

---

# FedEve: On Bridging the Client Drift and Period Drift for Cross-device Federated Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1      Federated learning (FL) is a machine learning paradigm that allows multiple  
2      clients to collaboratively train a shared model without exposing their private data.  
3      Data heterogeneity is a fundamental challenge in FL, which can result in poor  
4      convergence and performance degradation. *Client drift* has been recognized as one  
5      of the factors contributing to this issue resulting from the multiple local updates  
6      in FedAvg. However, in cross-device FL, a different form of drift arises due to  
7      the partial client participation, but it has not been studied well. This drift, we  
8      referred as *period drift*, occurs as participating clients at each communication  
9      round may exhibit distinct data distribution that deviates from that of all clients.  
10     It could be more harmful than client drift since the optimization objective shifts  
11     with every round. In this paper, we investigate the interaction between period drift  
12     and client drift, finding that period drift can have a particularly detrimental effect  
13     on cross-device FL as the degree of data heterogeneity increases. To tackle these  
14     issues, we propose a predict-observe framework and present an instantiated method,  
15     FedEve, where these two types of drift can compensate each other to mitigate their  
16     overall impact. We provide theoretical evidence that our approach can reduce the  
17     variance of model updates. Extensive experiments demonstrate that our method  
18     outperforms alternatives on non-iid data in cross-device settings.

19     

## 1 Introduction

20     Federated learning is a decentralized machine learning approach that enables multiple clients to  
21     collaboratively train a shared model without exposing their private data [McMahan et al., 2017]. In  
22     this paradigm, each client independently trains a local model using its own data and subsequently  
23     sends the model updates to a central server. The server then periodically aggregates these updates  
24     to improve the global model until it reaches convergence. There are two primary settings in FL:  
25     cross-silo and cross-device [Kairouz et al., 2021]. Cross-silo FL typically involves large organizations  
26     (small number of clients), where most clients actively participate in every round of training [Chen  
27     and Chao, 2021, Lin et al., 2020]. In contrast, cross-device FL focuses on scenarios like smartphones  
28     (huge number of clients, e.g., millions), where only a limited number of clients participate in each  
29     round [Li et al., 2020b, Reddi et al., 2020], due to communication bandwidth, client availability, and  
30     other issues. This paper primarily focuses on the cross-device setting with partial client participation  
31     since we discover and then solve its unique challenge — “*period drift*”.

32     Distinguished from traditional distributed optimization, the statistical heterogeneity of data has been  
33     acknowledged as a fundamental challenge in FL [Li et al., 2020a, Chen and Chao, 2021, Lin et al.,  
34     2020]. This data heterogeneity refers to the violation of the independent and identically distributed  
35     (non-iid) data assumption across clients, which can result in poor convergence and performance  
36     degradation when using FEDAVG . *Client drift* is recognized as one of the factors contributing to this  
37     issue and attracts numerous efforts to address it [Karimireddy et al., 2021, Li et al., 2020b, Reddi et al.,

38 This phenomenon is characterized by clients who, after multiple local updates, progress too far  
 39 towards minimizing their local objective, consequently diverging from the shared direction. However,  
 40 in cross-device FL, a different form of drift exists and could be more detrimental to the training  
 41 process than client drift, which has not been extensively studied. *This drift occurs periodically as*  
 42 *different clients participate in each communication round, and these participating clients as a group*  
 43 *may exhibit distinct data distribution that deviates from the overall distribution of all clients.* Period  
 44 drift is fundamentally different from the noise in SGD, as further described in Appendix C.

45 This deviation could potentially lead to  
 46 slow and unstable convergence, as the  
 47 optimization objective shifts with every  
 48 round. For simplicity, we refer to this  
 49 phenomenon as *period drift*. Despite both  
 50 period drift and client drift being rooted  
 51 in data heterogeneity, they stem from  
 52 different causes (as illustrated in Figure 1).  
 53 Client drift results from multiple local  
 54 updates and the non-iid, while period drift  
 55 arises due to partial client participation and the non-iid. The concept of period drift and client drift  
 56 is shown in the following equation:

$$w^* \neq w_N^* \neq w_S^* \neq \bar{w} = \frac{1}{|\mathcal{S}|} \sum w_k^* \neq w_k^*, \quad (1)$$

**Period Drift**      **Client Drift**

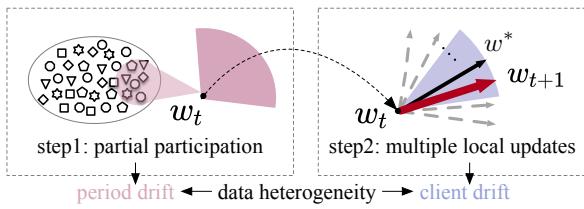


Figure 1: The generation of period drift and client drift

57 which is further described in detail in Appendix B.1. These two drifts, when occurring simultaneously,  
 58 significantly increase the complexity of achieving stable and efficient convergence, particularly  
 59 challenging the effectiveness of cross-device FL systems.

