaper under double-blind review

000

001

002003004

010 011

012

013

014

016

017

018

019

021

024

025

026

027 028 029

031

033

034

035

037

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Federated Learning (FL) enables decentralized clients to collaboratively train a global model without sharing raw data. However, most existing FL frameworks assume that clients train on static local datasets collected in advance or that the data follows a fixed underlying distribution, which limits their applicability in dynamic environments where data evolves over time. A parallel line of research, online FL, removes all assumptions and adopts an adversarial perspective, but this approach is often overly pessimistic and neglects the structured, partially predictable nature of real-world data dynamics. To bridge this gap, we propose SFedPO, a streaming federated learning framework that incorporates a prediction oracle to capture the temporal evolution of client-side data distributions. We theoretically analyze the convergence bounds of SFedPO and develop two practical sampling strategies: a Distribution-guided Data Sampling (DDS) strategy that dynamically selects training data under limited storage by balancing historical reuse and distribution adaptation, and a Shift-aware Aggregation Weights (SAW) mechanism that modulates global aggregation based on client-specific sampling behaviors. We further establish robustness guarantees under prediction errors. Extensive experiments demonstrate that SFedPO effectively adapts to streaming scenarios with distribution shifts and significantly outperforms existing methods.

1 Introduction

Federated Learning (FL) (McMahan et al., 2017) has emerged as a prominent distributed learning paradigm that enables multiple clients to collaboratively train a shared global model while keeping data local. By leveraging local computation and periodically exchanging model parameters via a central server, FL eliminates the need for centralized data collection. This makes it particularly well-suited for real-world scenarios where data is inherently distributed, such as mobile sensing (Jiang et al., 2020) and edge intelligence (Mills et al., 2019).

Most existing FL frameworks assume that clients train on static local datasets collected in advance or that their data follows a fixed underlying distribution (McMahan et al., 2017; Li et al., 2020b; Wang et al., 2020b; Wang & Ji, 2022; Ye et al., 2023). In contrast, online FL (Mitra et al., 2021; M Ghari & Shen, 2022) discards such assumptions and instead adopts a potentially adversarial modeling perspective. It aims to minimize cumulative regret (Kwon et al., 2023; Patel et al., 2023) and provides robust theoretical guarantees under worst-case scenarios. While static local datasets or a fixed data distribution could be a strong assumption, the adversarial perspective represents the opposite extreme and is overly pessimistic.

Consider a streaming FL setting where clients (e.g., mobile devices or UAVs) continuously acquire new data over time. On one hand, this results in time-varying and non-stationary data distributions. On the other hand, the evolution of clients' data often follows structured and partially predictable patterns (Huynh et al., 2025). For instance, in the case of UAVs, their mobility trajectories typically follow pre-defined routes or scheduled missions. As a result, the data they collect exhibits spatial and temporal regularities. While assuming static datasets or fixed distributions could be inappropriate, the adversarial perspective may be overly pessimistic. This raises a natural question: Can we design a new FL framework that leverages partial predictions about data distribution shifts to guide both client-side sampling and server-side aggregation?

Addressing this question introduces significant challenges in both theoretical analysis and system design. Theoretically, the time-varying nature of clients' data renders the local objective functions time-dependent, complicating convergence analysis and requiring new tools to characterize the evolving optimization landscape. From a system design perspective, the key challenge lies in accommodating heterogeneous client behaviors, particularly their sampling strategies in response to distribution shifts.

Motivated by the discussions above, we propose SFedPO, a streaming federated learning framework with a prediction oracle to address temporal distribution shifts. Rather than assuming static datasets or adversarial dynamics, SFedPO operates in environments where clients continuously collect new data and incorporates a prediction oracle that provides prior knowledge about the temporal evolution of clients' data distributions. Building upon this, we introduce a distribution-guided data sampling strategy that selectively reuses and updates local data in response to distribution shifts. This allows clients to maintain a representative memory buffer under limited storage constraints. Furthermore, we develop an aggregation algorithm that adapts to heterogeneous client behaviors, including their sampling strategies. Our main contributions are summarized as follows:

- We propose SFedPO, a novel streaming FL framework incorporating a prediction oracle to model the temporal evolution of clients' data distributions. This framework provides a principled approach that bridges the gap between static-data assumptions and adversarial modeling commonly found in existing FL paradigms.
- We provide a theoretical convergence analysis of SFedPO and develop two core components: Distribution-guided Data Sampling (DDS) for local training, and Shift-aware Aggregation Weights (SAW) for client-adaptive global aggregation. In addition, we conduct a robustness analysis that quantifies the impact of oracle prediction errors on convergence guarantees.
- We perform extensive experiments, demonstrating that SFedPO effectively adapts to streaming scenarios with distribution shifts and significantly outperforms existing methods.

2 RELATED WORKS

We categorize existing FL literature into three main lines based on their assumptions about the data.

Federated Learning with Static or Stationary Data. Most traditional FL methods assume that clients train on static datasets or that their data follows a fixed distribution (McMahan et al., 2017; Li et al., 2020b; Wang et al., 2020b; Wang & Ji, 2022; Ye et al., 2023). Some recent studies extend FL to streaming scenarios. For instance, Marfoq et al. (2023) formalize FL over data streams and propose a meta-algorithm similar to vanilla FedAvg (McMahan et al., 2017) through a weighted empirical risk minimization design. ODE (Gong et al., 2023) introduces a date evaluation metric based on inference accuracy for on-device data selection under storage constraints. However, these works still assume that the data follows a stationary distribution.

Federated Learning under Distribution Shifts. Some studies have extended FL to streaming settings where data continuously arrive over time (Huynh et al., 2025; Marfoq et al., 2023; Gong et al., 2023; Liu et al., 2023). Wang et al. (2023a) assume the existence of a long-term distribution in the local data stream of clients and propose the cache update strategy to align the data distribution in the local cache to the underlying long-term distribution. Fed-HIST (Zhang et al., 2024) avoids the problem of raw data storage by retrieving model-based historical representations through similarity comparison. In parallel, the concept drift problem has been explored in FL. Most concept drift adaptation methods typically modify the model architecture (Chen et al., 2021), optimization strategy (Panchal et al., 2023; Canonaco et al., 2021), or client clustering (Li et al., 2024; Jothimurugesan et al., 2023; Chen et al., 2024) in response to detected shifts. Furthermore, Federated Continual Learning (FCL) (Yang et al., 2024; Guo et al., 2021; Dong et al., 2022) has been developed to mitigate catastrophic forgetting under sequential task arrivals, typically by leveraging parameter decomposition (Yoon et al., 2021), generative replay (Wuerkaixi et al., 2024; Qi et al., 2023), or knowledge distillation (Huang et al., 2022; Usmanova et al., 2021). However, existing approaches across these lines do not explicitly model or predict the temporal evolution of data distributions. In contrast, our work leverages a predictive oracle and develops a theoretically grounded, distribution-guided data sampling strategy for streaming FL. See details in Appendix B.

Online Federated Learning. Online FL removes assumptions on the underlying data distribution and instead aims to minimize regret under potentially arbitrary or adversarial data streams. FedOMD (Mitra et al., 2021) studies online federated optimization against adversarially revealed loss functions using online mirror descent, achieving sublinear regret. Ganguly & Aggarwal (2023) further combine FedAvg and FedOMD within a multiscale framework to adapt to non-stationary environments, establishing dynamic regret bounds under general convex losses. To address system-level challenges such as device heterogeneity and availability variations, ASO-Fed (Chen et al., 2020) introduces an asynchronous online FL framework based on continuous local updates and asynchronous aggregation.

3 Problem Formulation

3.1 Modeling Dynamic Data Distributions

We consider a federated learning system consisting of N clients, denoted by the set $\mathcal{N}=\{1,2,\ldots,N\}$, and these clients are coordinated by a central server. The learning process unfolds in R communication rounds, indexed by $r\in\{1,\ldots,R\}$. Each round is further divided into T fine-grained time steps, indexed by $t\in\{1,\ldots,T\}$. We define a time step as the granularity at which one client's data distribution may evolve.

Instead of assuming a static dataset or a stationary distribution, clients in our system continuously receive new data generated from a dynamic distribution. We model the distributions of new data as governed by a latent state space $\mathcal{M} = \{1, 2, \ldots, M\}$. Specifically, each state $m \in \mathcal{M}$ corresponds to a stationary data distribution \mathcal{D}_m . At each time step t in round r, client $n \in \mathcal{N}$ is associated with a latent state $m_{n,t}^{(r)} \in \mathcal{M}$. The ground-truth state distribution of client n is denoted as $\pi_n = (\pi_{n,1}, \ldots, \pi_{n,M})$, where $\pi_{n,m}$ denotes the probability that client n is in state m. While the true state transition dynamics are unknown, we assume the existence of a **prediction oracle** that provides a prediction over clients' latent states. That is, the oracle for client n outputs a prediction vector $\hat{\pi}_n = (\hat{\pi}_{n,1}, \ldots, \hat{\pi}_{n,M})$, where $\hat{\pi}_{n,m}$ denotes the estimate of $\pi_{n,m}$.

In streaming settings with limited local storage, clients must manage their training data continuously by retaining only a subset (possibly empty) of newly encountered samples. Considering this issue, we consider a distribution-guided data sampling mechanism, which explicitly adjusts the composition of each client's local dataset based on the observed data distribution. Specifically, we define a sampling strategy $\alpha = \{\alpha_{n,m}\}_{n \in \mathcal{N}, m \in \mathcal{M}}$, where $\alpha_{n,m} \in [0,1]$ represents the ratio of new samples from state m that client n incorporates into its local dataset. Let $\mathcal{D}_{n,t}^{(r)}$ denote the effective training distribution maintained by client n at time t in round t, which is updated as a convex combination of the previous local distribution and the current state's distribution, i.e.,

$$\mathcal{D}_{n,t}^{(r)} = (1 - \alpha_{n,m_{n,t}^{(r)}}) \cdot \mathcal{D}_{n,t-1}^{(r)} + \alpha_{n,m_{n,t}^{(r)}} \cdot \mathcal{D}_{m_{n,t}^{(r)}}. \tag{1}$$

This update captures the trade-off between preserving historical samples and adapting to recent distribution shifts. Given the local training distribution $\mathcal{D}_{n,t}^{(r)}$, we define the local loss function of client n at time step t in round r as

$$F_{n,t}^{(r)}(\mathbf{x}) \triangleq \mathbb{E}_{\xi \sim \mathcal{D}_{n,t}^{(r)}}[f(\mathbf{x};\xi)], \tag{2}$$

where $f(\mathbf{x}; \xi)$ denotes the sample-wise loss function.

3.2 FEDERATED TRAINING PROCESS

Our learning objective is to train a global model $\mathbf{x}^* \in \mathbb{R}^d$ that generalizes well across the full space of data distributions encountered by all clients over time. Specifically, we define the global objective as a weighted sum over all M possible distributions:

$$\mathbf{x}^* = \arg\min_{\mathbf{x}} F(\mathbf{x}) = \sum_{m=1}^{M} w_m \cdot F_m(\mathbf{x}), \tag{3}$$

where $F_m(\mathbf{x}) = \mathbb{E}_{\xi \sim \mathcal{D}_m}[f(\mathbf{x}; \xi)]$ denotes the expected loss under distribution \mathcal{D}_m , and w_m reflects the population-level importance or frequency of state m across all clients and time steps.

The training proceeds over R communication rounds. In each round r, a subset of clients S_r is sampled independently with availability probabilities q_n (Wang & Ji, 2022; Xiang et al., 2024), and each selected client initializes its local model with the current global model $\bar{\mathbf{x}}^{(r)}$. As defined earlier, each round is divided into T time steps, we align local model updates with these time steps: at each time step t in round r, client $n \in S_r$ samples a mini-batch of data from its updated local distribution $\mathcal{D}_{n,t}^{(r)}$ and performs one step of Stochastic Gradient Descent (SGD) as follows:

$$\mathbf{x}_{n,t}^{(r)} = \mathbf{x}_{n,t-1}^{(r)} - \eta \cdot \mathbf{g}_{n,t}^{(r)},\tag{4}$$

where η is the learning rate, and $\mathbf{g}_{n,t}^{(r)} \triangleq \nabla f(\mathbf{x}_{n,t-1}^{(r)}; \xi_{n,t})$ is the stochastic gradient of $F_{n,t}^{(r)}(\mathbf{x}_{n,t-1}^{(r)})$ evaluated on samples $\xi_{n,t} \sim \mathcal{D}_{n,t}^{(r)}$. After completing the local updates, each client uploads its final model $\mathbf{x}_{n,T}^{(r)}$ to the server. The server then aggregates the received models using a weighted average:

$$\bar{\mathbf{x}}^{(r+1)} = \sum_{n \in S_r} p_n \cdot \mathbf{x}_{n,E}^{(r)},\tag{5}$$

where p_n is the aggregation weight for the client n.

To support effective learning under temporal distribution shifts, our framework incorporates two core components: (1) a distribution-guided data sampling strategy that dynamically adjusts local datasets, and (2) an aggregation algorithm that accounts for the heterogeneous client behaviors, including sampling strategies and available probabilities.

4 THEORETICAL ANALYSIS

In this section, we analyze the convergence bound of the proposed federated learning framework under temporal distribution shifts, which highlights the effect of sampling strategy α and aggregation weights $\{p_n\}_{n\in\mathcal{N}}$. We begin by stating several assumptions on the local loss functions $F_{n,t}^{(r)}$ and the state-specific loss function F_m .

Assumption 1. (L-smoothness). All loss functions involved in the optimization are L-smooth. That is, there exists a constant L > 0 such that

$$\|\nabla F_{n,t}^{(r)}(\mathbf{x}) - \nabla F_{n,t}^{(r)}(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\| \text{ and } \|\nabla F_m(\mathbf{x}) - \nabla F_m(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|,$$

for all
$$\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$$
, $n \in \mathcal{N}$, $t \in \{1, \dots, T\}$, $r \in \{1, \dots, R\}$, and $m \in \mathcal{M}$.

Assumption 2. (Unbiased gradient and bounded variance). For each client n, the stochastic gradient is an unbiased estimator of the full gradient and has bounded variance. Formally,

$$\mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{D}_{n,t}^{(r)}}[\nabla f(\mathbf{x};\boldsymbol{\xi})] = \nabla F_{n,t}^{(r)}(\mathbf{x}), \ \mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{D}_{n,t}^{(r)}}[\|\nabla f(\mathbf{x};\boldsymbol{\xi}) - \nabla F_{n,t}^{(r)}(\mathbf{x})\|^2 |\mathbf{x}] \le \sigma^2,$$

where ξ is sampled from the n-th client's local data distribution $\mathcal{D}_{n,t}^{(r)}$ uniformly at random.

To characterize the heterogeneity of local objectives, many prior works adopt standard dissimilarity assumptions (Wang & Ji, 2022; Wang et al., 2020b), which bound the deviation between the local and global gradients:

$$\frac{1}{M} \sum_{m=1}^{M} \|\nabla F_m(\mathbf{x}) - \nabla F(\mathbf{x})\|^2 \le \beta^2 \|\nabla F(\mathbf{x})\|^2 + \zeta^2.$$

While this formulation captures the average heterogeneity across all latent states, it lacks granularity for state-specific characterization. Therefore, we decompose it into two assumptions: one that bounds the overall gradient heterogeneity across all latent states and another that controls the state-specific gradient heterogeneity. These can be viewed as variants of the above dissimilarity assumptions.

Assumption 3. There exists a constant G > 0 such that $\sum_{m=1}^{M} \|\nabla F_m(\mathbf{x})\|^2 \le G$, for all $\mathbf{x} \in \mathbb{R}^d$. **Assumption 4.** For each latent state $m \in \mathcal{M}$, there exists a constant $d_m > 0$ such that $\|\nabla F_m(\mathbf{x}) - \nabla F(\mathbf{x})\|^2 \le d_m$, for all $\mathbf{x} \in \mathbb{R}^d$.

We now present the main convergence result under the assumptions stated above. The theorem provides an upper bound on the optimization error, taking into account the effects of data dynamics, sampling strategies, and aggregation weights.

Figure 1: An overview of the proposed SFedPO framework. On the client side, each client determines its DDS strategy based on a prediction oracle $\hat{\pi}_n$ and state-wise heterogeneity bounds $\{d_m\}_{m=1}^M$. This produces sampling ratios $\hat{\alpha}_{n,m}$ for constructing the local dataset and updating the local model. On the server side, global aggregation is performed with the well-designed SAW, which adaptively adjusts the aggregation weights \hat{p}_n according to each client's heterogeneity score \hat{s}_n .

Theorem 1. Let Assumptions 1-4 hold. Suppose the data sampling strategy α and the aggregation weights $\{p_n\}_{n\in\mathcal{N}}$ are fixed. If the learning rate satisfies $LT\eta \leq \min\{\frac{1}{2M}, \sqrt{\frac{1}{5}}\}$. Then, the optimization error will be bounded as follows:

$$\min_{r} \mathbb{E} \|\nabla F(\bar{\mathbf{x}}^{(r)})\|^{2} \le \frac{18}{\sum_{n=1}^{N} p_{n} q_{n}} \left[\frac{1}{T \eta R} F(\bar{\mathbf{x}}^{(1)}) + L \eta \sigma^{2} \sum_{n=1}^{N} p_{n}^{2} q_{n} + \frac{5L \eta \sigma^{2}}{3} \sum_{n=1}^{N} p_{n} q_{n} + \frac{5L \eta \sigma^$$

$$\frac{5}{3}\sum_{n=1}^{N}p_{n}q_{n}\Big(\frac{\beta_{n}}{\alpha_{n}}(1-\gamma_{n})+\frac{2G\gamma_{n}\beta_{n}'}{1-(1-\alpha_{n})^{2}}+2G\gamma_{n}\sum_{m=1}^{M}(\frac{\pi_{n,m}\alpha_{n,m}}{\alpha_{n}}-w_{m})^{2}\Big)+\frac{1}{R}\sum_{n=1}^{N}\frac{7G\gamma_{n}q_{n}}{(1-(1-\alpha_{n})^{2T})}\Big],$$

where

$$\alpha_n = \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m}, \quad \beta_n = \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} d_m, \quad \gamma_n = \frac{1}{T} \sum_{t=1}^{T} (1 - \alpha_n)^t, \tag{7}$$

$$\beta_n' = 2\sum_{m=1}^M \pi_{n,m} \alpha_{n,m}^2 - \sum_{m=1}^M \pi_{n,m}^2 \alpha_{n,m}^2 + \alpha_n^2.$$
 (8)

Proof. See details in Appendix E.

5 THE SFEDPO FRAMEWORK

Based on the convergence analysis, we develop two coordinated modules that jointly minimize the upper bound of the optimization error in streaming federated environments.

Distribution-guided Data Sampling (DDS). To ensure effective local training under temporally evolving data, DDS leverages client-specific latent state distributions to guide sample selection. As demonstrated in our theoretical analysis, the convergence bound is strongly influenced by the data sampling strategy, particularly through quantities such as α_n , β_n , and β'_n . We treat α_n as a fixed hyperparameter that governs the average update ratio of the client buffer, referred to as the *sampling budget*. Accordingly, we formulate a constrained optimization problem to determine the optimal client-specific sampling ratios $\{\alpha_{n,m}\}_{m=1}^{M}$, leading to the solution (see Appendix F.1 for details):

$$\alpha_{n,m} \propto \frac{\lambda_1 \alpha_n + \frac{\nu_m - \mu_m}{\pi_{n,m}} + 4G\gamma_n w_m - (1 - \gamma_n) d_m}{1 + \frac{1 - \alpha_n}{\alpha_n} \pi_{n,m}},\tag{9}$$

where λ_1, μ_m and ν_m are multipliers. To simplify the solution and derive a closed-form, interpretable expression, we relax the KKT conditions by eliminating the dual variables, setting $\mu_m=0$ and $\nu_m=0$, under the assumption that the optimal values of $\alpha_{n,m}$ lie strictly within the open interval (0,1). This assumption is practically reasonable. In realistic streaming environments, it is neither desirable nor feasible for a client to completely discard previously stored data $(\alpha_{n,m}=1)$ or to entirely ignore new samples from a given state $(\alpha_{n,m}=0)$. Either extreme leads to inefficient storage

usage and poor adaptability to evolving data distributions. Under this relaxed setting, we define the following score function for each state $m \in \mathcal{M}$:

$$score_n(m) = \frac{w_m - a_1 d_m + b_1}{1 + \frac{1 - \alpha_n}{\alpha_n} \pi_{n,m}},$$
(10)

where $a_1 = \frac{1-\gamma_n}{4G\gamma_n}$ and $b_1 = \frac{\lambda_1}{4G\gamma_n}$ are tunable constants. Intuitively, the score increases with the weight in the global objective function w_m , and decreases with the state-specific heterogeneity bound d_m or the probability $\pi_{n,m}$ of client n being in state m. To ensure feasibility and numerical stability, we apply a ReLU operation to filter out negative scores, and normalize the resulting values with respect to the state distribution π_n :

$$\alpha_{n,m} = \frac{\alpha_n \cdot \text{ReLU}(\text{score}_n(m))}{\sum_{m'=1}^{M} \pi_{n,m'} \cdot \text{ReLU}(\text{score}_n(m'))}$$
(11)

To deal with the extreme case that $\alpha_{n,m} > 1$, we extend (11) and compute $\alpha_{n,m}$ using Algorithm 2.

Shift-aware Aggregation Weights (SAW). To enable global aggregation that adapts to heterogeneous client behaviors, we optimize the aggregation weights p_n via a relaxed surrogate objective (see Appendix F.2 for details).

$$\min_{p_n} L\eta \sigma^2 \sum_{n=1}^{N} p_n^2 q_n + \frac{5}{3} \sum_{n=1}^{N} p_n q_n s_n - \lambda_2 \sum_{n=1}^{N} p_n q_n$$
s.t.
$$\sum_{n=1}^{N} p_n = 1, \ p_n \ge 0,$$

where λ_2 is a balancing hyperparameter, and

$$s_n = L\eta\sigma^2 + \frac{\beta_n}{\alpha_n}(1 - \gamma_n) + \frac{2G\gamma_n\beta_n'}{1 - (1 - \alpha_n)^2} + 2G\gamma_n\sum_{m=1}^M \left(\frac{1}{\alpha_n}\pi_{n,m}\alpha_{n,m} - w_m\right)^2.$$
 (12)

We refer to s_n as the *heterogeneity score* of client n, as it captures the combined effect of its latent state distribution π_n and data sampling strategy $\{\alpha_{n,m}\}_{m=1}^M$. The closed-form solution is given by:

$$p_n \propto \frac{1}{q_n} - a_2 \cdot s_n + b_2, \tag{13}$$

where a_2 and b_2 are tunable constants. Intuitively, clients with smaller s_n are assigned higher aggregation weights. The term $1/q_n$ ensures fairness with respect to availability. We propose to determine more distinguishing aggregation weights for each client n by leveraging available probability q_n and heterogeneity score s_n as follows:

$$p_n = \frac{\text{ReLU}(\frac{1}{q_n} - a_2 \cdot s_n + b_2)}{\sum_{n'=1}^{N} \text{ReLU}(\frac{1}{q_{n'}} - a_2 \cdot s_{n'} + b_2)}.$$
 (14)

Practical Implementation with Prediction Oracle. The DDS and SAW modules are theoretically derived under the true latent state distribution π_n for each client n, which is not directly observable in practice. To make the framework practically implementable, we utilize a prediction oracle that provides an estimated distribution $\hat{\pi}_n$. Accordingly, the sampling strategy is adapted as follows:

$$\widehat{\operatorname{score}}_n(m) = \frac{w_m - a_1 d_m + b_1}{1 + \frac{\alpha_n}{1 - \alpha_n} \hat{\pi}_{n,m}}, \quad \hat{\alpha}_{n,m} = \frac{\operatorname{ReLU}(\widehat{\operatorname{score}}_n(m))}{\sum_{m'=1}^M \hat{\pi}_{n,m'} \cdot \operatorname{ReLU}(\widehat{\operatorname{score}}_n(m'))}.$$

Similarly, the aggregation weights are computed using $\hat{\pi}_n$ in place of π_n in all related terms.

$$\hat{p}_n = \frac{\text{ReLU}(\frac{1}{q_n} - a_2 \cdot \hat{s}_n + b_2)}{\sum_{n'=1}^N \text{ReLU}(\frac{1}{q_{n'}} - a_2 \cdot \hat{s}_{n'} + b_2)}.$$

The overall procedure of SFedPO is illustrated in Fig. 1, highlighting the key components on both the client and server sides. The complete procedure is summarized in Algorithm 1.

While the above implementation enables SFedPO to operate using a prediction oracle, it is critical to understand how prediction errors may affect the theoretical performance. To this end, we perform a robustness analysis that quantifies the impact of prediction errors on the convergence guarantee. Specifically, we define the prediction error for each client n as $\delta_n := \|\hat{\pi}_n - \pi_n\|_1$, and analyze how this error propagates into the convergence bound through the data sampling strategy α and aggregation weights $\{p_n\}_{n\in\mathcal{N}}$, both of which are functions of π_n in theory but implemented based on $\hat{\pi}_n$ in practice. Recall the convergence bound in Theorem 1, we focus on the following term:

$$\mathcal{B}(\alpha_{n,m}, p_n) := \frac{1}{\sum_{n=1}^{N} p_n q_n} \left[L \eta \sigma^2 \sum_{n=1}^{N} p_n^2 q_n + \frac{5}{3} \sum_{n=1}^{N} p_n q_n s_n \right].$$
 (15)

Theorem 2 (Robustness to prediction errors). Assume that all sampling scores and aggregation weights remain strictly positive under both the true and estimated state distributions. Let the data sampling and aggregation strategies be constructed using the estimated distributions $\hat{\pi}_n$ provided by a prediction oracle. Then, the convergence degradation compared to the ideal strategy using the true distributions π_n is bounded as

$$|\mathcal{B}(\hat{\alpha}_{n,m},\hat{p}_n) - \mathcal{B}(\alpha_{n,m},p_n)| \le \mathcal{O}\left(\sum_{n=1}^N \delta_n\right),\tag{16}$$

where $\delta_n := \|\hat{\pi}_n - \pi_n\|_1$ denotes the estimation error of the prediction oracle for client n. Here \mathcal{O} hides absolute constants.

Proof. See details in Appendix G.
$$\Box$$

6 EXPERIMENTS

6.1 EXPERIMENTAL SETUP

Datasets and models. We conduct comprehensive experiments on four public benchmark datasets: Fashion-MNIST (Xiao et al., 2017), CIFAR-10 (Krizhevsky et al., 2009), CINIC-10 (Darlow et al., 2018), and HAM10000 (Tschandl et al., 2018). We adopt LeNet-5 (LeCun et al., 1998) for Fashion-MNIST, AlexNet (Krizhevsky et al., 2012) for CIFAR-10 and CINIC-10, and a customized CNN for HAM10000. Details of datasets and models are provided in Appendix H.1.

Federated settings. We simulate a federated environment consisting of N=30 clients. In each communication round, clients are selected independently based on their availability probability $\{q_n\}_{n=1}^N$. For each client n, we draw $q_n \sim \mathcal{N}(0.2, 0.01^2)$, clipped to [0.01, 1] to avoid degenerate participation. To capture temporally evolving local distributions, we introduce M=60 latent states, each representing a possible data distribution encountered over time. To simulate intra-state heterogeneity, we organize these states into 6 clusters (10 states per cluster), each associated with a Dirichlet partitioning strategy with distinct concentration parameters $\alpha \in \{0.05, 0.1, 0.2, 0.5, 1.0, 100.0\}$. Based on this latent state space, we consider two heterogeneity scenarios. (1) Full-access (mild heterogeneity): Each client has non-zero access probability $\pi_{n,m} > 0$ for all states $m \in \{1, \dots, M\}$. (2) Partial-access (extreme heterogeneity): Each client is restricted to a randomly sampled subset of 10 states, with $\pi_{n,m} = 0$ for others. Moreover, 50% of clients are initialized with latent states drawn from high-heterogeneity clusters ($\alpha \in \{0.05, 0.1\}$), thereby amplifying both spatial and temporal heterogeneity. These two scenarios allow us to rigorously test the robustness of SFedPO under both mild and extreme data heterogeneity. Details are provided in Appendix H.1.

Baselines. We compare SFedPO against representative methods from four categories: (1) Online FL: FedOGD (Kwon et al., 2023) and FedOMD (Mitra et al., 2021); (2) FL for concept drift: AdapFedAvg (Canonaco et al., 2021), FedDrift (Jothimurugesan et al., 2023), and Flash (Panchal et al., 2023); (3) FCL: FedEWC (FedOGD with EWC (Kirkpatrick et al., 2017) applied to clients) and FLwF-2T (Usmanova et al., 2021); (4) Data selection methods in FL: Importance Sampling (IS) (Li et al., 2021), ODE (Gong et al., 2023), and DRSR (Wang et al., 2023a).

6.2 MAIN RESULTS

Performance: test accuracy. Table 1 presents the test accuracy of SFedPO and all baselines on four datasets under two scenarios. The results demonstrate that SFedPO consistently outperforms

Table 1: Test accuracy (%, mean±std on 5 trials) comparison of our SFedPO framework to other baselines in full and partial scenarios on several datasets.

Method	Fashion-MNIST		CIFAR-10		CINIC-10		HAM10000	
	full	partial	full	partial	full	partial	full	partial
FedOGD	85.04±0.46	82.67±0.68	51.14±2.23	41.06±2.09	43.86±0.90	34.40±1.91	54.11±0.51	45.28±0.62
FedOMD	85.03±0.41	82.67±0.70	51.06±2.22	40.98±2.08	43.84±0.91	34.37±1.91	54.14±0.55	45.21±0.67
AdapFedAvg	84.95±0.42	82.61±0.74	48.86±2.33	39.37±2.25	42.94±0.99	33.22±1.96	52.24±0.53	41.65±1.77
FedDrift	85.06±0.40	82.62±0.72	51.07±2.40	41.23±2.14	43.88±0.84	34.57±1.73	54.25±0.43	45.35±0.67
Flash	85.65±0.83	83.89±1.05	63.87±3.44	56.08±0.93	50.80±0.92	44.34±2.72	64.70±0.74	54.70±1.99
FedEWC	85.12±0.46	82.77±0.69	51.18±2.23	41.12±2.06	43.88±0.92	34.46±1.92	54.38±0.53	45.43±0.75
FLwF-2T	86.46±0.10	84.92±0.64	49.69±0.83	40.01±1.09	43.65±0.56	36.51±0.69	56.88±0.40	49.28±1.27
IS	83.36±0.73	80.27±1.33	52.03±1.87	42.66±2.02	37.92±0.82	30.86±1.39	55.32±1.99	42.99±3.18
ODE	83.33±0.42	80.21±0.79	53.52±1.71	40.69±1.57	38.02±0.90	31.31±1.28	55.32±0.89	43.31±2.64
DRSR	85.75±0.21	83.04±0.57	57.64±1.11	48.42±0.73	43.49±0.59	37.90±1.00	59.72±0.55	52.42±1.47
SFedPO	87.60±0.06	86.77±0.42	67.45±0.16	63.53±0.80	51.00±0.32	47.71±0.59	68.63±0.43	64.23±0.31

Table 2: Modularity. Accuracy (%) of classic federated learning methods with SFedPO and their improvement over the originals without SFedPO.

Method	Fashion	-MNIST	CIFA	AR-10	CINI	C-10	HAM	HAM10000	
	full	partial	full	partial	full	partial	full	partial	
FedAvg FedProx FedCurv FedNTD FedEXP	87.55(+0.29) 87.23(+0.34) 87.62(+0.42) 87.50(+0.09) 85.07(+1.35)	86.63(+0.13) 86.29(+0.22) 86.76(+0.25) 87.09(+0.40) 85.50(+4.17)	67.72(+1.86) 65.21(+1.88) 67.60(+1.75) 65.63(+2.09) 65.71(+1.93)	63.83(+4.01) 61.69(+4.58) 63.48(+3.65) 62.43(+3.05) 66.57(+7.76)	50.90(+1.40) 50.08(+1.66) 50.97(+1.49) 52.61(+1.08) 43.65(+1.14)	47.90(+1.18) 46.73(+0.96) 48.04(+1.20) 50.45(+0.80) 48.48(+7.82)	68.88(+1.61) 67.14(+1.80) 68.94(+1.19) 68.56(+1.20) 65.72(+2.14)	64.53(+1.52) 62.68(+1.57) 65.03(+2.37) 65.82(+0.79) 64.71(+6.72)	

all baselines across different settings. For instance, on the CIFAR-10 dataset, SFedPO surpasses all other methods by at least 3.58% in the full-access scenario and 7.39% in the partial-access scenario. We further observe that the performance gain of SFedPO in the partial-access scenario is more pronounced than in the full-access scenario, indicating that our method is particularly effective under extreme heterogeneity. This validates the effectiveness of our distribution-guided data sampling and aggregation strategies, which adaptively respond to state-specific and client-specific variation.

Modularity: improvements over FL methods. Our proposed SFedPO exhibits strong modularity and can be easily integrated into a wide range of classical FL methods as a plug-and-play module to cope with streaming data scenarios. To evaluate its effectiveness in this setting, we apply SFedPO's data sampling and aggregation strategies to several representative FL methods, including FedAvg (McMahan et al., 2017), FedProx (Li et al., 2020a), FedCurv (Shoham et al., 2019), FedNTD (Lee et al., 2022), and FedEXP (Jhunjhunwala et al., 2023). For comparison, we consider the original versions of these methods under same streaming settings, where each client performs uniform data sampling across latent states (i.e., $\alpha_{n,m} = \alpha_n$ for all m), and the server performs uniform model averaging over the participating clients (i.e., $p_n = \frac{1}{|S_r|}$ for each $n \in S_r$). As shown in Table 2, integrating SFedPO consistently improves the test accuracy of all methods across different datasets and both full- and partial-access scenarios. The gains are particularly notable under partial access, where client heterogeneity is more severe, confirming that SFedPO enhances the robustness and adaptability of FL methods under dynamic data distributions.

6.3 ABLATION STUDY

Effects of different configurations. We first evaluate the robustness of SFedPO to different configuration parameters in the FL environment. Specifically, we vary three core parameters, including time step $(T \in \{2, 5, 8, 10\})$, training round $(R \in \{50, 100, 150, 200\})$, and data capacity of clients $(D \in \{250, 500, 750, 1000\})$, then respectively show the performance of our method and four baselines in Fig. 2a, Fig. 2b, and Fig. 2c. The experimental results demonstrate that our method outperforms all baseline approaches across different parameters. All the experiments in this part are conducted in the partial scenario on the CIFAR-10 dataset.

Effectiveness and robustness of modules. We evaluate the stability and effectiveness of the two core components in SFedPO: Distribution-guided Data Sampling (DDS) and Shift-aware Aggregation Weights (SAW). For the DDS module, we first vary the sampling budget ($\alpha_n \in$

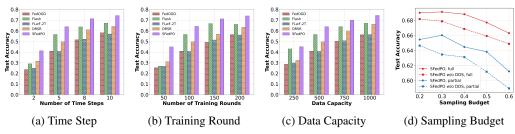


Figure 2: The impact of four key parameters on performance.

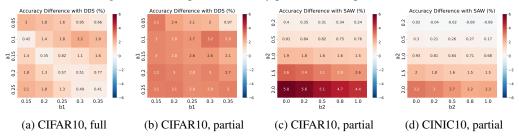


Figure 3: Hyperparameter sensitivity analysis for $\{a_1, b_1\}$ and $\{a_2, b_2\}$.

 $\{0.2, 0.3, 0.4, 0.5, 0.6\}$) and compare the performance of the standard SFedPO with an ablated version where clients sample data uniformly across latent states (i.e., $\alpha_{n,m}=\alpha_n$ for all m). Fig.2d demonstrates that the DDS module consistently improves training performance across all sampling budgets. We then tune the hyperparameters of $\{a_1,b_1\}$ in (86) across diverse datasets and scenarios. Fig. 3a and Fig. 3b illustrate that DDS consistently improves performance across a wide range of hyperparameter settings $(a_1 \in [0.05, 0.25], b_1 \in [0.15, 0.35])$. For the SAW module, we investigate the impact of parameters a_2, b_2 in (13). As shown in Fig. 3c and Fig. 3d, SAW brings noticeable accuracy gains over uniform aggregation $(p_n = 1/|S_r|)$, with stable performance across a wide range of parameter values $(a_2 \in [0.5, 2.0], b_2 \in [0.0, 1.0])$. These results collectively confirm that both DDS and SAW are not only effective but also resilient to hyperparameter changes.

Effects of Prediction Error. To evaluate the robustness of SFedPO against errors in the prediction oracle, we conduct an empirical study aligned with our theoretical analysis in Theorem 2. To simulate prediction errors, we introduce additive perturbations to the ground-truth π_n to obtain noisy estimates $\hat{\pi}_n$, and apply SFedPO based on $\hat{\pi}_n$ to perform local data sampling and global aggregation. Specifically, we perturb each state probability by a random noise uniformly drawn from $[-\epsilon, \epsilon]$, followed by renormalization to ensure $\sum_m \hat{\pi}_{n,m} = 1$. We vary ϵ from 0.00 to 0.10 to simulate increasing levels of oracle error, and report the resulting model accuracy in Table 3. We observe that SFedPO exhibits stable performance under varying degrees of perturbation on the CIFAR-10 dataset.

Table 3: Accuracy (%, mean±std on 5 trials) under different degrees of perturbation.

Epsilon	0.00	0.02	0.04	0.06	0.08	0.10
full	67.45 _{±0.16}	$66.90_{\pm0.56}$	$66.96_{\pm0.45}$	$66.97_{\pm0.31}$	$67.14_{\pm0.47}$	$66.95_{\pm0.32}$
partial	63.52 _{±0.82}	$62.63_{\pm0.16}$	$62.36_{\pm0.20}$	$62.49_{\pm0.14}$	$62.48_{\pm0.15}$	$61.90_{\pm0.15}$

7 CONCLUSIONS

We propose SFedPO, a streaming federated learning framework for dynamic environments with evolving local data distributions. Departing from conventional FL assumptions of static datasets or a stationary distribution, SFedPO incorporates a prediction oracle to capture the temporal evolution of client-side data distributions. Guided by theoretical convergence analysis, we develop two key components: a Distribution-guided Data Sampling (DDS) strategy that balances data reuse and distribution adaptation under storage constraints, and a Shift-aware Aggregation Weights (SAW) mechanism that adjusts global aggregation in response to client-specific sampling behavior. We further establish robustness guarantees under prediction errors. Extensive experiments demonstrate that SFedPO effectively adapts to streaming scenarios with distribution shifts and significantly outperforms existing methods.

REFERENCES

- Giuseppe Canonaco, Alex Bergamasco, Alessio Mongelluzzo, and Manuel Roveri. Adaptive federated learning in presence of concept drift. In *International Joint Conference on Neural Networks* (*IJCNN*). IEEE, 2021.
- Junbao Chen, Jingfeng Xue, Yong Wang, Zhenyan Liu, and Lu Huang. Classifier clustering and feature alignment for federated learning under distributed concept drift. In *Advances in Neural Information Processing Systems*, 2024.
- Yujing Chen, Yue Ning, Martin Slawski, and Huzefa Rangwala. Asynchronous online federated learning for edge devices with non-iid data. In *IEEE International Conference on Big Data (Big Data)*. IEEE, 2020.
- Yujing Chen, Zheng Chai, Yue Cheng, and Huzefa Rangwala. Asynchronous federated learning for sensor data with concept drift. In *IEEE International Conference on Big Data (Big Data)*, pp. 4822–4831. IEEE, 2021.
- Luke N Darlow, Elliot J Crowley, Antreas Antoniou, and Amos J Storkey. Cinic-10 is not imagenet or cifar-10. arXiv preprint arXiv:1810.03505, 2018.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE conference on computer vision and pattern recognition*, pp. 248–255. IEEE, 2009.
- Jiahua Dong, Lixu Wang, Zhen Fang, Gan Sun, Shichao Xu, Xiao Wang, and Qi Zhu. Federated class-incremental learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10164–10173, 2022.
- Bent Fuglede and Flemming Topsoe. Jensen-shannon divergence and hilbert space embedding. In *International symposium on Information theory (ISIT).*, pp. 31. IEEE, 2004.
- João Gama, Indré Žliobaitė, Albert Bifet, Mykola Pechenizkiy, and Abdelhamid Bouchachia. A survey on concept drift adaptation. *ACM computing surveys (CSUR)*, 46(4):1–37, 2014.
- Bhargav Ganguly and Vaneet Aggarwal. Online federated learning via non-stationary detection and adaptation amidst concept drift. *IEEE/ACM Transactions on Networking*, 32(1):643–653, 2023.
- Chen Gong, Zhenzhe Zheng, Fan Wu, Yunfeng Shao, Bingshuai Li, and Guihai Chen. To store or not? online data selection for federated learning with limited storage. In *Proceedings of the ACM Web Conference*, pp. 3044–3055, 2023.
- Yongxin Guo, Tao Lin, and Xiaoying Tang. Towards federated learning on time-evolving heterogeneous data. *arXiv preprint arXiv:2112.13246*, 2021.
- Yongxin Guo, Xiaoying Tang, and Tao Lin. Fedrc: Tackling diverse distribution shifts challenge in federated learning by robust clustering. In *International Conference on Machine Learning*. PMLR, 2024.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Chung-Hsuan Hu, Zheng Chen, and Erik G Larsson. Energy-efficient federated edge learning with streaming data: A lyapunov optimization approach. *IEEE Transactions on Communications*, 2024.
- Wenke Huang, Mang Ye, and Bo Du. Learn from others and be yourself in heterogeneous federated learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10143–10153, 2022.
- Tan-Khiem Huynh, Malcolm Egan, Giovanni Neglia, and Jean-Marie Gorce. Streaming federated learning with markovian data. *arXiv preprint arXiv:2503.18807*, 2025.
- Divyansh Jhunjhunwala, Shiqiang Wang, and Gauri Joshi. Fedexp: Speeding up federated averaging via extrapolation. *arXiv preprint arXiv:2301.09604*, 2023.