60 In this paper, we first investigate the impact of period drift and client drift, finding that period drift  
 61 can have a particularly detrimental effect on cross-device FL as the degree of data heterogeneity  
 62 increases (as demonstrated in detail in Section 3.2). *While the impacts of period drift and client drift*  
*are additive, we fortunately uncover a cooperative mechanism therby these two types of drift can*  
*compensate each other to mitigate their overall impact.* To achieve this, we propose a predict-observe  
 63 framework, where we consider at each round 1) the server optimization (e.g., momentum) as a  
 64 prediction of a update step of FL; 2) the clients' optimization (e.g., local SGD) as an observation of  
 65 this update step. Note that the vanilla FEDAVG is a special case in which the server does not make  
 66 any predictions and solely relies on the observation provided by clients. In this framework, period  
 67 drift and client drift are viewed as the noise respectively associated with prediction and observation.  
 68 We thereby incorporate a Bayesian filter to integrate prediction (with period drift) and observation  
 69 (with client drift) to achieve a better estimation of update step and reduce uncertainties. Based on  
 70 the predict-observe framework, we present an instantiated method, referred as FEDEVE , which  
 71 combines the prediction and observation through linear interpolation. The coefficient of this linear  
 72 combination indicate the relative confidence between prediction and observation, which is determined  
 73 by the variance of the period drift and client drift, thus produces a more precise estimation of updates.  
 74 FEDEVE does not increase the client storage or extra communication costs, and does not introduce  
 75 additional hyperparameter tunning, making it ideal for cross-device FL.

76 **Contributions** We summarize the primary contributions of this paper as follows:

- 77 • We analyze the impact of period drift and client drift for cross-device FL, and observe that period  
 78 drift has a particularly detrimental effect as the degree of data heterogeneity increases.
- 79 • We propose a predict-observe framework for cross-device FL that incorporates a Bayesian filter  
 80 to integrate server optimization and clients' optimization so that period drift and client drift can  
 81 compensate for each other.
- 82 • As an instantiation of the proposed framework, we present FEDEVE to combine prediction and  
 83 observation through linear interpolation based on the variance of the period drift and client drift.
- 84 • We provide theoretical evidence within our framework that FEDEVE can reduce the variance of  
 85 model updates. Extensive experiments demonstrate that our method outperforms alternatives on  
 86 non-iid data in cross-device settings.<sup>1</sup>

<sup>1</sup>For a smoother reading experience, please feel free to check out our reading guide in Appendix A.

89 **2 Methodology**

90 In this section, we discuss the problem of some methods (e.g., FEDAVG ) in cross-device FL, and  
 91 then propose our predict-observe framework and a method FEDEVE to deal with it.

92 **2.1 Typical Federated Learning Setup**

93 Federated learning, as described by McMahan et al. [2017], involves utilizing multiple clients and  
 94 a central server to optimize the overall learning objective. The goal is to minimize the following  
 95 objective function:

$$\min_w f(w) = \sum_{k=1}^N p_k F_k(w) = \mathbb{E}_k[F_k(w)], \quad (2)$$

96 where  $N$  is the number of clients,  $p_k \geq 0$ , and  $\sum_k p_k = 1$ . In general, the global objective is the ex-  
 97 pectation of the local objective over different data distributions  $\mathcal{D}_k$ , i.e.,  $F_k(w) = \mathbb{E}_{x_k \sim \mathcal{D}_k} f_k(w; x_k)$ ,  
 98 with  $n_k$  samples on each client  $k$  and weighted by  $p_k$ . We set  $p_k = \frac{n_k}{n}$ , where  $n = \sum_k n_k$  is the total  
 99 number of data points. In deep learning setting,  $F_k(w)$  is often non-convex. A common approach  
 100 to solve the objective (2) in federated settings is FEDAVG [McMahan et al., 2017]. For example, in  
 101 cross-device FL, a small subset  $\mathcal{S}_t$  ( $|\mathcal{S}_t| \ll N$ ) of the total clients are selected at each round (ideally  
 102 randomly, but possibly biased in practice), and then the server broadcasts its global model to the  
 103 selected client. In parallel, each of the selected clients runs SGD on their own loss function  $F_k(\cdot)$   
 104 for  $E$  number of epochs, and sends the resulting model to the server. The server then updates its  
 105 global model as the average of these local models and repeats this process until convergence. One  
 106 problem of FL is the non-iid data across clients, which can bring about “client drift” in the updates of  
 107 each client, resulting in slow and unstable convergence [Karimireddy et al., 2021]. Despite efforts to  
 108 address the problem of client drift [Karimireddy et al., 2021, Li et al., 2020b, Reddi et al., 2020], there  
 109 is a lack of research on the issue of period drift, i.e. the data distribution of selected clients at each  
 110 round may differ from the overall data distribution of all clients. Period drift along with client drift  
 111 can greatly impact the convergence of the learning process in FL, thus we propose a predict-observe  
 112 framework to deal with them.

113 **2.2 The Concept of Drift**

114 In contrast to conventional distributed optimization, federated learning possesses distinct characteris-  
 115 tics, such as client sampling, multiple local epochs, and non-iid data distribution. These attributes  
 116 may lead to a drift in the updates of global model, resulting in suboptimal performance. This drift  
 117 can be thought of as a noise term that is added to the true optimization states during the optimization  
 118 process. Thus, we can make the definition of period drift and client drift as:

119 **Definition 2.1** (Period Drift and Client Drift). In federated learning, two types of drift arise from  
 120 data heterogeneity:

- 121 • **Period Drift:** The deviation when the data distribution of selected clients differs from the overall  
 122 distribution. Formally:

$$\text{Period Drift} := \mathbb{E}_{\mathcal{S}_t} \left[ \left\| \frac{1}{|\mathcal{S}_t|} \sum_{k \in \mathcal{S}_t} \nabla F_k(w) - \nabla f(w) \right\|^2 \right], \quad (3)$$

123 where  $\mathbb{E}_{\mathcal{S}_t}$  is the expectation over client sampling,  $\mathcal{S}_t$  is the subset of clients selected at round  $t$ ,  
 124  $\nabla F_k(w)$  is the gradient of client  $k$ ’s objective, and  $\nabla f(w)$  is the gradient of the global objective.

- 125 • **Client Drift:** The deviation when the averaged optima of local objectives does not align with the  
 126 optimum of the averaged objective. Formally:

$$\text{Client Drift} := \mathbb{E}_{k \in \mathcal{S}_t} \left[ \left\| \nabla F_k(w_k^*) - \nabla F_k(w_{\mathcal{S}_t}^*) \right\|^2 \right], \quad (4)$$

127 where  $w_k^*$  is the local optimum for client  $k$  and  $w_{\mathcal{S}_t}^*$  is the optimum for the selected clients’  
 128 averaged objective.

129 These drifts are conceptually independent but together affect model convergence and performance.  
 130 See Appendix B.1 for detailed discussion.

131 **Assumption 2.2.** *The aggregated model parameters on the server  $w_{server}$ , can be represented  
 132 as the sum of the optimal parameters  $w^*$  and a drift (noise) that follows a normal distribution  
 133  $w_{drift} \sim \mathcal{N}(0, \sigma_{drift}^2)$ :*

$$w_{server} = w^* + w_{drift} \leftarrow noise, \quad (5)$$

134 where  $w^*$  represents the optimal parameters obtained through the use of stochastic gradient descent  
 135 (SGD),  $w_{drift}$  represents the noise term caused by factors such as client sampling, multiple local  
 136 epochs, and non-iid data distribution that we assume a normal distribution, and  $w_{server}$  represents  
 137 the aggregated model parameters also follows a normal distribution  $w_{server} \sim \mathcal{N}(w^*, \sigma_{drift}^2)$ ,  
 138 with the expectation of the aggregate model parameters being equal to the optimal parameters, i.e.  
 139  $\mathbb{E}[w_{server}] = w^*$ . Note that the assumption of Gaussian-like noise is natural, and its justification can  
 140 be found in Appendix D.1.