- Ji Chu Jiang, Burak Kantarci, Sema Oktug, and Tolga Soyata. Federated learning in smart city sensing: Challenges and opportunities. *Sensors*, 20(21):6230, 2020.
 - Yibo Jin, Lei Jiao, Zhuzhong Qian, Sheng Zhang, and Sanglu Lu. Budget-aware online control of edge federated learning on streaming data with stochastic inputs. *IEEE Journal on Selected Areas in Communications*, 39(12):3704–3722, 2021.
 - Ellango Jothimurugesan, Kevin Hsieh, Jianyu Wang, Gauri Joshi, and Phillip B Gibbons. Federated learning under distributed concept drift. In *International Conference on Artificial Intelligence and Statistics*, pp. 5834–5853. PMLR, 2023.
 - James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114 (13):3521–3526, 2017.
 - Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
 - Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
 - Dohyeok Kwon, Jonghwan Park, and Songnam Hong. Tighter regret analysis and optimization of online federated learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45 (12):15772–15789, 2023.
 - Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
 - Gihun Lee, Minchan Jeong, Yongjin Shin, Sangmin Bae, and Se-Young Yun. Preservation of the global knowledge by not-true distillation in federated learning. *Advances in Neural Information Processing Systems*, 35:38461–38474, 2022.
 - Anran Li, Lan Zhang, Juntao Tan, Yaxuan Qin, Junhao Wang, and Xiang-Yang Li. Sample-level data selection for federated learning. In *IEEE INFOCOM 2021-IEEE Conference on Computer Communications*, pp. 1–10. IEEE, 2021.
 - Minghao Li, Dmitrii Avdiukhin, Rana Shahout, Nikita Ivkin, Vladimir Braverman, and Minlan Yu. Federated learning clients clustering with adaptation to data drifts. *arXiv preprint arXiv:2411.01580*, 2024.
 - Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*, 2:429–450, 2020a.
 - Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. In *International Conference on Learning Representations*, 2020b.
 - Weijie Liu, Xiaoxi Zhang, Jingpu Duan, Carlee Joe-Wong, Zhi Zhou, and Xu Chen. Dynamite: Dynamic interplay of mini-batch size and aggregation frequency for federated learning with static and streaming datasets. *IEEE Transactions on Mobile Computing*, 23(7):7664–7679, 2023.
 - Jie Lu, Anjin Liu, Fan Dong, Feng Gu, Joao Gama, and Guangquan Zhang. Learning under concept drift: A review. *IEEE transactions on knowledge and data engineering*, 31(12):2346–2363, 2018.
 - Pouya M Ghari and Yanning Shen. Personalized online federated learning with multiple kernels. *Advances in Neural Information Processing Systems*, 35:33316–33329, 2022.
 - Dimitrios Michael Manias, Ibrahim Shaer, Li Yang, and Abdallah Shami. Concept drift detection in federated networked systems. In *IEEE global communications conference (GLOBECOM)*, pp. 1–6. IEEE, 2021.
 - Othmane Marfoq, Giovanni Neglia, Laetitia Kameni, and Richard Vidal. Federated learning for data streams. In *International Conference on Artificial Intelligence and Statistics*, pp. 8889–8924. PMLR, 2023.

- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pp. 1273–1282. PMLR, 2017.
 - Jed Mills, Jia Hu, and Geyong Min. Communication-efficient federated learning for wireless edge intelligence in iot. *IEEE Internet of Things Journal*, 7(7):5986–5994, 2019.
 - Aritra Mitra, Hamed Hassani, and George J Pappas. Online federated learning. In 2021 60th IEEE Conference on Decision and Control (CDC), pp. 4083–4090. IEEE, 2021.
 - Kunjal Panchal, Sunav Choudhary, Subrata Mitra, Koyel Mukherjee, Somdeb Sarkhel, Saayan Mitra, and Hui Guan. Flash: Concept drift adaptation in federated learning. In *International Conference on Machine Learning*, pp. 26931–26962. PMLR, 2023.
 - Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
 - Kumar Kshitij Patel, Lingxiao Wang, Aadirupa Saha, and Nathan Srebro. Federated online and bandit convex optimization. In *International Conference on Machine Learning*, pp. 27439–27460. PMLR, 2023.
 - Daiqing Qi, Handong Zhao, and Sheng Li. Better generative replay for continual federated learning. In *The Eleventh International Conference on Learning Representations*, 2023.
 - Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. In *International Conference on Learning Representations*, 2020.
 - Neta Shoham, Tomer Avidor, Aviv Keren, Nadav Israel, Daniel Benditkis, Liron Mor-Yosef, and Itai Zeitak. Overcoming forgetting in federated learning on non-iid data. *arXiv preprint arXiv:1910.07796*, 2019.
 - Ashraf Tahmasbi, Ellango Jothimurugesan, Srikanta Tirthapura, and Phillip B Gibbons. Driftsurf: Stable-state/reactive-state learning under concept drift. In *International Conference on Machine Learning*, pp. 10054–10064. PMLR, 2021.
 - Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific data*, 5(1):1–9, 2018.
 - Anastasiia Usmanova, François Portet, Philippe Lalanda, and German Vega. A distillation-based approach integrating continual learning and federated learning for pervasive services. *arXiv* preprint arXiv:2109.04197, 2021.
 - Heqiang Wang, Jieming Bian, and Jie Xu. On the local cache update rules in streaming federated learning. *IEEE Internet of Things Journal*, 11(6):10808–10816, 2023a.
 - Hongyi Wang, Mikhail Yurochkin, Yuekai Sun, Dimitris S. Papailiopoulos, and Yasaman Khazaeni. Federated learning with matched averaging. In *International Conference on Learning Representations*, 2020a.
 - Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. *Advances in neural information processing systems*, 33:7611–7623, 2020b.
 - Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual learning: Theory, method and application. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Naiyu Wang, Xuan Li, Zhitao Guan, and Shuai Yuan. Fedstream: A federated learning framework on heterogeneous streaming data for next-generation traffic analysis. *IEEE Transactions on Network Science and Engineering*, 11(3):2485–2496, 2023b.

- Shiqiang Wang and Mingyue Ji. A unified analysis of federated learning with arbitrary client participation. *Advances in Neural Information Processing Systems*, 35:19124–19137, 2022.
- Abudukelimu Wuerkaixi, Sen Cui, Jingfeng Zhang, Kunda Yan, Bo Han, Gang Niu, Lei Fang, Changshui Zhang, and Masashi Sugiyama. Accurate forgetting for heterogeneous federated continual learning. In *The Twelfth International Conference on Learning Representations*, 2024.
- Ming Xiang, Stratis Ioannidis, Edmund Yeh, Carlee Joe-Wong, and Lili Su. Efficient federated learning against heterogeneous and non-stationary client unavailability. *Advances in Neural Information Processing Systems*, 37:104281–104328, 2024.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv* preprint arXiv:1708.07747, 2017.
- Xin Yang, Hao Yu, Xin Gao, Hao Wang, Junbo Zhang, and Tianrui Li. Federated continual learning via knowledge fusion: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 36(8): 3832–3850, 2024.
- Rui Ye, Mingkai Xu, Jianyu Wang, Chenxin Xu, Siheng Chen, and Yanfeng Wang. Feddisco: Federated learning with discrepancy-aware collaboration. In *International Conference on Machine Learning*, 2023.
- Jaehong Yoon, Wonyong Jeong, Giwoong Lee, Eunho Yang, and Sung Ju Hwang. Federated continual learning with weighted inter-client transfer. In *International Conference on Machine Learning*, pp. 12073–12086. PMLR, 2021.
- Jianyi Zhang, Ang Li, Minxue Tang, Jingwei Sun, Xiang Chen, Fan Zhang, Changyou Chen, Yiran Chen, and Hai Li. Fed-CBS: A heterogeneity-aware client sampling mechanism for federated learning via class-imbalance reduction. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 41354–41381. PMLR, 2023.
- Ruirui Zhang, Yifei Zou, Zhenzhen Xie, Xiao Zhang, Peng Li, Zhipeng Cai, Xiuzhen Cheng, and Dongxiao Yu. Federating from history in streaming federated learning. In *Proceedings of the Twenty-fifth International Symposium on Theory, Algorithmic Foundations, and Protocol Design for Mobile Networks and Mobile Computing*, pp. 151–160, 2024.

Appendix

A The Use of LLMs Related Works Discussions **D** Notations and Technical Lemmas E Proof of Theorem 1 **Upper Bound Minimization** G Proof of Theorem 2 **H** More Experimental Details

A THE USE OF LLMS

756

757 758

759

760

761 762

763 764

765

766

767 768

769

770

771

772

773

774

775

776

777

778

779

781

782

783

784

785

786 787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806 807

808

This paper was proofread and linguistically polished with the assistance of a large language model (ChatGPT). All research ideas, methods, analyses, experiments, and conclusions are entirely the work of the authors.

B RELATED WORKS

This section provides a more detailed introduction to the second type of data stream paradigm, i.e., FL under distribution shifts, mainly from the perspectives of streaming FL, concept drift in FL, and federated continual learning.

Streaming Federated Learning. To better capture real-world scenarios where data arrives continuously, recent studies have extended FL to the streaming setting. From a theoretical perspective, Marfoq et al. (2023) formalize FL over data streams and propose a general FL algorithm through a weighted empirical risk minimization, but still assume stationary distributions. Huynh et al. (2025) investigate convergence under data streams modeled by non-stationary Markov processes. From a data sampling perspective, ODE (Gong et al., 2023) introduces a selection method based on a data evaluation metric under storage constraints, but also relies on the stationary assumption. Wang et al. (2023a) assume the existence of a true long-term distribution in the local data stream of clients and proposes the cache update strategy to align the data distribution in the local cache to the underlying long-term distribution. Fed-HIST (Zhang et al., 2024) avoids the problem of raw data storage by retrieving model-based historical representations through similarity comparison. DYNAMITE (Liu et al., 2023) optimizes batch size and aggregation frequency under dynamic conditions, but it uses reservoir sampling and does not adjust data selection based on distribution changes. On the system side, Jin et al. (2021) formulate a latency-minimizing FL scheduling problem with online and bandit algorithms under budget and network constraints. Hu et al. (2024) develop a Lyapunov-based resource management scheme for streaming FL with adaptive control over computation and communication under long-term energy constraints. FedStream (Wang et al., 2023b) tackles dual heterogeneity in data and arrival patterns with asynchronous aggregation and local adaptation, but lacks theoretical justification. Different from the above studies, our work considers temporally evolving data distributions and develops a distribution-aware data sampling strategy grounded in theoretical analysis.

Concept Drift in Federated Learning. While concept drift (Gama et al., 2014; Lu et al., 2018; Tahmasbi et al., 2021) has been widely studied in traditional machine learning, existing solutions often cannot be directly applied to FL due to the inherent heterogeneity across clients. Several works focus on detecting and responding to drift at the client level. FedConD (Chen et al., 2021) detects local drift based on historical model performance and adapts by adjusting the local regularization parameters. Flash (Panchal et al., 2023) detects drift via the magnitude of client updates and adapts the learning rate accordingly, while AdapFedAvg (Canonaco et al., 2021) passively adjusts learning rates to improve model plasticity under drift. Manias et al. (2021) detect drifted clients using dimensionality reduction and clustering on model updates, primarily aiming at client isolation rather than adaptation. Other approaches address drift through clustering-based strategies. Fielding (Li et al., 2024) detects concept drift via label distribution changes and selectively re-clusters clients to preserve cluster quality under heterogeneity. FedDrift (Jothimurugesan et al., 2023) formalizes staggered drift adaptation as a time-varying clustering problem and proposes hierarchical clustering algorithms guided by local drift detection. FedRC (Guo et al., 2024) proposes a bi-level optimization framework based on a clustering principle to address simultaneous feature, label distribution shifts, and concept shift in FL. FedCCFA (Chen et al., 2024) aligns client feature spaces under distributed concept drift by combining classifier clustering and entropy-based adaptive feature alignment. Most concept drift adaptation methods typically modify the model architecture (Chen et al., 2021), optimization strategy (Panchal et al., 2023; Canonaco et al., 2021), or client clustering (Li et al., 2024; Jothimurugesan et al., 2023; Chen et al., 2024) in response to detected shifts. In contrast, we propose a theoretically grounded client-specific data sampling strategy and model aggregation algorithm in streaming FL.

Federated Continual Learning. Federated continual learning (FCL) (Yang et al., 2024; Wang et al., 2024) addresses evolving data but primarily focuses on mitigating catastrophic forgetting across sequential tasks. Guo et al. (2021) propose a FCL framework based on approximating prior

local objectives, while FedWeIT (Yoon et al., 2021) decomposes model weights into global and task-specific components to enable selective inter-client knowledge transfer. AF-FCL (Wuerkaixi et al., 2024) introduces a selective generative replay method for FCL that emphasizes accurate forgetting to discard biased knowledge across heterogeneous clients. FedCIL (Qi et al., 2023) introduces model consolidation and consistency enforcement to stabilize training on non-IID streaming tasks without storing historical data. Huang et al. (2022) and Usmanova et al. (2021) both incorporate knowledge distillation strategies to enhance generalization and mitigate forgetting in federated settings. GLFC (Dong et al., 2022) further proposes class-aware gradient compensation and proxybased global model selection to handle class-incremental learning with dynamic client participation. While these FCL approaches focus on mitigating catastrophic forgetting across sequential tasks or enhancing model generalization under domain shifts, they do not explicitly model or predict the temporal evolution of data distributions. In contrast, our work leverages a predictive oracle to guide data sampling and optimization in streaming FL.

C DISCUSSIONS

C.1 DISCUSSION ON THE PREDICTION ORACLE

In the main body of our work, we assume the existence of a prediction oracle that provides each client n with an estimated state distribution $\hat{\pi}_n$. In this section, we further discuss the applicability and plausibility of this assumption, clarifying under what conditions such an oracle is reasonable and how it can be instantiated in practice.

The prediction oracle assumption becomes reasonable when two conditions are satisfied: (i) the client data naturally evolves over time, leading to state transitions in the underlying distribution, and (ii) these transitions exhibit structural patterns that can be learned from historical observations. As examples, we present three representative application scenarios:

- Mobility-driven environments. In UAV networks or mobile sensing platforms, client data distributions shift due to physical movement across regions. In such settings, the underlying state can be naturally defined by spatial regions (e.g., grid cells or points of interest). Moreover, mobility patterns such as constrained flight paths, periodic patrol routes, or Markovian movement models provide structured trajectories from which the state distribution can be effectively inferred. As a result, the prediction oracle can leverage historical mobility traces to estimate future state distributions with reasonable accuracy.
- Periodicity-driven environments. In many domains, data distributions evolve according to
 recurring temporal patterns. Typical examples include transportation systems, where traffic
 intensity varies between rush hours and off-peak periods, and environmental monitoring, where
 sensor readings change with daily or seasonal cycles. In such cases, states can be naturally
 defined based on temporal segments. Since these periodic structures repeat over time, state
 distributions can be reliably estimated using historical data through statistical models.
- Interaction-driven environments. In recommendation systems, user—item interactions evolve in real time, producing non-stationary feedback streams. Here, the state can be defined as a representation of the user's latent preference profile, which evolves as the user interacts with new items. For instance, preference shifts may correspond to transitions between clusters of interaction features (e.g., genres of movies, categories of products, or communities of social content). Although individual actions are often noisy, the aggregated behavior of users tends to reveal structural dynamics such as preference drifts, trending items, or temporal co-occurrence patterns. By clustering interaction features and modeling transitions across preference states, these dynamics can be captured using online learning or probabilistic estimators.

In all these scenarios, the oracle does not need to be perfectly accurate: even approximate estimates of the state distribution, obtained through lightweight predictors or Bayesian updates from past data, are sufficient for guiding the DDS and SAW mechanisms in our framework. This demonstrates that the prediction oracle is not a restrictive abstraction but rather a broadly applicable tool in streaming federated learning settings. Moreover, our experimental results (Section 6 and Appendix H) further confirm that the framework remains robust under various levels of prediction inaccuracies. To make this abstraction more concrete, we next illustrate how such an oracle can be instantiated in practice through a Bayesian estimation framework.

Bayesian Estimation Framework. Suppose that at time t, client n has observed a sequence of states $S_{n,t} = \{s_1, s_2, \ldots, s_t\}$, where each $s_i \in \{1, 2, \ldots, M\}$. We model the state generation process as a categorical distribution parameterized by π_n . To estimate π_n , we adopt a Bayesian approach by placing a Dirichlet prior:

$$\boldsymbol{\pi}_n \sim \mathrm{Dir}(\boldsymbol{\alpha}_0),$$

where $\alpha_0 = (\alpha_0^1, \alpha_0^2, \dots, \alpha_0^M)$ is a prior belief.

Given the observed state counts $\mathbf{c}_t = (c_t^1, c_t^2, \dots, c_t^M)$, where $c_t^m = \sum_{i=1}^t \mathbb{I}[s_i = m]$, the posterior distribution becomes:

$$\pi_n \mid \mathcal{S}_{n,t} \sim \text{Dir}(\alpha_0 + \mathbf{c}_t),$$

and the posterior mean is:

$$\hat{\pi}_{n,m} = \mathbb{E}[\pi_{n,m} \mid \mathcal{S}_{n,t}] = \frac{\alpha_0^m + c_t^m}{\sum_{j=1}^M (\alpha_0^j + c_t^j)}.$$

The posterior variance of each component π_n^m under the Dirichlet distribution is given by:

$$\operatorname{Var}[\pi_{n}^{m} \mid \mathcal{S}_{n,t}] = \frac{(\alpha_{0}^{m} + c_{t}^{m}) (S_{t} - \alpha_{0}^{m} - c_{t}^{m})}{S_{t}^{2}(S_{t} + 1)},$$

where
$$S_t = \sum_{j=1}^{M} (\alpha_0^j + c_t^j) = \sum_{j=1}^{M} \alpha_0^j + t$$
.

This expression reveals that the estimation error decreases as the number of observed samples increases (i.e., as c_t^m grows), leading to more confident and accurate estimates over time. The variance scales approximately as $\mathcal{O}(1/t)$.

The above procedure provides a concrete and theoretically grounded method for constructing a prediction oracle. It enables dynamic tracking of local state distributions with low overhead, and its estimation error diminishes with increasing historical information. This property is especially suitable for streaming federated learning settings, where clients accumulate observations over time and require increasingly accurate guidance for data sampling and model aggregation.

C.2 DISCUSSION ON THE HETEROGENEITY BOUNDS

While our theoretical framework defines d_m as an upper bound on the gradient variance under state m, estimating gradient variance directly is often impractical in streaming scenarios with limited and dynamic local data.

Inspired by prior works such as Fed-CBS (Zhang et al., 2023) and FedDisco (Ye et al., 2023), we approximate d_m using distance metrics of the class distribution, based on the intuition that skewed class distributions often induce unstable local updates, which in turn imply higher gradient variance. This approximation enables a practical and data-accessible estimation of d_m .

We further provide a theoretical insight to support this approximation. Specifically, the local objective gradient $\nabla F_m(x)$ on client m can be decomposed over its label distribution (Zhang et al., 2023):

$$\nabla F_m(x) = \sum_{c=1}^C p_m^c \nabla F_m^c(x),$$

where p_m^c is the proportion of class c on client m, and $\nabla F_m^c(x)$ is the gradient of the loss conditioned on class c. Then, we note that:

$$||\nabla F_m(x) - \nabla F(x)||^2 = ||\sum_{m=1}^C (p_m^c - \frac{1}{C})\nabla F_m^c(x)||^2 \le \sum_{m=1}^C (p_m^c - \frac{1}{C})^2 \cdot \sum_{m=1}^C ||\nabla F_m^c(x)||^2.$$

Here, $\sum_{c=1}^{C} (p_m^c - \frac{1}{C})^2$ can be viewed as a kind of distance between the class distribution and the uniform distribution. Therefore, we have a strong intuition that there is a linear relationship between d_m and the distance.

To support this approximation, we present an empirical analysis in Appendix H comparing four metrics—L1, L2, KL, and JS divergence—and find a consistent correlation between distributional divergence and update instability. These results validate the suitability of our approximation in realistic settings.

NOTATIONS AND TECHNICAL LEMMAS

D.1 NOTATIONS

Table 4 summarizes the notations appearing in this paper.

Table 4: Summary of key notations.

Symbol	Description
R, r	number, index of training rounds
N, n	number, index of clients
M, m	number, index of latent states
T, t	number, index of local update steps
S_r	set of participating clients in round r
\mathcal{D}_m	the data distribution of latent state m
$\mathcal{D}_{n,t}^{(r)}$	the data distribution of client n at time t in round r
$oldsymbol{\pi}_n$	$(\pi_{n,1},\ldots,\pi_{n,M})$ is a round-truth state distribution of client n
$\pi_{n,m}$	probability that client n is in state m
$\hat{\boldsymbol{\pi}}_n$	$(\hat{\pi}_{n,1},\ldots,\hat{\pi}_{n,M})$ is a prediction oracle of client n
$\hat{\pi}_{n,m}$	estimate of $\pi_{n,m}$
$\alpha_{n,m}$	the fraction of new samples from state m that client n samples
p_n	the aggregation weight
q_n	the available probability
η	learning rate (or stepsize)
L	L-smoothness constant (Asm. 1)
σ^2	upper bound on variance of stochastic gradients at each client (Asm. 2)
G	constant in Asm. 3 to bound the overall gradient heterogeneity across all latent states
d_m	constants in Asm. 4 to bound the state-specific gradient heterogeneity
$F/F_{n,t}^{(r)}$	global objective/local objective of client n at time t in the r -th round
F_m	expected loss function on the distribution \mathcal{D}_m of latent state m
w_m	the weight of F_m in the global objective
$ar{\mathbf{x}}^{(r)}$	global model parameters in the r -th round
$\mathbf{x}_{n,t}^{(r)}$	local model parameters of client n after t local steps in the r -th round
$\mathbf{g}_{n,t}^{(r)}$	$\mathbf{g}_{n,t}^{(r)} \triangleq \nabla f(\mathbf{x}_{n,t-1}^{(r)}; \xi_{n,t}) \text{ denotes the stochastic gradients of } F_{n,t}^{(r)} \text{ regarding } \mathbf{x}_{n,t-1}^{(r)}$

D.2 LEMMAS

Jensen's inequality. Let $h: \mathbb{R}^d \to \mathbb{R}$ be a convex function. For any vectors $\boldsymbol{x}_1, \dots, \boldsymbol{x}_n \in \mathbb{R}^d$ and any non-negative weights $\lambda_1, \dots, \lambda_n$ satisfying $\sum_{i=1}^n \lambda_i = 1$, it holds that

$$h\left(\sum_{i=1}^{n} \lambda_i \boldsymbol{x}_i\right) \le \sum_{i=1}^{n} \lambda_i h(\boldsymbol{x}_i). \tag{17}$$

As a special case with $h(x) = ||x||^2$, we obtain

$$\left\| \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{x}_{i} \right\|^{2} \leq \frac{1}{n} \sum_{i=1}^{n} \|\boldsymbol{x}_{i}\|^{2}.$$
 (18)

Lemma 1. Let $\{\mathbf{z}_t\}_{t=1}^T$ be a sequence of random vectors adapted to the underlying filtration $\{\mathcal{F}_t\}$ such that for all t, $\mathbb{E}[\mathbf{z}_t|\mathcal{F}_{t-1}] = 0$ and $\mathbb{E}[\|\mathbf{z}_t\|^2|\mathcal{F}_{t-1}] \leq \sigma^2$. Then, the following identity holds:

$$\mathbb{E}\Big[\Big\|\sum_{t=1}^{T}\mathbf{z}_{t}\Big\|^{2}\Big] \le T\sigma^{2}.\tag{19}$$

This lemma provides a standard bound on the variance of a martingale difference sequence and will be used to control the accumulation of stochastic gradient noise over time.

Lemma 2. Suppose we have a sequence $\{a_t\}$ satisfying the recursive inequality $a_t \leq (1-\alpha)a_{t-1} + \beta$, where $\alpha \in (0,1)$ and $\beta > 0$ are constants. Then, for any $t \geq 1$, the sequence is bounded as:

$$a_t \le a_0 (1 - \alpha)^t + \frac{\beta}{\alpha} (1 - (1 - \alpha)^t).$$
 (20)

Proof. We prove this by iteratively unfolding the recursion. Divide both sides of the inequality by $(1-\alpha)^t$:

$$\frac{a_t}{(1-\alpha)^t} \le \frac{a_{t-1}}{(1-\alpha)^{t-1}} + \frac{\beta}{(1-\alpha)^t}$$
(21)

$$\leq \frac{a_{t-2}}{(1-\alpha)^{t-2}} + \frac{\beta}{(1-\alpha)^{t-1}} + \frac{\beta}{(1-\alpha)^t}$$
 (22)

$$\leq a_0 + \sum_{i=1}^t \frac{\beta}{(1-\alpha)^i}.\tag{24}$$

Multiplying both sides by $(1 - \alpha)^t$ yields:

$$a_t \le (1 - \alpha)^t \left(a_0 + \sum_{i=1}^t \frac{\beta}{(1 - \alpha)^i} \right)$$
 (25)

$$= a_0 (1 - \alpha)^t + \beta \sum_{i=1}^t (1 - \alpha)^{t-i}$$
 (26)

$$= a_0 (1 - \alpha)^t + \frac{\beta}{\alpha} (1 - (1 - \alpha)^t).$$
 (27)

E PROOF OF THEOREM 1

 The global objective function is $F(\mathbf{x}) = \sum_{m=1}^{M} w_m F_m(\mathbf{x})$. We begin by analyzing the global model update across one communication round, which can be expressed as:

$$\Delta \mathbf{x} := \bar{\mathbf{x}}^{(r+1)} - \bar{\mathbf{x}}^{(r)} = -\eta \sum_{n \in S_r} p_n \sum_{t=1}^T \mathbf{g}_{n,t}^{(r)},$$
(28)

where $\mathbf{g}_{n,t}^{(r)}$ denotes the stochastic gradient computed by client n at time step t in round r. According to Assumption 2, these gradients are unbiased estimators of the local loss, i.e., $\mathbb{E}[\mathbf{g}_{n,t}^{(r)}] = \nabla F_{n,t}^{(r)}(\mathbf{x}_{n,t}^{(r)})$.

In the following, we focus on a single training round, and hence we drop the superscript (r) for clarity. For example, we write $\mathbf{x}_{n,t}$ for $\mathbf{x}_{n,t}^{(r)}$ and use $F_{n,t}$ to replace $F_{n,t}^{(r)}$. Moreover, let \mathbf{x} denote the initial local model $\mathbf{x}_{n,0}^{(r)}$. Unless otherwise stated, the expectation is conditioned on the global model $\bar{\mathbf{x}}^{(r)}$ and the participating client set S_r . Since $F(\cdot)$ is L-smooth, we get:

$$\mathbb{E}\left[F(\mathbf{x} + \Delta \mathbf{x}) - F(\mathbf{x})\right] \le \mathbb{E}\left[\left\langle \nabla F(\mathbf{x}), \Delta \mathbf{x} \right\rangle\right] + \frac{L}{2} \mathbb{E}\|\Delta \mathbf{x}\|^{2}.$$
 (29)

We now proceed to bound the two terms in the RHS of (29) separately.

Bounding $\mathbb{E}[\langle \nabla F(\mathbf{x}), \Delta \mathbf{x} \rangle]$ in (29).

$$A := \mathbb{E}\left[\langle \nabla F(\mathbf{x}), \Delta \mathbf{x} \rangle\right] = -\eta \sum_{n \in S_r} p_n \sum_{t=1}^T \mathbb{E}\left[\langle \nabla F(\mathbf{x}), \nabla F_{n,t}(\mathbf{x}_{n,t-1}) \rangle\right]$$

$$= -\frac{\eta}{2} \sum_{n \in S_r} p_n \sum_{t=1}^T \|\nabla F(\mathbf{x})\|^2 - \frac{\eta}{2} \sum_{n \in S_r} p_n \sum_{t=1}^T \mathbb{E}\|\nabla F_{n,t}(\mathbf{x}_{n,t-1})\|^2$$

$$+ \frac{\eta}{2} \sum_{n \in S_r} p_n \sum_{t=1}^T \|\nabla F(\mathbf{x})\|^2 - \frac{\eta}{2} \sum_{n \in S_r} p_n \sum_{t=1}^T \mathbb{E}\|\nabla F_{n,t}(\mathbf{x}_{n,t-1})\|^2$$

$$\leq -\frac{\eta}{2} \sum_{n \in S_r} p_n \sum_{t=1}^T \|\nabla F(\mathbf{x})\|^2 - \frac{\eta}{2} \sum_{n \in S_r} p_n \sum_{t=1}^T \mathbb{E}\|\nabla F_{n,t}(\mathbf{x}_{n,t-1})\|^2$$

$$+ \frac{\eta}{2} \sum_{n \in S_r} p_n \sum_{t=1}^T \left(2\mathbb{E}\|\nabla F_{n,t}(\mathbf{x}_{n,t-1}) - \nabla F_{n,t}(\mathbf{x})\|^2 + 2\mathbb{E}\|\nabla F_{n,t}(\mathbf{x}) - \nabla F(\mathbf{x})\|^2\right)$$

$$\leq -\frac{\eta}{2} \sum_{n \in S_r} p_n \sum_{t=1}^T \|\nabla F(\mathbf{x})\|^2 - \frac{\eta}{2} \sum_{n \in S_r} p_n \sum_{t=1}^T \mathbb{E}\|\nabla F_{n,t}(\mathbf{x}_{n,t-1})\|^2$$

$$+ \frac{\eta}{2} \sum_{n \in S_r} p_n \sum_{t=1}^T (2L^2 \mathbb{E}\|\mathbf{x}_{n,t-1} - \mathbf{x}\|^2 + 2\mathbb{E}\|\nabla F_{n,t}(\mathbf{x}) - \nabla F(\mathbf{x})\|^2)$$

$$= -\frac{T\eta}{2} \sum_{n \in S_r} p_n \|\nabla F(\mathbf{x})\|^2 - \frac{\eta}{2} \sum_{n \in S_r} p_n \sum_{t=1}^T \mathbb{E}\|\nabla F_{n,t}(\mathbf{x}_{n,t-1})\|^2$$

$$+ L^2 \eta \sum_{n \in S_r} p_n \sum_{t=1}^T \mathbb{E}\|\mathbf{x}_{n,t-1} - \mathbf{x}\|^2 + \eta \sum_{t=1}^T p_n \sum_{t=1}^T \mathbb{E}\|\nabla F_{n,t}(\mathbf{x}) - \nabla F(\mathbf{x})\|^2,$$
(34)

where (31) applies the identity $\langle a,b\rangle=\frac{1}{2}\|a\|^2+\frac{1}{2}\|b\|^2-\frac{1}{2}\|a-b\|^2$, (32) uses the inequality $\|a-b\|^2\leq 2\|a-c\|^2+2\|b-c\|^2$ by inserting the intermediate term $\nabla F_{n,t}(\mathbf{x})$, and then (33) uses L-smoothness.

Bounding $\frac{L}{2}\mathbb{E}\|\Delta\mathbf{x}\|^2$ in (29).

$$B := \frac{L}{2} \mathbb{E} \|\Delta \mathbf{x}\|^2 = \frac{L}{2} \eta^2 \mathbb{E} \left\| \sum_{n \in S_-} p_n \sum_{t=1}^T \mathbf{g}_{n,t} \right\|^2$$

$$(35)$$

$$\leq L\eta^{2} \mathbb{E} \left\| \sum_{n \in S_{r}} p_{n} \sum_{t=1}^{T} \left(\mathbf{g}_{n,t} - \nabla F_{n,t}(\mathbf{x}_{n,t-1}) \right) \right\|^{2} + L\eta^{2} \mathbb{E} \left\| \sum_{n \in S_{r}} p_{n} \sum_{t=1}^{T} \nabla F_{n,t}(\mathbf{x}_{n,t-1}) \right\|^{2}$$
(36)

$$\leq L\eta^{2} \sum_{n \in S_{n}} p_{n}^{2} \cdot T\sigma^{2} + L\eta^{2} |S_{r}| \sum_{n \in S_{n}} p_{n}^{2} \cdot T \sum_{t=1}^{T} \mathbb{E} \|\nabla F_{n,t}(\mathbf{x}_{n,t-1})\|^{2}, \tag{37}$$

where (36) uses the inequality $||a+b||^2 \le 2||a||^2 + 2||b||^2$. The first term in (37) is based on Lemma 1 and uses the fact that clients are independent of each other. The second term in (37) uses the Jensen's inequality. Combining bounds for A and B, and substituting back into (29), we obtain:

$$\mathbb{E}\left[F(\mathbf{x} + \Delta \mathbf{x}) - F(\mathbf{x})\right]$$

$$\leq -\frac{T\eta}{2} \sum_{n \in S_r} p_n \|\nabla F(\mathbf{x})\|^2 + LT\eta^2 \sigma^2 \sum_{n \in S_r} p_n^2 + \sum_{n \in S_r} \left(LT\eta^2 |S_r| p_n^2 - \frac{\eta}{2} p_n\right) \sum_{t=1}^T \mathbb{E}\|\nabla F_{n,t}(\mathbf{x}_{n,t-1})\|^2$$

$$+ L^2 \eta \sum_{n \in S_r} p_n \sum_{t=1}^T \mathbb{E}\|\mathbf{x}_{n,t-1} - \mathbf{x}\|^2 + \eta \sum_{n \in S_r} p_n \sum_{t=1}^T \mathbb{E}\|\nabla F_{n,t}(\mathbf{x}) - \nabla F(\mathbf{x})\|^2$$

$$\leq -\frac{T\eta}{2} \sum_{n \in S_r} p_n \|\nabla F(\mathbf{x})\|^2 + LT\eta^2 \sigma^2 \sum_{n \in S_r} p_n^2 + L^2 \eta \sum_{n \in S_r} p_n \sum_{t=1}^T \mathbb{E}\|\mathbf{x}_{n,t-1} - \mathbf{x}\|^2$$

$$+ \eta \sum_{n \in S_r} p_n \sum_{t=1}^T \mathbb{E}\|\nabla F_{n,t}(\mathbf{x}) - \nabla F(\mathbf{x})\|^2, \tag{39}$$

where (39) holds under the condition that for all $n \in S_r$ we have $LT|S_r|\eta^2 p_n^2 \le \frac{\eta}{2}p_n$, which is naturally satisfied when $LT\eta \le \frac{1}{2M}$.

Bounding
$$\sum_{t=1}^T \mathbb{E} \|\mathbf{x}_{n,t-1} - \mathbf{x}\|^2$$
 in (39).

$$C := \sum_{t=1}^{T} \mathbb{E} \|\mathbf{x}_{n,t-1} - \mathbf{x}\|^{2} = \sum_{t=1}^{T} \mathbb{E} \| - \eta \sum_{i=1}^{t-1} \mathbf{g}_{n,i} \|^{2}$$

$$\leq 4\eta^{2} \sum_{t=1}^{T} \mathbb{E} \left[\left\| \sum_{i=1}^{t-1} \mathbf{g}_{n,i} - \sum_{i=1}^{t-1} \nabla F_{n,i}(\mathbf{x}_{n,i-1}) \right\|^{2} + \left\| \sum_{i=1}^{t-1} \nabla F_{n,i}(\mathbf{x}_{n,i-1}) - \sum_{i=1}^{t-1} \nabla F_{n,i}(\mathbf{x}) \right\|^{2}$$

$$+ \left\| \sum_{i=1}^{t-1} \nabla F_{n,i}(\mathbf{x}) - \sum_{i=1}^{t-1} \nabla F(\mathbf{x}) \right\|^{2} + \left\| \sum_{i=1}^{t-1} \nabla F(\mathbf{x}) \right\|^{2} \right]$$

$$\leq 4\eta^{2} \sum_{t=1}^{T} (t-1)\sigma^{2} + 4\eta^{2} \sum_{t=1}^{T} (t-1)L^{2} \sum_{i=1}^{t-1} \mathbb{E} \|\mathbf{x}_{n,i-1} - \mathbf{x}\|^{2}$$

$$+ 4\eta^{2} \sum_{t=1}^{T} (t-1) \sum_{t=1}^{t-1} \mathbb{E} \|\nabla F_{n,i}(\mathbf{x}) - \nabla F(\mathbf{x}) \|^{2} + 4\eta^{2} \sum_{t=1}^{T} (i-1)^{2} \|\nabla F(\mathbf{x}) \|^{2}$$

$$(40)$$

$$\leq 2T^{2}\eta^{2}\sigma^{2} + 2L^{2}T^{2}\eta^{2} \cdot C + 2T^{2}\eta^{2} \sum_{t=1}^{T} \mathbb{E} \|\nabla F_{n,t}(\mathbf{x}) - \nabla F(\mathbf{x})\|^{2} + \frac{4}{3}T^{3}\eta^{2} \|\nabla F(\mathbf{x})\|^{2}, \quad (43)$$

where (41) uses the Jensen's inequality. The first term in (42) is based on Lemma 1. The second term in (37) uses the Jensen's inequality and L-smoothness. The third term in (42) uses the Jensen's inequality. (43) uses the fact that $\sum_{t=1}^{T} (t-1) \leq \frac{T^2}{2}$ and $\sum_{t=1}^{T} (t-1)^2 \leq \frac{T^3}{3}$.

After rearranging the preceding inequality and using $2L^2T^2\eta^2 < 1$, we get

$$C \le \frac{1}{1 - 2L^2 T^2 \eta^2} \left[2T^2 \eta^2 \sigma^2 + 2T^2 \eta^2 \sum_{t=1}^T \mathbb{E} \left\| \nabla F_{n,i}(\mathbf{x}) - \nabla F(\mathbf{x}) \right\|^2 + \frac{4}{3} T^3 \eta^2 \| \nabla F(\mathbf{x}) \|^2 \right]. \tag{44}$$

We denote that $c=LT\eta$. Then, plugging the bound of C into (39), we get

$$\mathbb{E}\left[F(\mathbf{x}+\Delta\mathbf{x}) - F(\mathbf{x})\right] \leq \left(\frac{4c^2}{3(1-2c^2)} - \frac{1}{2}\right)T\eta \sum_{n \in S_r} p_n \|\nabla F(\mathbf{x})\|^2 + LT\eta^2 \sigma^2 \sum_{n \in S_r} p_n^2 + \frac{2L^2T^2\eta^3\sigma^2}{1-2c^2} \sum_{n \in S} p_n + \left(\eta + \frac{2L^2T^2\eta^3}{1-2c^2}\right) \sum_{n \in S} p_n \sum_{t=1}^T \mathbb{E}\|\nabla F_{n,t}(\mathbf{x}) - \nabla F(\mathbf{x})\|^2. \tag{45}$$

Bounding $\|\nabla F_{n,t}(\mathbf{x}) - \nabla F(\mathbf{x})\|^2$ in (45).

Recall that the local data distribution evolves according to $\mathcal{D}_{n,t} = (1 - \alpha_{n,m_{n,t}}) \cdot \mathcal{D}_{n,t-1} + \alpha_{n,m_{n,t}} \cdot \mathcal{D}_{m_{n,t}}$, where $m_{n,t} \in \mathcal{M}$ denotes the latent state of client n at time step t, and $\alpha_{n,m_{n,t}} \in [0,1]$ is the sampling ratio defined in Section 3. Due to the linearity of expectation, the corresponding loss function also follows a similar recursive relationship:

$$F_{n,t} = (1 - \alpha_{n,m_{n,t}})F_{n,t-1} + \alpha_{n,m_{n,t}}F_{m_{n,t}}, \tag{46}$$

where $F_{m_{n,t}}$ denotes the loss associated with $\mathcal{D}_{m_{n,t}}$. We apply the Jensen's inequality and obtain:

$$\|\nabla F_{n,t}(\mathbf{x}) - \nabla F(\mathbf{x})\|^{2} \le (1 - \alpha_{n,m_{n,t}}) \|\nabla F_{n,t-1}(\mathbf{x}) - \nabla F(\mathbf{x})\|^{2} + \alpha_{n,m_{n,t}} \|\nabla F_{m_{n,t}}(\mathbf{x}) - \nabla F(\mathbf{x})\|^{2}.$$
(47)

We now take expectation on both sides of (47) with respect to the latent state $m_{n,t}$, under the ground-truth state distribution. By using the Assumption 4, we get:

$$\mathbb{E} \|\nabla F_{n,t}(\mathbf{x}) - \nabla F(\mathbf{x})\|^2 \tag{48}$$

$$\leq \left(1 - \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m}\right) \mathbb{E} \left\| \nabla F_{n,t-1}(\mathbf{x}) - \nabla F(\mathbf{x}) \right\|^{2} + \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} \mathbb{E} \left\| \nabla F_{m}(\mathbf{x}) - \nabla F(\mathbf{x}) \right\|^{2}$$
(49)

$$\leq \left(1 - \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m}\right) \mathbb{E} \left\| \nabla F_{n,t-1}(\mathbf{x}) - \nabla F(\mathbf{x}) \right\|^2 + \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} d_m. \tag{50}$$

We now apply Lemma 2 to the recursive bound derived in (50). Let $a_t = \mathbb{E} \|\nabla F_{n,t}(\mathbf{x}) - \nabla F(\mathbf{x})\|^2$, $\alpha = \alpha_n \triangleq \sum_{m=1}^M \pi_{n,m} \alpha_{n,m}$, and $\beta = \beta_n \triangleq \sum_{m=1}^M \pi_{n,m} \alpha_{n,m} d_m$. Then, we get

$$\sum_{t=1}^{T} \mathbb{E} \left\| \nabla F_{n,t}(\mathbf{x}) - \nabla F(\mathbf{x}) \right\|^{2} \le \frac{\beta_{n}}{\alpha_{n}} \left(T - \sum_{t=1}^{T} (1 - \alpha_{n})^{t} \right) + \sum_{t=1}^{T} (1 - \alpha_{n})^{t} \mathbb{E} \left\| \nabla F_{n,0}(\mathbf{x}) - \nabla F(\mathbf{x}) \right\|^{2}. \tag{51}$$

Plugging this result into the bound in (45), we obtain:

$$\mathbb{E}\left[F(\mathbf{x} + \Delta \mathbf{x}) - F(\mathbf{x})\right] \\
\leq \left(\frac{4c^{2}}{3(1 - 2c^{2})} - \frac{1}{2}\right)T\eta \sum_{n \in S_{r}} p_{n} \|\nabla F(\mathbf{x})\|^{2} + LT\eta^{2}\sigma^{2} \sum_{n \in S_{r}} p_{n}^{2} + \frac{2L^{2}T^{2}\eta^{3}\sigma^{2}}{1 - 2c^{2}} \sum_{n \in S_{r}} p_{n} \\
+ \left(T\eta + \frac{2L^{2}T^{3}\eta^{3}}{1 - 2c^{2}}\right) \sum_{n \in S_{r}} p_{n} \left[\frac{\beta_{n}}{\alpha_{n}}(1 - \gamma_{n}) + \gamma_{n}\mathbb{E}\|\nabla F_{n,0}(\mathbf{x}) - \nabla F(\mathbf{x})\|^{2}\right]. \tag{52}$$

where $\gamma_n = \frac{1}{T} \sum_{t=1}^{T} (1 - \alpha_n)^t \in (0, 1)$.

Bounding $\mathbb{E}ig\|
abla F_{n,0}(\mathbf{x}) -
abla F(\mathbf{x})ig\|^2$ in (52).

We observe that the local loss function $F_{n,t}^{(r)}$ can be expressed as a weighted mixture of loss functions associated with latent states $m \in \mathcal{M}$. Specifically, we suppose $F_{n,t}^{(r)} = \langle \vec{w}_{n,r,t}, \vec{F} \rangle$, where $\vec{w}_{n,r,t} = (w_{n,r,t}^{(1)}, w_{n,r,t}^{(2)}, \ldots, w_{n,r,t}^{(M)})$ denotes the mixture weights over the latent states and $\vec{F} = (F_1, F_2, \ldots, F_M)$ represents the vector of corresponding state-specific loss functions. Based on the recursive update rule of the local loss function in (46), the weight vector $\vec{w}_{n,r,t}$ evolves as

$$\vec{\mathbf{w}}_{n,r,t} = (1 - \alpha_{n,m_{n,t}})\vec{\mathbf{w}}_{n,r,t-1} + \alpha_{n,m_{n,t}}\vec{\mathbf{e}}_{m_{n,t}}, \tag{53}$$

where \vec{e}_m denotes the one-hot unit vector corresponding to latent state m, with its m-th component equal to 1 and all others set to 0. By applying the Cauchy-Schwarz inequality, we get:

$$\|\nabla F_{n,0}^{(r)}(\bar{\mathbf{x}}^{(r)}) - \nabla F(\bar{\mathbf{x}}^{(r)})\|^2 \le \|\vec{\boldsymbol{w}}_{n,r,0} - \vec{\boldsymbol{w}}\|^2 \cdot \sum_{m=1}^{M} \|\nabla F_m(\bar{\mathbf{x}}^{(r)})\|^2 \le G\|\vec{\boldsymbol{w}}_{n,r,0} - \vec{\boldsymbol{w}}\|^2.$$
 (54)

where $\vec{w} = (w_1, w_2, \dots, w_M)$ denotes the weight vector in the global objective, and Assumption 3 is applied.

To further bound the discrepancy between the client-specific mixture weights and the global target weights, we decompose the $\|\vec{w}_{n,r,t} - \vec{w}\|^2$ into two parts using the Jensen's inequality:

$$\|\vec{w}_{n,r,t} - \vec{w}\|^{2} \le 2 \|\vec{w}_{n,r,t} - \frac{1}{\alpha_{n}} \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} \vec{e}_{m}\|^{2} + 2 \|\frac{1}{\alpha_{n}} \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} \vec{e}_{m} - \vec{w}\|^{2}.$$
 (55)

The first term in equation (55) quantifies the deviation between the actual weight vector and a fixed weight vector, which is fully determined by the distribution-aware data sampling strategy. The second term describes the inherent mismatch between the weights determined by the strategy and the global target weights. We deal with the first term in the following.

Bounding
$$\left\|\vec{w}_{n,r,t} - \frac{1}{\alpha_n}\sum_{m=1}^{M}\pi_{n,m}\alpha_{n,m}\vec{e}_m\right\|^2$$
 in (55).

$$\mathbb{E} \left\| \vec{\boldsymbol{w}}_{n,r,t} - \frac{1}{\alpha_n} \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} \vec{\boldsymbol{e}}_m \right\|^2$$
 (56)

$$= \mathbb{E} \|\vec{\boldsymbol{w}}_{n,r,t}\|^2 - \frac{2}{\alpha_n} \mathbb{E} \left[\left\langle \vec{\boldsymbol{w}}_{n,r,t}, \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} \vec{\boldsymbol{e}}_m \right\rangle \right] + \frac{1}{\alpha_n^2} \sum_{m=1}^{M} \pi_{n,m}^2 \alpha_{n,m}^2$$
 (57)

$$= \mathbb{E} \left\| (1 - \alpha_{n,m_{n,t}}) \vec{\boldsymbol{w}}_{n,r,t-1} + \alpha_{n,m_{n,t}} \vec{\boldsymbol{e}}_{m_{n,t}} \right\|^2 - \frac{2}{\alpha_n} \mathbb{E} \left[\left\langle (1 - \alpha_{n,m_{n,t}}) \vec{\boldsymbol{w}}_{n,r,t-1} + \alpha_{n,m_{n,t}} \vec{\boldsymbol{e}}_{m_{n,t}}, \right. \right. \right]$$

$$\sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} \vec{e}_{m} \rangle \Big] + \frac{1}{\alpha_{n}^{2}} \sum_{m=1}^{M} \pi_{n,m}^{2} \alpha_{n,m}^{2}$$
(58)

$$= \sum_{m=1}^{M} \pi_{n,m} (1 - \alpha_{n,m})^{2} \mathbb{E} \|\vec{\boldsymbol{w}}_{n,r,t-1}\|^{2} + 2 \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} (1 - \alpha_{n,m}) \mathbb{E} [\langle \vec{\boldsymbol{w}}_{n,r,t-1}, \vec{\boldsymbol{e}}_{m} \rangle] + \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m}^{2}$$

$$-\frac{2(1-\alpha_n)}{\alpha_n} \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} \mathbb{E}\left[\langle \vec{\boldsymbol{w}}_{n,r,t-1}, \vec{\boldsymbol{e}}_m \rangle\right] - \frac{2}{\alpha_n} \sum_{m=1}^{M} \pi_{n,m}^2 \alpha_{n,m}^2 + \frac{1}{\alpha_n^2} \sum_{m=1}^{M} \pi_{n,m}^2 \alpha_{n,m}^2$$
(59)

$$= \sum_{m=1}^{M} \pi_{n,m} (1 - \alpha_{n,m})^{2} \mathbb{E} \|\vec{\boldsymbol{w}}_{n,r,t-1}\|^{2} + 2 \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} \left(1 - \alpha_{n,m} - \frac{1 - \alpha_{n}}{\alpha_{n}}\right) \mathbb{E} \left[\langle \vec{\boldsymbol{w}}_{n,r,t-1}, \vec{\boldsymbol{e}}_{m} \rangle \right]$$

$$+\sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m}^{2} + \left(\frac{1}{\alpha_{n}^{2}} - \frac{2}{\alpha_{n}}\right) \sum_{m=1}^{M} \pi_{n,m}^{2} \alpha_{n,m}^{2}$$
(60)

$$= (1 - \alpha_n)^2 \mathbb{E} \|\vec{\boldsymbol{w}}_{n,r,t-1}\|^2 - \frac{2(1 - \alpha_n)^2}{\alpha_n} \sum_{m=1}^M \pi_{n,m} \alpha_{n,m} \mathbb{E} \left[\langle \vec{\boldsymbol{w}}_{n,r,t-1}, \vec{\boldsymbol{e}}_m \rangle \right] + \frac{(1 - \alpha_n)^2}{\alpha_n^2} \sum_{m=1}^M \pi_{n,m}^2 \alpha_{n,m}^2$$

$$+ \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m}^{2} - \sum_{m=1}^{M} \pi_{n,m}^{2} \alpha_{n,m}^{2} + 2 \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} (\alpha_{n} - \alpha_{n,m}) \mathbb{E} [\langle \vec{w}_{n,r,t-1}, \vec{e}_{m} \rangle]$$

$$+ \left(\sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m}^{2} - \alpha_{n}^{2} \right) \mathbb{E} \|\vec{\boldsymbol{w}}_{n,r,t-1}\|^{2}$$
(61)

$$= (1 - \alpha_n)^2 \mathbb{E} \left\| \vec{w}_{n,r,t-1} - \frac{1}{\alpha_n} \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} \vec{e}_m \right\|^2 + \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m}^2 - \sum_{m=1}^{M} \pi_{n,m}^2 \alpha_{n,m}^2$$

$$+\sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} (\alpha_{n,m} - \alpha_n) \mathbb{E} \left[\|\vec{\boldsymbol{w}}_{n,r,t-1}\|^2 - 2\langle \vec{\boldsymbol{w}}_{n,r,t-1}, \vec{\boldsymbol{e}}_m \rangle \right]$$
(62)

$$= (1 - \alpha_n)^2 \mathbb{E} \left\| \vec{\boldsymbol{w}}_{n,r,t-1} - \frac{1}{\alpha_n} \sum_{m=1}^M \pi_{n,m} \alpha_{n,m} \vec{\boldsymbol{e}}_m \right\|^2 - \sum_{m=1}^M \pi_{n,m}^2 \alpha_{n,m}^2$$

$$+\sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} (\alpha_{n,m} - \alpha_n) \mathbb{E}\left[\left\| \vec{\boldsymbol{w}}_{n,r,t-1} - \vec{\boldsymbol{e}}_{m} \right\|^{2} \right] + \alpha_{n}^{2}$$

$$(63)$$

$$\leq (1 - \alpha_n)^2 \mathbb{E} \left\| \vec{\boldsymbol{w}}_{n,r,t-1} - \frac{1}{\alpha_n} \sum_{m=1}^M \pi_{n,m} \alpha_{n,m} \vec{\boldsymbol{e}}_m \right\|^2 + 2 \sum_{m=1}^M \pi_{n,m} \alpha_{n,m}^2 - \sum_{m=1}^M \pi_{n,m}^2 \alpha_{n,m}^2 + \alpha_n^2. \tag{64}$$

Here, (61) splits the first term in (60) by applying the identity $\sum_{m=1}^{M} \pi_{n,m} (1 - \alpha_{n,m})^2 = (1 - \alpha_n)^2 + \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m}^2 - \alpha_n^2$, and splits the second term in (60) by applying the identity $1 - \alpha_{n,m} - \frac{1 - \alpha_n}{\alpha_n} = (\alpha_n - \alpha_{n,m}) - \frac{(1 - \alpha_n)^2}{\alpha_n}$. (64) uses the fact that $\|\vec{w}_{n,r,t-1} - \vec{e}_m\|^2 \le 2$.

Then, we set $a_t = \mathbb{E} \|\vec{w}_{n,r,t} - \frac{1}{\alpha_n} \sum_{m=1}^M \pi_{n,m} \alpha_{n,m} \vec{e}_m\|^2$, and apply Lemma 2 to this sequence. Let $\alpha = \alpha_n' \triangleq 1 - (1 - \alpha_n)^2$ and $\beta = \beta_n' \triangleq 2 \sum_{m=1}^M \pi_{n,m} \alpha_{n,m}^2 - \sum_{m=1}^M \pi_{n,m}^2 \alpha_{n,m}^2 + \alpha_n^2$.

We consider the number of times client n has participated in training up to round r, denoted as $\tau_n(r)$. Let r' < r denote the last round that client n participated before round r. Then, since the local weight vector is only updated when the client participates, we have:

$$\mathbb{E} \left\| \vec{\boldsymbol{w}}_{n,r,0} - \frac{1}{\alpha_n} \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} \vec{\boldsymbol{e}}_m \right\|^2 = \mathbb{E} \left\| \vec{\boldsymbol{w}}_{n,r',T} - \frac{1}{\alpha_n} \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} \vec{\boldsymbol{e}}_m \right\|^2 \\
\leq \frac{\beta'_n}{1 - (1 - \alpha_n)^2} + (1 - \alpha_n)^{2T \cdot \tau_n(r-1)} \cdot \mathbb{E} \left\| \vec{\boldsymbol{w}}_{n,0,0} - \frac{1}{\alpha_n} \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} \vec{\boldsymbol{e}}_m \right\|^2 \\
\leq \frac{\beta'_n}{1 - (1 - \alpha_n)^2} + 2(1 - \alpha_n)^{2T \cdot \tau_n(r-1)}, \tag{66}$$

where the second inequality uses the fact that both $\vec{w}_{n,0,0}$ and $\frac{1}{\alpha_n} \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} \vec{e}_m$ are probability vectors supported on the simplex, and hence their squared distance is upper bounded by 2.

Take expectation into the two sides of (52) about the selected client set in the round r, and combining the bounds in (55) and (66), we get

$$\mathbb{E}\left[F(\mathbf{x} + \Delta \mathbf{x}) - F(\mathbf{x})\right] \\
\leq \left(\frac{4c^{2}}{3(1 - 2c^{2})} - \frac{1}{2}\right)T\eta \sum_{n=1}^{N} p_{n}q_{n} \|\nabla F(\mathbf{x})\|^{2} + LT\eta^{2}\sigma^{2} \sum_{n=1}^{N} p_{n}^{2}q_{n} + \frac{2L^{2}T^{2}\eta^{3}\sigma^{2}}{1 - 2c^{2}} \sum_{n=1}^{N} p_{n}q_{n} \\
+ \frac{T\eta}{1 - 2c^{2}} \sum_{n=1}^{N} p_{n}q_{n} \left[\frac{\beta_{n}}{\alpha_{n}}(1 - \gamma_{n}) + \gamma_{n}G\|\vec{w}_{n,r,0} - \vec{w}\|^{2}\right] \\
\leq \left(\frac{4c^{2}}{3(1 - 2c^{2})} - \frac{1}{2}\right)T\eta \sum_{n=1}^{N} p_{n}q_{n} \|\nabla F(\mathbf{x})\|^{2} + LT\eta^{2}\sigma^{2} \sum_{n=1}^{N} p_{n}^{2}q_{n} + \frac{2L^{2}T^{2}\eta^{3}\sigma^{2}}{1 - 2c^{2}} \sum_{n=1}^{N} p_{n}q_{n} \\
+ \frac{T\eta}{1 - 2c^{2}} \sum_{n=1}^{N} p_{n}q_{n} \left[\frac{\beta_{n}}{\alpha_{n}}(1 - \gamma_{n}) + \gamma_{n} \cdot 2G\left(\frac{\beta_{n}'}{1 - (1 - \alpha_{n})^{2}} + 2(1 - \alpha_{n})^{2T\tau_{n}(r-1)}\right) \\
+ \left\|\frac{1}{\alpha_{n}} \sum_{m=1}^{M} \pi_{n,m}\alpha_{n,m}\vec{e}_{m} - \vec{w}\right\|^{2}\right]. \tag{68}$$

Taking the expectation over the randomness in previous training rounds, conditioned on $\bar{\mathbf{x}}^{(r)}$, we obtain:

$$\mathbb{E}\left[F(\bar{\mathbf{x}}^{(r+1)}) - F(\bar{\mathbf{x}}^{(r)})\right] = \mathbb{E}\left[\mathbb{E}\left[F(\bar{\mathbf{x}}^{(r)} + \Delta \mathbf{x}) - F(\bar{\mathbf{x}}^{(r)})|\bar{\mathbf{x}}^{(r)}\right]\right] \\
\leq \left(\frac{4c^{2}}{3(1-2c^{2})} - \frac{1}{2}\right)T\eta \sum_{n=1}^{N} p_{n}q_{n}\mathbb{E}\|\nabla F(\bar{\mathbf{x}}^{(r)})\|^{2} + LT\eta^{2}\sigma^{2} \sum_{n=1}^{N} p_{n}^{2}q_{n} + \frac{2L^{2}T^{2}\eta^{3}\sigma^{2}}{1-2c^{2}} \sum_{n=1}^{N} p_{n}q_{n} \\
+ \frac{T\eta}{1-2c^{2}} \sum_{n=1}^{N} p_{n}q_{n} \left[\frac{\beta_{n}}{\alpha_{n}}(1-\gamma_{n}) + \gamma_{n} \cdot 2G\left(\frac{\beta_{n}'}{1-(1-\alpha_{n})^{2}} + 2\mathbb{E}\left[(1-\alpha_{n})^{2T\cdot\tau_{n}(r-1)}\right]\right] \\
+ \left\|\frac{1}{\alpha_{n}} \sum_{m=1}^{M} \pi_{n,m}\alpha_{n,m}\vec{e}_{m} - \vec{w}\right\|^{2}\right]. \tag{69}$$

We suppose $r \ge 2$ and denote $\tau_n(0) = 0$, we get:

$$\mathbb{E}\left[\left(1-\alpha_n\right)^{2T\cdot\tau_n(r-1)}\right] = \mathbb{E}\left[\mathbb{E}\left[\left(1-\alpha_n\right)^{2T\tau_n(r-1)}\middle|\tau_n(r-2)\right]\right]$$
(70)

$$= \mathbb{E}\left[(1 - p_n)(1 - \alpha_n)^{2T\tau_n(r-2)} + p_n(1 - \alpha_n)^{2T(\tau_n(r-2)+1)} \right]$$
(71)

$$= (1 - p_n + p_n (1 - \alpha_n)^{2T}) \mathbb{E}[(1 - \alpha_n)^{2T \tau_n (r - 2)}]$$
(72)

$$= (1 - p_n + p_n (1 - \alpha_n)^{2T})^{r-1} \mathbb{E}[(1 - \alpha_n)^{2T \tau_n(0)}]$$
(73)

$$= (1 - p_n + p_n (1 - \alpha_n)^{2T})^{r-1}. \tag{74}$$

Finally, we conclude that

$$\min_{r} \mathbb{E} \|\nabla F(\bar{\mathbf{x}}^{(r)})\|^{2} \le \frac{1}{R} \sum_{r=1}^{R} \|\nabla F(\bar{\mathbf{x}}^{(r)})\|^{2}$$

$$\leq \frac{1}{\left(\frac{1}{2} - \frac{4c^2}{3(1-2c^2)}\right) \sum_{n=1}^{N} p_n q_n} \left[\frac{1}{T\eta R} F(\bar{\mathbf{x}}^{(1)}) + L\eta \sigma^2 \sum_{n=1}^{N} p_n^2 q_n + \frac{2cL\eta \sigma^2}{1-2c^2} \sum_{n=1}^{N} p_n q_n + \frac{1}{1-2c^2} \sum_{n=1}^{N} p_n q_n \right] \\
\left(\frac{\beta_n}{\alpha_n} (1 - \gamma_n) + 2G\gamma_n \cdot \frac{\beta_n'}{1 - (1 - \alpha_n)^2} + 4G\gamma_n \cdot \frac{1}{R} \sum_{n=1}^{R} \left(1 - p_n + p_n (1 - \alpha_n)^{2T} \right)^{r-1} \right) \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n}{\alpha_n} (1 - \gamma_n) + 2G\gamma_n \cdot \frac{\beta_n'}{1 - (1 - \alpha_n)^2} + 4G\gamma_n \cdot \frac{1}{R} \sum_{n=1}^{R} \left(1 - p_n + p_n (1 - \alpha_n)^{2T} \right)^{r-1} \right) \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n}{\alpha_n} (1 - \gamma_n) + 2G\gamma_n \cdot \frac{\beta_n'}{1 - (1 - \alpha_n)^2} + 4G\gamma_n \cdot \frac{1}{R} \sum_{n=1}^{R} \left(1 - p_n + p_n (1 - \alpha_n)^{2T} \right)^{r-1} \right) \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n'}{1 - (1 - \alpha_n)^2} + 4G\gamma_n \cdot \frac{1}{R} \sum_{n=1}^{R} \left(1 - p_n + p_n (1 - \alpha_n)^{2T} \right)^{r-1} \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n'}{1 - (1 - \alpha_n)^2} + 4G\gamma_n \cdot \frac{1}{R} \sum_{n=1}^{R} \left(1 - p_n + p_n (1 - \alpha_n)^{2T} \right)^{r-1} \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n'}{1 - (1 - \alpha_n)^2} + 4G\gamma_n \cdot \frac{1}{R} \sum_{n=1}^{R} \left(1 - p_n + p_n (1 - \alpha_n)^{2T} \right)^{r-1} \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n'}{1 - (1 - \alpha_n)^2} + 4G\gamma_n \cdot \frac{1}{R} \sum_{n=1}^{R} \left(1 - p_n + p_n (1 - \alpha_n)^{2T} \right)^{r-1} \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n'}{1 - (1 - \alpha_n)^2} + 4G\gamma_n \cdot \frac{1}{R} \sum_{n=1}^{R} \left(1 - p_n + p_n (1 - \alpha_n)^{2T} \right)^{r-1} \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n'}{1 - (1 - \alpha_n)^2} + 4G\gamma_n \cdot \frac{1}{R} \sum_{n=1}^{R} \left(1 - p_n + p_n (1 - \alpha_n)^{2T} \right)^{r-1} \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n'}{1 - (1 - \alpha_n)^2} + 4G\gamma_n \cdot \frac{1}{R} \sum_{n=1}^{R} \left(\frac{\beta_n'}{1 - (1 - \alpha_n)^2} + \frac{\beta_n'}{1 - (1 - \alpha_n)^2} \right)^{r-1} \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n'}{1 - (1 - \alpha_n)^2} + \frac{\beta_n'}{1 - (1 - \alpha_n)^2} \right)^{r-1} \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n'}{1 - (1 - \alpha_n)^2} + \frac{\beta_n'}{1 - (1 - \alpha_n)^2} \right)^{r-1} \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n'}{1 - (1 - \alpha_n)^2} + \frac{\beta_n'}{1 - (1 - \alpha_n)^2} \right)^{r-1} \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n'}{1 - (1 - \alpha_n)^2} + \frac{\beta_n'}{1 - (1 - \alpha_n)^2} \right)^{r-1} \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n'}{1 - (1 - \alpha_n)^2} + \frac{\beta_n'}{1 - (1 - \alpha_n)^2} \right)^{r-1} \\
= \frac{1}{(1 - \alpha_n)^2} \left(\frac{\beta_n'} + \frac{\beta_n'}{1 - (1 - \alpha_n)^2} \right)^{r-1} \\
= \frac{1}{(1 - \alpha_n)^2} \left$$

$$+2G\gamma_n \sum_{m=1}^{M} \left(\frac{1}{\alpha} \pi_{n,m} \alpha_{n,m} - w_m\right)^2\right)$$

$$(75)$$

$$\leq \frac{1}{\left(\frac{1}{2} - \frac{4c^2}{3(1 - 2c^2)}\right) \sum_{n=1}^{N} p_n q_n} \left[\frac{1}{T\eta R} F(\bar{\mathbf{x}}^{(1)}) + L\eta \sigma^2 \sum_{n=1}^{N} p_n^2 q_n + \frac{2cL\eta \sigma^2}{1 - 2c^2} \sum_{n=1}^{N} p_n q_n + \frac{1}{1 - 2c^2} \sum_{n=1}^{N} p_n q_n \right]$$

$$\left(\frac{\beta_n}{\alpha_n}(1-\gamma_n) + 2G\gamma_n \cdot \frac{\beta_n'}{1-(1-\alpha_n)^2} + 4G\gamma_n \cdot \frac{1}{R} \cdot \frac{1}{p_n(1-(1-\alpha_n)^{2T})}\right)$$

$$+2G\gamma_n\sum_{m=1}^{M}\left(\frac{1}{\alpha_n}\pi_{n,m}\alpha_{n,m}-w_m\right)^2\right)\right] \tag{76}$$

$$\leq \frac{18}{\sum_{n=1}^{N} p_n q_n} \left[\frac{1}{T \eta R} F(\bar{\mathbf{x}}^{(1)}) + L \eta \sigma^2 \sum_{n=1}^{N} p_n^2 q_n + \frac{5L \eta \sigma^2}{3} \sum_{n=1}^{N} p_n q_n + \frac{5}{3} \sum_{n=1}^{N} p_n q_n \right]$$

$$\left(\frac{\beta_n}{\alpha_n}(1-\gamma_n) + 2G\gamma_n \cdot \frac{\beta_n'}{1-(1-\alpha_n)^2} + 4G\gamma_n \cdot \frac{1}{R} \cdot \frac{1}{p_n(1-(1-\alpha_n)^{2T})}\right)$$

$$+2G\gamma_n\sum_{m=1}^{M}\left(\frac{1}{\alpha_n}\pi_{n,m}\alpha_{n,m}-w_m\right)^2\right),\tag{77}$$

where

$$\alpha_n = \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m}, \quad \beta_n = \sum_{m=1}^{M} \pi_{n,m} \alpha_{n,m} d_m, \quad \gamma_n = \sum_{t=1}^{T} (1 - \alpha_n)^t,$$
 (78)

$$\beta_n' = 2\sum_{m=1}^M \pi_{n,m} \alpha_{n,m}^2 - \sum_{m=1}^M \pi_{n,m}^2 \alpha_{n,m}^2 + \alpha_n^2, \tag{79}$$

(77) holds when $c \le \min\{\frac{1}{2M}, \sqrt{\frac{1}{5}}\}$.

By choosing a sufficiently small learning rate $\eta = \frac{1}{\sqrt{R}}$, we derive the following convergence bound:

$$\min_{r} \mathbb{E} \|\nabla F(\bar{\mathbf{x}}^{(r)})\|^{2} \le \mathcal{O}(\frac{1}{\sqrt{R}}) + \mathcal{O}(\frac{1}{R}) + \epsilon, \tag{80}$$

where R is the number of global rounds, and ϵ is a non-vanishing residual term induced by the distribution shift from streaming data, i.e., the objective inconsistency (Wang & Ji, 2022).

The $\mathcal{O}(1/\sqrt{R})$ term is consistent with the convergence rate of several federated learning baselines under smooth non-convex assumptions, such as FedAdam (Reddi et al., 2020), FLASH (Panchal et al., 2023), and FedDisco (Ye et al., 2023).

F UPPER BOUND MINIMIZATION

Generally, a tighter bound corresponds to a better result. Thus, we optimize the upper bound in (77) from two aspects: the client-specific data sampling strategy $\{\alpha_{n,m}\}_{n\in\mathcal{N},m\in\mathcal{M}}$ and the server-side aggregation algorithm $\{p_n\}_{n\in\mathcal{N}}$. Directly solving the minimization of the upper bound results in a complicated expression. To simplify the expression, we treat α_n as a fixed hyperparameter that governs the average update ratio of the client buffer, referred to as the *sampling budget*.

F.1 DISTRIBUTION-GUIDED DATA SAMPLING STRATEGY

Our client-specific data sampling strategy is obtained through solving the following optimization problem:

$$\min_{\alpha_{n,m}} \frac{\beta_n}{\alpha_n} (1 - \gamma_n) + 2G\gamma_n \cdot \frac{\beta_n'}{1 - (1 - \alpha_n)^2} + 2G\gamma_n \sum_{m=1}^M \left(\frac{1}{\alpha_n} \pi_{n,m} \alpha_{n,m} - w_m\right)^2,$$
s.t.
$$\sum_{m=1}^M \pi_{n,m} \alpha_{n,m} = \alpha_n, \ \alpha_{n,m} \in [0, 1].$$
(81)

To solve this optimization problem, one condition of the optimal solution is that the derivative of the following function equals zero:

$$Q(\alpha_{n,m}) = \frac{\beta_n}{\alpha_n} (1 - \gamma_n) + 2G\gamma_n \cdot \frac{\beta'_n}{1 - (1 - \alpha_n)^2} + 2G\gamma_n \sum_{m=1}^{M} \left(\frac{1}{\alpha_n} \pi_{n,m} \alpha_{n,m} - w_m \right)^2$$
(82)

$$-\lambda_{1}\left(\sum_{m=1}^{M}\pi_{n,m}\alpha_{n,m} - \alpha_{n}\right) + \sum_{m=1}^{M}\mu_{m}\alpha_{n,m} - \nu_{m}\left(\sum_{m=1}^{M}\alpha_{n,m} - 1\right),\tag{83}$$

where λ_1 , μ_m and ν_m are multipliers. Then, we have

$$\frac{1 - \gamma_n}{\alpha_n} \pi_{n,m} d_m + \frac{2G\gamma_n}{1 - (1 - \alpha_n)^2} (4\pi_{n,m} - 2\pi_{n,m}^2) \alpha_{n,m} + \frac{4G\gamma_n \pi_{n,m}}{\alpha_n} \left(\frac{1}{\alpha_n} \pi_{n,m} \alpha_{n,m} - w_m\right) - \lambda_1 \pi_{n,m} + \mu_m - \nu_m = 0$$
(84)

To derive the sampling ratio $\alpha_{n,m}$, we only consider the active states with $\pi_{n,m} > 0$. Therefore, we obtain

$$\alpha_{n,m} \propto \frac{\lambda \alpha_n + \frac{\nu_m - \mu_m}{\pi_{n,m}} + 4G\gamma_n w_m - (1 - \gamma_n) d_m}{1 + \frac{1 - \alpha_n}{\alpha_n} \pi_{n,m}}$$
(85)

To simplify the solution and derive a closed-form, interpretable expression, we relax the KKT conditions by eliminating the dual variables, setting $\mu_m=0$ and $\nu_m=0$, under the assumption that the optimal values of $\alpha_{n,m}$ lie strictly within the open interval (0,1). This assumption is practically reasonable. In realistic streaming environments, it is neither desirable nor feasible for a client to completely discard previously stored data $(\alpha_{n,m}=1)$ or to entirely ignore new samples from a given state $(\alpha_{n,m}=0)$. Either extreme leads to inefficient storage usage and poor adaptability to evolving data distributions. Under this relaxed setting, we define the following score function:

$$score_n(m) = \frac{w_m - a_1 d_m + b_1}{1 + \frac{1 - \alpha_n}{\alpha_n} \pi_{n,m}},$$
(86)

where $a_1 = \frac{1-\gamma_n}{4G\gamma_n}$ and $b_1 = \frac{\lambda_1}{4G\gamma_n}$ are tunable constants. Intuitively, the score increases with the weight in the global objective function w_m , and decreases with the state-specific heterogeneity bound d_m or the probability $\pi_{n,m}$ of client n being in state m.

To ensure feasibility and numerical stability, we apply a ReLU operation to filter out negative scores, and normalize the resulting values with respect to the state distribution π_n :

$$\alpha_{n,m} = \frac{\alpha_n \cdot \text{ReLu}(\text{score}_n(m))}{\sum_{m'=1}^{M} \pi_{n,m'} \cdot \text{ReLu}(\text{score}_n(m'))}$$
(87)

F.2 AGGREGATION WEIGHT DETERMINATION

We now describe how to determine the client aggregation weights $\{p_n\}_{n=1}^N$ by minimizing the dominant terms in the convergence upper bound (77), discarding constant factors irrelevant to p_n . Specifically, we formulate the following constrained optimization:

$$\min_{p_n} \frac{1}{\sum_{n=1}^{N} p_n q_n} \left[L \eta \sigma^2 \sum_{n=1}^{N} p_n^2 q_n + \frac{5}{3} \sum_{n=1}^{N} p_n q_n \left(L \eta \sigma^2 + \frac{\beta_n}{\alpha_n} (1 - \gamma_n) \right) + 2G \gamma_n \cdot \frac{\beta_n'}{1 - (1 - \alpha_n)^2} + 2G \gamma_n \sum_{m=1}^{M} \left(\frac{1}{\alpha_n} \pi_{n,m} \alpha_{n,m} - w_m \right)^2 \right]$$
s.t.
$$\sum_{n=1}^{N} p_n = 1, \ p_n \ge 0. \tag{88}$$

To simplify the expression, we denote

$$s_n = L\eta\sigma^2 + \frac{\beta_n}{\alpha_n}(1 - \gamma_n) + 2G\gamma_n \cdot \frac{\beta_n'}{1 - (1 - \alpha_n)^2} + 2G\gamma_n \sum_{n=1}^{M} \left(\frac{1}{\alpha_n}\pi_{n,m}\alpha_{n,m} - w_m\right)^2$$
 (89)

$$T_0 = \sum_{n=1}^{N} p_n q_n, \ T_1 = L \eta \sigma^2 \sum_{n=1}^{N} p_n^2 q_n, \ T_2 = \frac{5}{3} \sum_{n=1}^{N} p_n q_n s_n.$$
 (90)

We apply a standard surrogate objective (Ye et al., 2023), transforming from minimizing $(T_1+T_2)/T_0$ to minimizing $T_1+T_2-\lambda_2T_0$, where $\lambda_2>0$ is a balancing hyper-parameter. Therefore, the optimization is

$$\min_{p_n} L\eta \sigma^2 \sum_{n=1}^{N} p_n^2 q_n + \frac{5}{3} \sum_{n=1}^{N} p_n q_n s_n - \lambda_2 \sum_{n=1}^{N} p_n q_n$$
s.t.
$$\sum_{n=1}^{N} p_n = 1, \ p_n \ge 0,$$
(91)

Let ν and μ_n be the Lagrangian multipliers associated with the equality and inequality constraints, respectively. The Lagrangian is:

$$L(p_n) = L\eta\sigma^2 \sum_{n=1}^{N} p_n^2 q_n + \frac{5}{3} \sum_{n=1}^{N} p_n q_n s_n - \lambda_2 \sum_{n=1}^{N} p_n q_n + \sum_{n=1}^{N} \mu_n p_n - \nu \left(\sum_{n=1}^{N} p_n - 1\right).$$
 (92)

Taking the derivative of $L(\cdot)$ with respect to p_n and setting it to zero yields:

$$2L\eta\sigma^2 p_n q_n + \frac{5}{3}q_n s_n - \lambda_2 q_n + \mu_n - \nu = 0$$
(93)

Solving this, we obtain:

$$p_n = \frac{1}{2L\eta\sigma^2} \left(\lambda - \frac{5}{3}s_n + \frac{\nu - \mu_n}{q_n}\right). \tag{94}$$

If we relax the KKT multipliers and absorb them into constants a_2 and b_2 , we arrive at the interpretable approximation:

$$p_n \propto \frac{1}{q_n} - a_2 \cdot s_n + b_2,\tag{95}$$

where a_2 and b_2 are tunable constants. Intuitively, s_n is a metric to measure the heterogeneity of client-specific data distribution. It takes into account a client-specific data sampling strategy. Clients with smaller s_n are assigned higher aggregation weights. The term $1/q_n$ ensures fairness with respect to availability.

THE SFEDPO ALGORITHM

1512

1513 1514

1515

1516 1517

1518

1519

1520

1521

1522 1523

1524

1525

1526

1527

1540

1541 1542 1543

1545

1546 1547 1548

1561

1564 1565 This section provides the complete pseudocode of our proposed SFedPO framework, as described in Section 5. Algorithm 1 outlines the end-to-end federated optimization process in the streaming setting with a prediction oracle and dynamic data sampling.

We observe that although Eq. (11) preserves the proportionality between sampling ratios and utility scores, it does not inherently ensure that $\hat{\alpha}_{n,m} \in [0,1]$. To resolve this, Algorithm 2 employs a budget-aware projection strategy: it iteratively clips any $\hat{\alpha}_{n,m}$ exceeding 1 and adjusts the remaining budget allocation across the other states. This iterative procedure guarantees that all sampling ratios remain feasible and that the total allocation strictly satisfies the sampling budget constraint α_n .

Algorithm 1: SFedPO: Streaming Federated Learning with Prediction Oracle

Require: Initial global model $\bar{\mathbf{x}}^{(1)}$, learning rate η , number of rounds R, local steps T, prediction oracle \mathcal{O} , sampling budget α_n

- 1: **for** each round r = 1 to R **do**
- Server selects a client set S_r
 - 3: for each selected client $n \in S_r$ in parallel do
- 4: $\hat{\boldsymbol{\pi}}_n \leftarrow \mathcal{O}(n,r)$

// predicted state distribution

- 1530 5: Compute $\hat{\alpha}_{n,m}$ using Algorithm 2.
- Initialize $\mathbf{x}_{n,0}^{(r)} \leftarrow \bar{\mathbf{x}}^{(r)}$ 1531 6:
- 1532 for each local step t = 1 to T do 7:
- Receive new data distribution $\{\mathcal{D}_m\}_{m=1}^M$ from streaming source. 1533 8:
- Update local dataset with a ratio of $\alpha_{n,m_{n,t}^{(r)}}$ 1534 9: 1535
 - Sample minibatch $\xi_{n,t} \sim \mathcal{D}_{n,t}^{(r)}$, and train the model: $\mathbf{x}_{n,t}^{(r)} = \mathbf{x}_{n,t-1}^{(r)} \eta \cdot \mathbf{g}_{n,t}^{(r)}$ 10:
- 1536 11: 1537
 - Send update $\Delta \mathbf{x}_n^{(r)} = \mathbf{x}_{n,E}^{(r)} \bar{\mathbf{x}}^{(r)}$ to server 12:
- 1538 13: 1539
 - 14: Server computes aggregation weights \hat{p}_n using:

$$\hat{p}_n = \frac{\text{ReLU}\left(\frac{1}{q_n} - a_2 \cdot \hat{s}_n + b_2\right)}{\sum_{j \in S_r} \text{ReLU}\left(\frac{1}{q_j} - a_2 \cdot \hat{s}_j + b_2\right)}$$

15: Update global model:

$$\bar{\mathbf{x}}^{(r+1)} = \bar{\mathbf{x}}^{(r)} + \eta \sum_{n \in S_r} \hat{p}_n \cdot \Delta \mathbf{x}_n^{(r)}$$

16: **end for**

```
Algorithm 2: Computation of \hat{\alpha}_{n,m}
```

```
1582
                 Require: Number of clients N and states M, sampling budget \alpha_n, estimated state distribution \hat{\pi}_{n,m},
1583
1584
                          scôre_n(m)
                    1: Initialize remaining set \mathcal{R} \leftarrow \{m \mid \hat{\pi}_{n,m} > 0\}, residual budget \tilde{\alpha}_n \leftarrow \alpha_n
1585
                    2: \ Set \ \texttt{found\_one} \leftarrow \textbf{True}
1586
                    3: while found_one and \mathcal{R} \neq \emptyset do
1587
                              \begin{array}{l} \text{found\_one} \leftarrow \textbf{False} \\ S \leftarrow \sum_{m \in \mathcal{R}} \hat{\pi}_{n,m} \cdot \text{ReLU}(\widehat{\text{score}}_n(m)) \\ \textbf{if } S = 0 \textbf{ then} \end{array}
1588
                    5:
1589
                    6:
1590
                    7:
                                    for each m \in \mathcal{R} do
1591
                                         \hat{\alpha}_{n,m} \leftarrow \frac{\tilde{\alpha}_n}{\sum_{m \in \mathcal{R}} \hat{\pi}_{n,m}}
                    8:
1592
                    9:
                                    end for
1593
                  10:
                                    break
1594
                  11:
                               else
                                    \begin{split} & \textbf{for } \operatorname{each} m \in \mathcal{R} \ \textbf{do} \\ & \hat{\alpha}_{n,m} \leftarrow \min \left( \frac{\tilde{\alpha}_n \cdot \operatorname{ReLU}(\widehat{\operatorname{score}}_n(m))}{S}, 1 \right) \\ & \textbf{if } \hat{\alpha}_{n,m} = 1 \ \textbf{then} \end{split}
1595
                  12:
1596
                 13:
1597
                  14:
1598
                  15:
                                              \tilde{\alpha}_n \leftarrow \tilde{\alpha}_n - \hat{\pi}_{n,m}
1599
                  16:
                                              remove m from \mathcal{R}, set found_one \leftarrow True
1600
                  17:
1601
                  18:
                                         end if
1602
                 19:
                                    end for
                 20:
                               end if
1603
                 21: end while
1604
```

G PROOF OF THEOREM 2

Based on the convergence analysis in Appendix E, we develop a theoretical strategy under the ground-truth state distribution $\pi_n = (\pi_{n,1}, \dots, \pi_{n,M})$ for each client $n \in \mathcal{N}$. In practical deployments, however, this is not directly observable and must be approximated using a prediction oracle, which provides an estimated state distribution $\hat{\pi}_n = (\hat{\pi}_{n,1}, \dots, \hat{\pi}_{n,M})$. This naturally raises the question:

How does the prediction error affect the convergence behavior of the overall algorithm?

To address this question, we conduct a robustness analysis that quantifies the impact of prediction error on the convergence guarantee derived in Appendix E. Specifically, we define the prediction error for client n as $\delta_n := \|\hat{\pi}_n - \pi_n\|_1$, and study how this error propagates into the convergence bound through the data sampling strategy α and aggregation weights $\{p_n\}_{n\in\mathcal{N}}$, both of which are functions of π_n in theory but implemented based on $\hat{\pi}_n$ in practice.

Recall the convergence bound in Theorem 1, we focus on the following term:

$$\mathcal{B}(\alpha_{n,m}, p_n) := \frac{1}{\sum_{n=1}^{N} p_n q_n} \left[L \eta \sigma^2 \sum_{n=1}^{N} p_n^2 q_n + \frac{5}{3} \sum_{n=1}^{N} p_n q_n s_n \right]$$
(96)

where

$$s_n = L\eta\sigma^2 + \frac{\beta_n}{\alpha_n}(1 - \gamma_n) + \frac{2G\gamma_n\beta_n'}{1 - (1 - \alpha_n)^2} + 2G\gamma_n\sum_{m=1}^M \left(\frac{1}{\alpha_n}\pi_{n,m}\alpha_{n,m} - w_m\right)^2.$$
 (97)

In the realistic setting, these ratios are implemented as $\hat{\alpha}_{n,m}$ using the estimated distribution $\hat{\pi}_n$. Following the relaxed design in Section 5, the sampling ratios are determined by:

$$\widehat{\mathrm{score}}_n(m) = \frac{w_m - a_1 d_m + b_1}{1 + \frac{\alpha_n}{1 - \alpha_n} \widehat{\pi}_{n,m}}, \text{ and } \widehat{\alpha}_{n,m} = \frac{\mathrm{ReLU}(\widehat{\mathrm{score}}_n(m))}{\sum_{m'=1}^M \widehat{\pi}_{n,m'} \cdot \mathrm{ReLU}(\widehat{\mathrm{score}}_n(m'))}.$$

Our goal is to analyze the difference between the convergence bound under the theoretical strategy $(\alpha_{n,m},p_n)$ and the practical strategy $(\hat{\alpha}_{n,m},\hat{p}_n)$ computed using prediction oracle, and to establish a robustness bound in terms of the prediction error δ_n .

G.1 IMPACT ON SAMPLING STRATEGY.

We now analyze how the prediction error $\delta_n = \|\hat{\pi}_n - \pi_n\|_1$ affects the resulting data sampling strategy. Let $\alpha_n = \{\alpha_{n,m}\}_{m=1}^M$ and $\hat{\alpha}_n = \{\hat{\alpha}_{n,m}\}_{m=1}^M$ be the sampling ratios computed using the true and estimated state distributions, respectively. Our goal is to bound the ℓ_1 deviation between the two vectors, i.e., $\|\hat{\alpha}_n - \alpha_n\|_1$.

Let $S_n := \sum_{j=1}^M \pi_{n,j} \cdot \operatorname{ReLU}(\operatorname{score}_n(j))$ and $\hat{S}_n := \sum_{j=1}^M \hat{\pi}_{n,j} \cdot \operatorname{ReLU}(\widehat{\operatorname{score}}_n(j))$. We suppose that the score function is L_s -Lipschitz in $\pi_{n,m}$, i.e.,

$$|\widehat{\text{score}}_n(m) - \text{score}_n(m)| \le L_s |\widehat{\pi}_{n,m} - \pi_{n,m}|. \tag{98}$$

We can derive the following bound:

$$|\hat{\alpha}_{n,m} - \alpha_{n,m}| \le \left| \frac{\text{ReLU}(\widehat{\text{score}}_{n}(m))}{\hat{S}_{n}} - \frac{\text{ReLU}(\text{score}_{n}(m))}{\hat{S}_{n}} \right| + \left| \frac{\text{ReLU}(\text{score}_{n}(m))}{\hat{S}_{n}} - \frac{\text{ReLU}(\text{score}_{n}(m))}{\hat{S}_{n}} \right|$$

$$\le \frac{L_{s}|\hat{\pi}_{n,m} - \pi_{n,m}|}{\hat{S}_{n}} + \text{ReLU}(\text{score}_{n}(m)) \cdot \left| \frac{1}{\hat{S}_{n}} - \frac{1}{S_{n}} \right|.$$
(99)

We further estimate the deviation between the \hat{S}_n and S_n :

$$|\hat{S}_{n} - S_{n}| = \left| \sum_{j=1}^{M} \hat{\pi}_{n,j} \cdot \text{ReLU}(\widehat{\text{score}}_{n}(j)) - \sum_{j=1}^{M} \pi_{n,j} \cdot \text{ReLU}(\text{score}_{n}(j)) \right|$$

$$\leq \sum_{j=1}^{M} \left| \hat{\pi}_{n,j} \cdot \text{ReLU}(\widehat{\text{score}}_{n}(j)) - \pi_{n,j} \cdot \text{ReLU}(\text{score}_{n}(j)) \right|$$

$$\leq \sum_{j=1}^{M} \left(|\hat{\pi}_{n,j} - \pi_{n,j}| \cdot \text{ReLU}(\text{score}_{n}(j)) + \hat{\pi}_{n,j} \cdot |\text{ReLU}(\widehat{\text{score}}_{n}(j)) - \text{ReLU}(\text{score}_{n}(j)) \right|$$

$$\leq \delta_{n} \cdot \|\text{ReLU}(\text{score}_{n})\|_{\infty} + L_{s} \cdot \delta_{n}.$$
(100)

Thus, we obtain

$$\left| \frac{1}{\hat{S}_n} - \frac{1}{S_n} \right| = \frac{|\hat{S}_n - S_n|}{\hat{S}_n \cdot S_n} \le \frac{\delta_n \cdot (L_s + \|\text{ReLU}(\text{score}_n)\|_{\infty})}{\hat{S}_n \cdot S_n}, \tag{101}$$

and hence,

$$|\hat{\alpha}_{n,m} - \alpha_{n,m}| \le \frac{L_s |\hat{\pi}_{n,m} - \pi_{n,m}|}{\hat{S}_n} + \text{ReLU}(\text{score}_n(m)) \cdot \frac{\delta_n \cdot (L_s + ||\text{ReLU}(\text{score}_n)||_{\infty})}{\hat{S}_n \cdot S_n}.$$
(102)

Summing over all $m \in \mathcal{M}$, and using the fact that ReLU scores are bounded and $S_n, \hat{S}_n \geq \epsilon$ for some small constant, we obtain:

$$\|\hat{\boldsymbol{\alpha}}_n - \boldsymbol{\alpha}_n\|_1 \le C_1 \cdot \delta_n,\tag{103}$$

where C_1 depends on the Lipschitz constant L_s , the upper bound of the ReLU scores, and the lower bounds of the S_n and \hat{S}_n .

G.2 IMPACT ON AGGREGATION WEIGHTS.

Recall the client-specific heterogeneity score:

$$s_n = L\eta \sigma^2 + \frac{\beta_n}{\alpha_n} (1 - \gamma_n) + \frac{2G\gamma_n \beta_n'}{1 - (1 - \alpha_n)^2} + 2G\gamma_n \sum_{m=1}^M \left(\frac{1}{\alpha_n} \pi_{n,m} \alpha_{n,m} - w_m \right)^2.$$

The first term is a constant and unaffected by prediction error. We now analyze the impact of prediction error $\delta_n := \|\hat{\pi}_n - \pi_n\|_1$ on the remaining three terms.

Bounding $\frac{\beta_n}{\alpha_n}(1-\gamma_n)$. Let

$$\beta_n = \sum_{m=1}^M \pi_{n,m} \alpha_{n,m} d_m, \quad \hat{\beta}_n = \sum_{m=1}^M \hat{\pi}_{n,m} \hat{\alpha}_{n,m} d_m.$$

Then, using triangle inequality, we have

$$|\hat{\beta}_{n} - \beta_{n}| \leq \sum_{m=1}^{M} |\hat{\pi}_{n,m} \hat{\alpha}_{n,m} - \pi_{n,m} \alpha_{n,m}| \cdot |d_{m}|$$

$$\leq \sum_{m=1}^{M} (|\hat{\pi}_{n,m} - \pi_{n,m}| \cdot |\hat{\alpha}_{n,m}| + \pi_{n,m} \cdot |\hat{\alpha}_{n,m} - \alpha_{n,m}|) d_{m}.$$

$$\leq D_{\max} \cdot (||\hat{\alpha}_{n} - \alpha_{n}||_{1} + \delta_{n}), \qquad (104)$$

where $D_{\max} := \max_m d_m$. This gives:

$$\left| \frac{\hat{\beta}_n}{\alpha_n} (1 - \gamma_n) - \frac{\beta_n}{\alpha_n} (1 - \gamma_n) \right| \le \frac{(1 - \gamma_n) D_{\text{max}}}{\alpha_n} \cdot (C_1 + 1) \cdot \delta_n =: C_2 \delta_n.$$
 (105)

Bounding $\frac{2G\gamma_n\beta'_n}{1-(1-\alpha_n)^2}$. Recall

$$\beta'_n = 2\sum_{m=1}^M \pi_{n,m}\alpha_{n,m}^2 - \sum_{m=1}^M \pi_{n,m}^2 \alpha_{n,m}^2 + \alpha_n^2.$$

Let $\hat{\beta}'_n$ be the corresponding quantity computed under $\hat{\pi}_{n,m}$ and $\hat{\alpha}_{n,m}$:

$$\hat{\beta}'_n = 2\sum_{m=1}^M \hat{\pi}_{n,m} \hat{\alpha}_{n,m}^2 - \sum_{m=1}^M \hat{\pi}_{n,m}^2 \hat{\alpha}_{n,m}^2 + \alpha_n^2$$

Note that α_n is assumed to be fixed, so it cancels in the subtraction. Then we have:

$$|\hat{\beta}_n' - \beta_n'| \le 2 \sum_{m=1}^M |\hat{\pi}_{n,m} \hat{\alpha}_{n,m}^2 - \pi_{n,m} \alpha_{n,m}^2| + \sum_{m=1}^M |\hat{\pi}_{n,m}^2 \hat{\alpha}_{n,m}^2 - \pi_{n,m}^2 \alpha_{n,m}^2|.$$
 (106)

We now bound each term separately. First, observe that

$$\left| \hat{\pi}_{n,m} \hat{\alpha}_{n,m}^{2} - \pi_{n,m} \alpha_{n,m}^{2} \right| \leq \left| \hat{\pi}_{n,m} (\hat{\alpha}_{n,m}^{2} - \alpha_{n,m}^{2}) \right| + \left| (\hat{\pi}_{n,m} - \pi_{n,m}) \alpha_{n,m}^{2} \right|$$

$$\leq \hat{\pi}_{n,m} \cdot 2\alpha_{\max} |\hat{\alpha}_{n,m} - \alpha_{n,m}| + \alpha_{\max}^{2} |\hat{\pi}_{n,m} - \pi_{n,m}|,$$
(107)

where we assume that $\alpha_{n,m}$, $\hat{\alpha}_{n,m} \leq \alpha_{\max} \leq 1$.

Similarly, for the second term:

$$\left| \hat{\pi}_{n,m}^{2} \hat{\alpha}_{n,m}^{2} - \pi_{n,m}^{2} \alpha_{n,m}^{2} \right| \leq \left| \hat{\pi}_{n,m}^{2} (\hat{\alpha}_{n,m}^{2} - \alpha_{n,m}^{2}) \right| + \left| (\hat{\pi}_{n,m}^{2} - \pi_{n,m}^{2}) \alpha_{n,m}^{2} \right|$$

$$\leq \hat{\pi}_{n,m}^{2} \cdot 2\alpha_{\max} |\hat{\alpha}_{n,m} - \alpha_{n,m}| + 2\pi_{\max} \alpha_{\max}^{2} |\hat{\pi}_{n,m} - \pi_{n,m}|,$$
(108)

where we assume that $\pi_{n,m}$, $\hat{\pi}_{n,m} \leq \pi_{\max} \leq 1$.

Summing over all m, and applying the triangle inequality:

$$|\hat{\beta}'_{n} - \beta'_{n}| \leq 4\pi_{\max}\alpha_{\max} \|\hat{\alpha}_{n} - \alpha_{n}\|_{1} + 2\alpha_{\max}^{2} \|\hat{\pi}_{n} - \pi_{n}\|_{1}$$

$$+ 2\pi_{\max}^{2}\alpha_{\max} \|\hat{\alpha}_{n} - \alpha_{n}\|_{1} + 2\pi_{\max}\alpha_{\max}^{2} \|\hat{\pi}_{n} - \pi_{n}\|_{1}$$

$$\leq (4\pi_{\max}\alpha_{\max} + 2\pi_{\max}^{2}\alpha_{\max})C_{1}\delta_{n} + (2\alpha_{\max}^{2} + 2\pi_{\max}\alpha_{\max}^{2})\delta_{n}.$$
(110)

Hence, we conclude:

$$\left| \frac{2G\gamma_{n}\hat{\beta}'_{n}}{1 - (1 - \alpha_{n})^{2}} - \frac{2G\gamma_{n}\beta'_{n}}{1 - (1 - \alpha_{n})^{2}} \right| \\
\leq \frac{2G\gamma_{n}}{1 - (1 - \alpha_{n})^{2}} \left[(4\pi_{\max}\alpha_{\max} + 2\pi_{\max}^{2}\alpha_{\max})C_{1} + (2\alpha_{\max}^{2} + 2\pi_{\max}\alpha_{\max}^{2}) \right] \cdot \delta_{n} =: C_{3} \cdot \delta_{n}. \tag{111}$$

Bounding
$$2G\gamma_n \sum_{m=1}^M \left(\frac{1}{\alpha_n} \pi_{n,m} \alpha_{n,m} - w_m\right)^2$$
. Let
$$z_m = \frac{1}{\alpha_n} \pi_{n,m} \alpha_{n,m}, \quad \hat{z}_m = \frac{1}{\alpha_n} \hat{\pi}_{n,m} \hat{\alpha}_{n,m}.$$

Then, we have

$$\left| \hat{z}_m - z_m \right| = \left| \frac{1}{\alpha_n} (\hat{\pi}_{n,m} \hat{\alpha}_{n,m} - \pi_{n,m} \alpha_{n,m}) \right| \le \frac{1}{\alpha_n} (\left| \hat{\pi}_{n,m} - \pi_{n,m} \right| \cdot \hat{\alpha}_{n,m} + \pi_{n,m} \cdot \left| \hat{\alpha}_{n,m} - \alpha_{n,m} \right|)$$

$$\le \frac{1}{\alpha_n} (\alpha_{\text{max}} |\hat{\pi}_{n,m} - \pi_{n,m}| + \pi_{\text{max}} |\hat{\alpha}_{n,m} - \alpha_{n,m}|). \tag{112}$$

Summing over all m, and applying the triangle inequality:

$$\left| \sum_{m=1}^{M} (\hat{z}_{m} - w_{m})^{2} - \sum_{m=1}^{M} (z_{m} - w_{m})^{2} \right| \leq \sum_{m=1}^{M} |\hat{z}_{m}^{2} - z_{m}^{2}| + 2 \sum_{m=1}^{M} |w_{m}| |\hat{z}_{m} - z_{m}|$$

$$\leq \sum_{m=1}^{M} |\hat{z}_{m} - z_{m}| \cdot |\hat{z}_{m} + z_{m}| + 2 \sum_{m=1}^{M} |w_{m}| \cdot |\hat{z}_{m} - z_{m}|$$

$$\leq \sum_{m=1}^{M} \frac{4}{\alpha_{n}} (\alpha_{\max} |\hat{\pi}_{n,m} - \pi_{n,m}| + \pi_{\max} |\hat{\alpha}_{n,m} - \alpha_{n,m}|)$$

$$\leq \frac{4(\alpha_{\max} + C_{1}\pi_{\max})}{\alpha_{n}} \cdot \delta_{n}, \tag{113}$$

where we use the fact that $\hat{z}_m, z_m \in [0, 1]$ and $w_m \leq 1$, for all m.

Therefore, we have:

$$\left| 2G\gamma_n \sum_{m=1}^{M} (\hat{z}_m - w_m)^2 - 2G\gamma_n \sum_{m=1}^{M} (z_m - w_m)^2 \right| \le 2G\gamma_n \cdot \frac{4(\alpha_{\max} + C_1 \pi_{\max})}{\alpha_n} \cdot \delta_n =: C_4 \cdot \delta_n.$$
(114)

Combining bounds (105), (111), and (114), we obtain the total deviation of s_n :

$$|\hat{s}_n - s_n| \le (C_2 + C_3 + C_4) \cdot \delta_n =: C_s \cdot \delta_n,$$
 (115)

We now analyze how the prediction error δ_n affects the final aggregation weights. Given the aggregation weights:

$$p_n = \frac{\text{ReLU}(\frac{1}{q_n} - a_2 \cdot s_n + b_2)}{\sum_{i=1}^{N} \text{ReLU}(\frac{1}{q_i} - a_2 \cdot s_i + b_2)}, \quad \hat{p}_n = \frac{\text{ReLU}(\frac{1}{q_n} - a_2 \cdot \hat{s}_n + b_2)}{\sum_{i=1}^{N} \text{ReLU}(\frac{1}{q_i} - a_2 \cdot \hat{s}_i + b_2)}.$$

Let us denote:

$$\psi_n := \text{ReLU}(\frac{1}{q_n} - a_2 s_n + b_2), \quad \hat{\psi}_n := \text{ReLU}(\frac{1}{q_n} - a_2 \hat{s}_n + b_2),$$

$$Z := \sum_{i=1}^{N} \psi_i, \quad \hat{Z} := \sum_{i=1}^{N} \hat{\psi}_i.$$

Then, the deviation between \hat{p}_n and p_n can be bounded by:

$$|\hat{p}_{n} - p_{n}| = \left| \frac{\hat{\psi}_{n}}{\hat{Z}} - \frac{\psi_{n}}{Z} \right|$$

$$\leq \left| \frac{\hat{\psi}_{n} - \psi_{n}}{\hat{Z}} \right| + \left| \psi_{n} \cdot \left(\frac{1}{\hat{Z}} - \frac{1}{Z} \right) \right|$$

$$\leq \frac{a_{2}|\hat{s}_{n} - s_{n}|}{\hat{Z}} + \psi_{n} \cdot \frac{|Z - \hat{Z}|}{Z\hat{Z}}, \tag{116}$$

where we used the fact that ReLU is 1-Lipschitz.

Note that $|Z - \hat{Z}| \leq \sum_{i=1}^{N} |\psi_i - \hat{\psi}_i| \leq a_2 \sum_{i=1}^{N} |s_i - \hat{s}_i| \leq a_2 C_s \cdot \sum_{i=1}^{N} \delta_i$.

Thus,

$$|\hat{p}_n - p_n| \le \frac{a_2 C_s \delta_n}{\hat{Z}} + \frac{\psi_n \cdot a_2 C_s \sum_{i=1}^N \delta_i}{Z \hat{Z}},$$
 (117)

where C_s is a bound such that $|s_n - \hat{s}_n| \le C_s \delta_n$ as shown in (115).

Therefore, under mild conditions on $\psi_n \ge \epsilon$ and $Z, \hat{Z} \ge N\epsilon$ for some $\epsilon > 0$, we conclude:

$$|\hat{p}_n - p_n| \le C_5 \cdot \delta_n + C_6 \cdot \sum_{i=1}^N \delta_i, \tag{118}$$

for constants C_5, C_6 depending on a_2, C_s , and bounds on ψ_n and Z.

G.3 IMPACT ON CONVERGENCE BOUND.

We now analyze how the prediction error $\delta_n = \|\hat{\pi}_n - \pi_n\|_1$ propagates to the convergence upper bound through sampling strategy $\{\alpha_{n,m}\}$ and aggregation weights $\{\hat{p}_n\}$.

Recall

$$\mathcal{B}(\alpha_{n,m}, p_n) = \frac{1}{\sum_{n=1}^{N} p_n q_n} \left[L \eta \sigma^2 \sum_{n=1}^{N} p_n^2 q_n + \frac{5}{3} \sum_{n=1}^{N} p_n q_n s_n \right],$$

$$\mathcal{B}(\hat{\alpha}_{n,m}, \hat{p}_n) = \frac{1}{\sum_{n=1}^{N} \hat{p}_n q_n} \left[L \eta \sigma^2 \sum_{n=1}^{N} \hat{p}_n^2 q_n + \frac{5}{3} \sum_{n=1}^{N} \hat{p}_n q_n \tilde{s}_n \right],$$

where

$$\tilde{s}_n = L\eta\sigma^2 + \frac{\tilde{\beta}_n}{\alpha_n}(1 - \gamma_n) + \frac{2G\gamma_n\tilde{\beta}_n'}{1 - (1 - \alpha_n)^2} + 2G\gamma_n\sum_{m=1}^M \left(\frac{1}{\alpha_n}\pi_{n,m}\alpha_{n,m} - w_m\right)^2,\tag{119}$$

$$\tilde{\beta}_n = \sum_{m=1}^M \pi_{n,m} \hat{\alpha}_{n,m} d_m \tag{120}$$

$$\tilde{\beta}'_{n} = 2\sum_{m=1}^{M} \pi_{n,m} \hat{\alpha}_{n,m}^{2} - \sum_{m=1}^{M} \pi_{n,m}^{2} \hat{\alpha}_{n,m}^{2} + \alpha_{n}^{2}.$$
(121)

Then, the deviation can be bounded as:

$$|\mathcal{B}(\hat{\alpha}_{n,m},\hat{p}_n) - \mathcal{B}(\alpha_{n,m},p_n)| \le \left| \frac{1}{\sum_{n=1}^{N} \hat{p}_n q_n} - \frac{1}{\sum_{n=1}^{N} p_n q_n} \right| \cdot \mathcal{B}_{\max} + \frac{1}{\sum_{n=1}^{N} p_n q_n} \cdot \Delta_{\text{num}}, \quad (122)$$

where $\mathcal{B}_{\mathrm{max}}$ denotes a uniform upper bound on the numerator, and

$$\Delta_{\text{num}} = L\eta\sigma^2 \cdot \left| \sum_{n=1}^{N} (\hat{p}_n^2 - p_n^2) q_n \right| + \frac{5}{3} \cdot \left| \sum_{n=1}^{N} \hat{p}_n q_n \tilde{s}_n - \sum_{n=1}^{N} p_n q_n s_n \right|.$$
 (123)

We now bound each part of Δ_{num} :

Bounding
$$\sum_{n=1}^{N} (\hat{p}_n^2 - p_n^2) q_n$$
. Let $\delta_n^p := \hat{p}_n - p_n$. Then

$$\left| \sum_{n=1}^{N} (\hat{p}_{n}^{2} - p_{n}^{2}) q_{n} \right| = \left| \sum_{n=1}^{N} \left[(p_{n} + \delta_{n}^{p})^{2} - p_{n}^{2} \right] q_{n} \right| = \left| \sum_{n=1}^{N} \left[2p_{n} \delta_{n}^{p} + (\delta_{n}^{p})^{2} \right] q_{n} \right|$$

$$\leq 2 \sum_{n=1}^{N} p_{n} q_{n} |\delta_{n}^{p}| + \sum_{n=1}^{N} q_{n} (\delta_{n}^{p})^{2}$$

$$\leq 2C_{5} \sum_{n=1}^{N} p_{n} q_{n} \delta_{n} + 2C_{6} \left(\sum_{n=1}^{N} p_{n} q_{n} \right) \cdot \left(\sum_{n=1}^{N} \delta_{n} \right) + \mathcal{O}(\delta_{n}^{2}). \tag{124}$$

Bounding $\sum_{n=1}^{N} \hat{p}_n q_n \tilde{s}_n - \sum_{n=1}^{N} p_n q_n s_n$. We fist bound $|\tilde{s}_n - s_n|$. Similar to calculating the bound of $|\hat{s}_n - s_n|$, we have

$$|\tilde{\beta}_n - \beta_n| \le \sum_{m=1}^M \pi_{n,m} |\hat{\alpha}_{n,m} - \alpha_{n,m}| \cdot |d_m| \le D_{\max} \cdot ||\hat{\alpha}_n - \alpha_n||_1, \tag{125}$$

$$|\tilde{\beta}'_{n} - \beta'_{n}| \leq 2 \sum_{m=1}^{M} \pi_{n,m} |\hat{\alpha}_{n,m}^{2} - \alpha_{n,m}^{2}| + \sum_{m=1}^{M} \pi_{n,m}^{2} |\hat{\alpha}_{n,m}^{2} - \alpha_{n,m}^{2}|$$

$$\leq (4\pi_{\max}\alpha_{\max} + 2\pi_{\max}^{2}\alpha_{\max}) ||\hat{\alpha}_{n} - \alpha_{n}||_{1}.$$
(126)

Therefore, there exists a const \tilde{C}_s such that

$$|\tilde{s}_n - s_n| \le \tilde{C}_s \cdot \delta_n. \tag{127}$$

We decompose the difference:

$$\left| \sum_{n=1}^{N} \hat{p}_{n} q_{n} \tilde{s}_{n} - \sum_{n=1}^{N} p_{n} q_{n} s_{n} \right| = \left| \sum_{n=1}^{N} (\hat{p}_{n} - p_{n}) q_{n} s_{n} + \sum_{n=1}^{N} \hat{p}_{n} q_{n} (\tilde{s}_{n} - s_{n}) \right|$$

$$\leq \sum_{n=1}^{N} |\delta_{n}^{p}| q_{n} s_{n} + \sum_{n=1}^{N} \hat{p}_{n} q_{n} \tilde{C}_{s} \delta_{n}$$

$$\leq \sum_{n=1}^{N} (C_{5} \cdot \delta_{n} + C_{6} \cdot \sum_{n=1}^{N} \delta_{i}) q_{n} s_{n} + \sum_{n=1}^{N} \hat{p}_{n} q_{n} \tilde{C}_{s} \delta_{n}.$$
(128)

Bounding the Denominator Term. We now analyze the denominator difference:

$$\left| \frac{1}{\sum_{n=1}^{N} \hat{p}_n q_n} - \frac{1}{\sum_{n=1}^{N} p_n q_n} \right| = \left| \frac{Y - \hat{Y}}{Y \cdot \hat{Y}} \right|, \quad \text{where} \quad Y := \sum_{n=1}^{N} p_n q_n, \ \hat{Y} := \sum_{n=1}^{N} \hat{p}_n q_n.$$

Assuming $Y, \hat{Y} \ge \epsilon > 0$, we have

$$\left|\frac{1}{\hat{Y}} - \frac{1}{Y}\right| \le \frac{1}{\epsilon^2} \cdot |Y - \hat{Y}| = \frac{1}{\epsilon^2} \cdot \left|\sum_{n=1}^{N} (p_n - \hat{p}_n)q_n\right|. \tag{130}$$

Using the previously derived bound $|\hat{p}_n - p_n| \le C_5 \cdot \delta_n + C_6 \cdot \sum_{i=1}^N \delta_i$, we obtain

$$|Y - \hat{Y}| \le \sum_{n=1}^{N} |p_n - \hat{p}_n| \cdot q_n \le \sum_{n=1}^{N} (C_5 \cdot \delta_n + C_6 \cdot \sum_{i=1}^{N} \delta_i) q_n$$

$$= C_5 \sum_{n=1}^{N} q_n \delta_n + C_6 \left(\sum_{i=1}^{N} \delta_i \right) \cdot \sum_{n=1}^{N} q_n.$$
(131)

Hence,

$$\left| \frac{1}{\hat{Y}} - \frac{1}{Y} \right| \le \frac{C_5 \cdot \|q\|_1 + C_6 \cdot \|q\|_1}{\epsilon^2} \cdot \sum_{n=1}^N \delta_n, \tag{132}$$

Finally, substituting (124), (129), and (132) into (122), we obtain:

$$|\mathcal{B}(\hat{\alpha}_{n,m},\hat{p}_{n}) - \mathcal{B}(\alpha_{n,m},p_{n})| \leq \left(\frac{\mathcal{B}_{\max}(C_{5} + C_{6})\|q\|_{1}}{\epsilon^{2}}\right) \cdot \sum_{n=1}^{N} \delta_{n}$$

$$+ \frac{1}{\epsilon} \left[L\eta \sigma^{2} \left(2C_{5} \sum_{n=1}^{N} p_{n} q_{n} \delta_{n} + 2C_{6} \sum_{n=1}^{N} p_{n} q_{n} \cdot \sum_{n=1}^{N} \delta_{n} \right) + \frac{5}{3} \left(\sum_{n=1}^{N} C_{5} \delta_{n} q_{n} s_{n} + C_{6} \sum_{i=1}^{N} \delta_{i} \sum_{n=1}^{N} q_{n} s_{n} + \sum_{n=1}^{N} \hat{p}_{n} q_{n} \tilde{C}_{s} \delta_{n} \right) \right] + \mathcal{O}(\delta_{n}^{2})$$

$$\leq \mathcal{O}\left(\sum_{i=1}^{N} \delta_{n}\right). \tag{133}$$

H MORE EXPERIMENTAL DETAILS

H.1 EXPERIMENTAL SETUP

Datasets and models. We conduct comprehensive experiments on four public benchmark datasets, including Fashion-MNIST (Xiao et al., 2017), CIFAR-10 (Krizhevsky et al., 2009), CINIC-10 (Darlow et al., 2018), and HAM10000 (Tschandl et al., 2018). Fashion-MNIST comprises 28×28 grayscale images of 70,000 fashion products with 10 classes, and there are 60,000 training images and 10,000 testing images. The CIFAR-10 dataset consists of 50,000 training images and 10,000 testing images, each with a size of $3 \times 32 \times 32$. CINIC-10 is an extension of CIFAR-10 via the addition of downsampled ImageNet (Deng et al., 2009) images. The HAM10000 dataset is an image dataset used for skin lesion classification in the medical field. For Fashion-MNIST, We use a LeNet-5 (LeCun et al., 1998) with two 5×5 convolutional layers, each followed by ReLU activation and max pooling, and three fully connected layers. For CIFAR-10 and CINIC-10, we use an 8-layer AlexNet (Krizhevsky et al., 2012) with a size of 136 MB. For HAM10000, we use a customized CNN consisting of three 3×3 convolutional layers (with ReLU activation and max pooling) and two fully connected layers with dropout regularization.

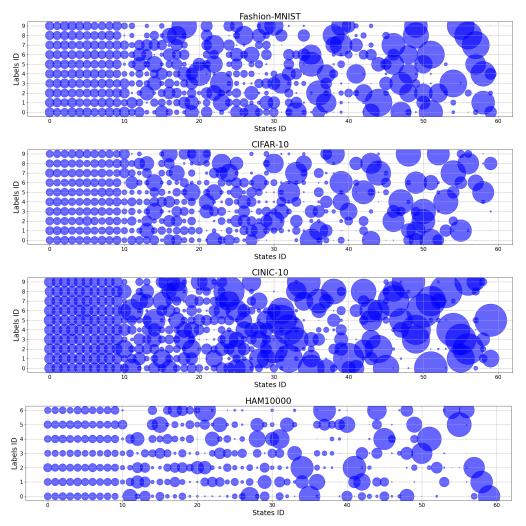


Figure 4: Visualization of the class distribution under stratified state space on the four datasets, where the bubble size represents the number of class samples per state.

Intra-state heterogeneity. We organize the 60 latent states into 6 clusters (10 states per cluster), each associated with a Dirichlet partitioning strategy (Wang et al., 2020a) with distinct concentration

parameters $\alpha \in \{0.05, 0.1, 0.2, 0.5, 1.0, 100.0\}$. We visualize the class distribution under stratified state space on the four datasets in Fig. 4.

Two heterogeneity scenarios. We consider two heterogeneity scenarios. (1) Full-access (mild heterogeneity): Each client has non-zero access probability $\pi_{n,m} > 0$ for all states $m \in \{1,\ldots,M\}$. (2) Partial-access (extreme heterogeneity): Each client is restricted to a randomly sampled subset of 10 states, with $\pi_{n,m} = 0$ for others. Moreover, 50% of clients are initialized with latent states drawn from high-heterogeneity clusters ($\alpha \in \{0.05, 0.1\}$). We visualize the state probability $\pi_{n,m}$ of all clients under two scenarios in Fig. 5.

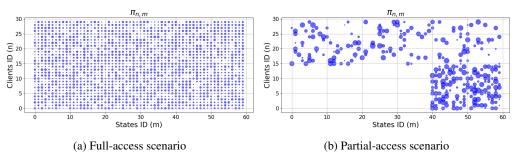


Figure 5: Visualization of the state probability of all clients under two scenarios, where the bubble size reflects the probability of a client reaching that state.

Implementation details. The overall framework of SFedPO is implemented with Pytorch (Paszke et al., 2019), and all experiments are conducted on an Intel(R) Xeon(R) Platinum 8352V CPU and an NVIDIA A40 (48GB) GPU, and 256GB RAM. We run federated learning for 100 rounds. The number of time steps and batch size are 5 and 64, respectively. We use an SGD optimizer with a 0.01 learning rate, and the weight decay is set to 1e-4. For each client n, we set its sampling budget α_n to 0.5 and allocate a data capacity D of 500 samples. While our theoretical framework defines d_m as the upper bound on the gradient variance under state m, estimating such variance is often impractical in streaming scenarios with limited and evolving local data. Inspired by (Ye et al., 2023), we adopt a practical surrogate by assuming that d_m is proportional to the discrepancy between the class distribution and a uniform distribution. Accordingly, we use the KL divergence as a proxy measure for d_m in all experiments.

Client-specific sampling ratios. We visualize the client-specific sampling ratios $\{\alpha_{n,m}\}_{m=1}^M$ under two scenarios in Fig. 6. We observe that the client-specific sampling ratios exhibit a clustering pattern aligned with the clustered structure of the state space. Moreover, as the distribution of the state space becomes more imbalanced, the client-specific sampling ratios tend to decrease accordingly.

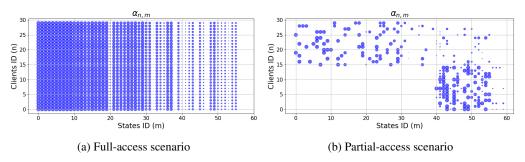


Figure 6: Visualization of the client-specific sampling ratios of all clients under two scenarios, where the bubble size reflects the value of sampling ratios.

Baselines. We compare SFedPO with the following methods:

- **FedOGD:** FedOGD is a vanilla online federated Learning method that processes sequential data and performs online optimization per round.
- FedOMD: FedOMD is an online federated Learning method, which performs online mirror descent and achieves sublinear regret.

- Adaptive-FedAvg: Adaptive-FedAvg adjusts local learning rates to improve model plasticity under concept drift in federated Learning.
- **FedDrift:** FedDrift uses a hierarchical clustering method to address the problem of concept drift adaptation in federated Learning.
- **Flash:** Flash synergizes client-side early-stopping to detect concept drifts with server-side drift-aware adaptive optimization to effectively adjust the global learning rate.
- **FedEWC:** EWC is a classic regularization-based method for continual learning that uses the Fisher information matrix to estimate the importance of parameters. In our streaming scenario, each client leverages EWC to retain information from the previous state.
- FLwF-2T: FLwF-2T is a distillation-based method that deals with catastrophic forgetting in federated continual learning. FLwF-2T enables the current model to distill knowledge from both the model trained in the previous state and the global model from the last round.
- **IS:** Important Sampling (IS) uses a higher gradient norm to reflect the informativeness of the data. In our streaming scenario, it selects the optimal data from the new state based on the metric and discards the old data accordingly.
- **ODE:** ODE introduces an online data selection method based on a data valuation metric: the projection of local gradients onto the global gradient. In our streaming scenario, it selects the optimal data from the new state based on the metric and discards the old data accordingly.
- **DRSR:** DRSR updates the local data of each client according to the distribution discrepancy between the long-term data distribution and the client's local data. In our streaming scenario, it selects informative samples from both the existing and incoming data based on the update rule.

Since some methods (i.e., FedOGD, FedOMD, Adaptive-FedAvg, FedDrift, Flash, FedEWC, FLwF-2T) do not incorporate any data sampling mechanism, we assume that each client fully replaces its local dataset with new data drawn from the current state distribution, thereby aligning with the streaming FL setting.

FL methods in modularity experiment. We apply SFedPO's data sampling and aggregation strategies to several representative FL methods: **FedAvg** (McMahan et al., 2017), which is the pioneering FL method; **FedProx** (Li et al., 2020a), a classic regularization-based FL method; **FedCurv** (Shoham et al., 2019), a FL method based on curvature adjustment regularization; **FedNTD** (Lee et al., 2022), which preserves the global knowledge by not-true distillation in FL; **FedEXP** (Jhunjhunwala et al., 2023), which tune the global learning rate via extrapolation to speed up the global convergence. In our streaming data scenario, these methods sample data based on state transitions in each round and train local models accordingly.

Algorithm-specific hyperparameters. The above baselines and FL methods adopt the same experimental setup as SFedPO, while the algorithm-specific hyperparameters are configured as follows:

- Flash: we tune the global learning rate $\eta_q \in \{0.001, 0.01, 0.1, 1.0\}$, and set it to 0.01.
- **FedEWC:** we tune the penalty coefficient $\lambda \in \{0.001, 0.01, 0.1, 1.0\}$, and set it to 0.1. λ is a scalar that balances the contribution of the regularization loss relative to the cross-entropy loss in the total objective.
- **FLwF-2T:** we tune the $\alpha \in \{0.001, 0.01, 0.1, 0.3, 0.7\}, \beta \in \{0.001, 0.01, 0.1, 0.3, 0.7\}$, and set them to $\alpha = 0.3, \beta = 0.3$.
- **FedProx:** we tune the $\mu \in \{0.001, 0.01, 0.1, 1.0\}$, and set it to 0.1.
- **FedCurv:** we tune the $\lambda \in \{0.001, 0.01, 0.1, 1.0\}$, and set it to 0.01.
- **FedNTD:** we tune the $\beta \in \{0.01, 0.1, 1.0, 2.0\}$, and set it to 1.0.

H.2 OTHER EXPERIMENTS

Performance comparison. Fig. 7 presents the convergence performance of SFedPO and all baselines on four datasets under two scenarios (here we use F-MNIST for Fashion-MNIST).

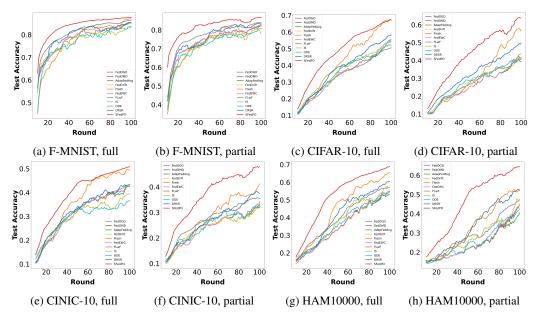


Figure 7: Performance comparison of SFedPO to other baselines in full and partial scenarios on several datasets.

Effects of different configurations under other datasets. We vary three core parameters, including time step $(T \in \{2, 5, 8, 10\})$, training round $(R \in \{50, 100, 150, 200\})$, and data capacity of clients $(D \in \{250, 500, 750, 1000\})$. Then, we respectively show the performance of our method and four baselines on Fashion-MNIST, CINIC-10, and HAM10000 in Fig. 8, Fig. 9, and Fig. 10. The experimental results demonstrate that our method outperforms all baseline approaches across different parameters under various datasets.

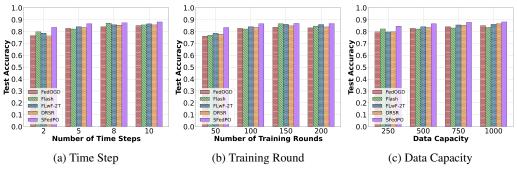


Figure 8: Fashion-MNIST.

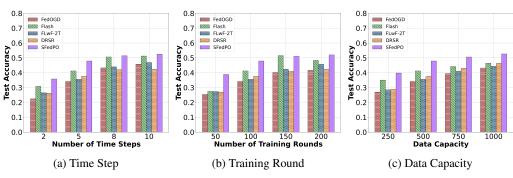


Figure 9: CINIC-10.

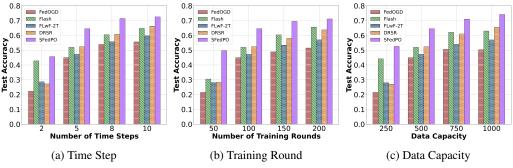


Figure 10: HAM10000.

Effects of number of clients N. In Table 5, we vary the number of clients $N \in \{10, 50, 100\}$ and compare the performance of SFedPO with several baselines in full and partial scenarios on CIFAR-10. The results show that SFedPO consistently outperforms the baselines across different N.

Table 5: Test accuracy (%, mean \pm std on 5 trials) comparison of our SFedPO framework to several baselines with different N in full and partial scenarios on CIFAR-10.

Method		CIFAR-10(full))	C	IFAR-10(partial)			
	N = 10	N = 50	N = 100	N = 10	N = 50	N = 100		
FedOGD Flash FLwF-2T	43.04±1.92 49.61±1.47 46.94±1.42	53.78±1.38 59.21±0.51 50.00±0.96	55.32±1.29 60.94±0.60 50.06±0.62	34.10±2.70 32.69±3.52 38.69±1.58	42.38±1.60 48.32±1.91 40.24±0.54	45.42±1.29 52.78±2.45 40.26±0.41		
DRSR SFedPO	51.45±1.99 63.51±0.90	58.74 ± 0.43 67.94 ± 0.23	$\frac{59.98 \pm 0.45}{68.62 \pm 0.26}$	43.43±1.22 55.66 ± 1.24	50.70 ± 1.06 65.57 ± 0.34	53.44±0.78 66.74±0.43		

Effects of model architectures. We evaluate SFedPO against several baselines in full and partial scenarios on CIFAR-10 under different model architectures, including CNN, ResNet-10 (He et al., 2016), ResNet-18 (He et al., 2016), and ResNet-34 (He et al., 2016). The results in Table 6 show that SFedPO consistently achieves the best performance, demonstrating the robustness of our method across model architectures.

Table 6: Test accuracy (%, mean±std on 5 trials) comparison of our SFedPO framework to several baselines with different model architectures in full and partial scenarios on CIFAR-10.

Method		CIFAR-	-10(full)			CIFAR-10(partial)			
	CNN	ResNet10	ResNet18	ResNet34	CNN	ResNet10	ResNet18	ResNet34	
FedOGD	65.13±0.62	45.46±1.05	47.16±1.53	50.05±1.32	56.95±1.21	37.18±1.45	40.11±2.22	41.67±1.07	
Flash	67.29 ± 0.56	49.99 ± 1.80	60.14 ± 3.41	55.26 ± 2.00	60.59±3.15	42.25 ± 1.49	51.74 ± 2.035	47.07 ± 1.55	
FLwF-2T	66.51±0.69	48.60 ± 0.20	49.04 ± 1.04	51.42 ± 1.53	59.02±1.50	40.60 ± 2.33	42.81 ± 1.19	44.74 ± 1.28	
DRSR	67.56±0.67	51.12 ± 1.13	54.32 ± 0.55	57.80 ± 1.03	61.11±0.63	43.91 ± 0.91	47.01 ± 1.11	48.62 ± 1.24	
SFedPO	74.10±0.42	63.69±1.06	66.64±0.87	68.69±0.33	70.56±0.71	59.38±0.83	61.50±1.46	64.39±0.96	

Effects of different discrepancy metrics. While our theoretical framework defines d_m as the upper bound on the gradient variance under state m, estimating such variance is often impractical in streaming scenarios with limited and evolving local data. Inspired by (Ye et al., 2023), we adopt a practical surrogate by assuming that d_m is proportional to the discrepancy between the class distribution and a uniform distribution. To investigate the impact of different discrepancy metrics, we evaluate four commonly used measures on the CIFAR-10 dataset: L1 & L2 distance, KL divergence, and Jensen–Shannon (JS) divergence (Fuglede & Topsoe, 2004). As shown in Table 7, all metrics yield comparable performance and consistently outperform baseline methods, highlighting the robustness of our framework to the choice of discrepancy measure.

Table 7: Accuracy (%) under different discrepancy metrics.

Metric	L1	L2	JS	KL	Flash	DRSR
full	67.56	67.58	67.56	68.18	67.00	58.36
partial	63.23	63.28	63.75	64.41	56.55	49.60

Effects of Prediction Error. To simulate prediction errors, we perturb each state probability $\hat{\pi}_n$ by a random noise uniformly drawn from $[-\epsilon, \epsilon]$, followed by renormalization to ensure $\sum_m \hat{\pi}_{n,m} = 1$. We vary ϵ from 0.00 to 0.10 to simulate increasing levels of oracle error. As shown in Figs. 11, SFedPO tends to maintain stable performance as the degree of perturbation varies on the CIFAR-10 dataset.

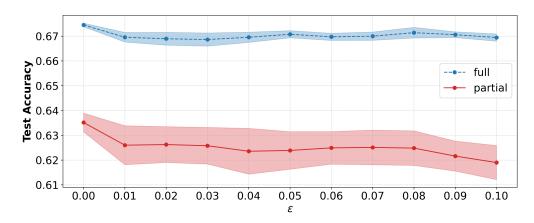


Figure 11: Performance under different degrees of perturbation.