### 141 2.3 The Predict-observe Framework

142 Initially, we establish the concept of period drift, represented by  $Q_t$ , and client drift, represented by  
 143  $R_t$  at the  $t$ -th communication round. We first make an assumption of independence concerning the  
 144 two types of drift, which states that the two drifts are independent of one another. This assumption  
 145 allows us to more accurately analyze the impact of each drift on the model's performance and devise  
 146 methods to mitigate their effects.

147 **Assumption 2.3.** *The initialization model parameters are independent of all period drifts  $Q_t$   
 148 and client drifts  $R_t$  at each communication round, that is  $w_0 \perp Q_0, Q_1, \dots, Q_t$  and  $w_0 \perp R_0, R_1, \dots, R_t$ .*

150 The justification and limitation of this assumption can be found in Appendix D.2. Since the clients  
 151 participating in each round in cross-device FL is only a small fraction of all clients, period drift can  
 152 be attributed to the discrepancy that the objective of selected clients at each round does not align with  
 153 the overall objective. Thus, an effective prediction of updates can potentially help reduce the period  
 154 drift. As formulated in Equation (5), we express the prediction of updates on the server as:

$$\hat{w}_{t+1} = g(w_t) + Q_t, \quad Q_t \sim \mathcal{N}(0, \sigma_{Q_t}^2), \quad (6)$$

155 where  $\hat{w}_{t+1}$  is the prediction model of  $(t+1)$ -th round as the output of predict function  $g(\cdot)$  with  
 156 the current model  $w_t$  as input. It is noteworthy that the period drift at the  $t$ -th round is represented  
 157 by  $Q_t$ , and just like the drift in assumption 2.2, it is assumed to follow a normal distribution  
 158  $\mathcal{N}(0, \sigma_{Q_t}^2)$ , characterized by a mean of zero and a variance of  $\sigma_{Q_t}^2$ . Client drift can be attributed to  
 159 the phenomenon that the averaged optima of objectives does not align with the optima of averaged  
 160 objectives. Thus, we consider the updates provided by these clients is a kind of observation of global  
 161 updates. As formulated in Equation (5), we express it as:

$$\tilde{w}_{t+1} = h(\hat{w}_{t+1}) + R_t, \quad R_t \sim \mathcal{N}(0, \sigma_{R_t}^2), \quad (7)$$

162 where  $\tilde{w}_{t+1}$  is the model of  $(t+1)$ -th round as the output of observe function  $h(\cdot)$  with the predict  
 163 model  $\hat{w}_{t+1}$  as input. Also, the client drift at the  $t$ -th round is represented by  $R_t$ , and just like the  
 164 drift in assumption 2.2, it is assumed to follow a normal distribution  $\mathcal{N}(0, \sigma_{R_t}^2)$ , characterized by a  
 165 mean of zero and a variance of  $\sigma_{R_t}^2$ . It is clear that standard FEDAVG is a special case since there  
 166 is no prediction for server optimization, and it solely relies on the observations provided by clients.  
 167 Furthermore, the period drift,  $Q_t$ , and the client drift,  $R_t$ , are represented as noise terms that are  
 168 incorporated into the prediction and observation functions. According to assumption (2.3), these  
 169 drifts are independent of the current model states, and the lemma of independence noise is posited:

170 **Lemma 2.4. (Independence of Noise).** *the noise present in the prediction and observation at each  
 171 communication round is independent of the current model state, specifically,  $w_t \perp Q_t$  and  $w_t \perp R_t$ .*

172 The complete proof of the independence of noise can be found in appendix D.3. The equations pre-  
 173 sented in equations 6 and 7 depict the prediction and observation of updates, respectively, taking into  
 174 account both period drift and client drift. In order to reconcile the discrepancy between the prediction  
 175 (including period drift) and observation (including client drift), a Bayesian filter is introduced to  
 176 allow for compensation between the two sources of drift. The prior probability of  $w_{t+1}$  is represented  
 177 by  $P(\hat{w}_{t+1})$ , and by combining the observation  $P(\tilde{w}_{t+1})$  and the likelihood  $P(\tilde{w}_{t+1} | \hat{w}_{t+1})$ , the  
 178 posterior probability  $P(w_{t+1} | \tilde{w}_{t+1})$  of  $w_{t+1}$  can be calculated as the new model at the  $(t+1)$ -th  
 179 round, as shown in Equation (8).

$$P(w_{t+1}) := P(\hat{w}_{t+1} | \tilde{w}_{t+1}) = \frac{P(\tilde{w}_{t+1} | \hat{w}_{t+1})P(\hat{w}_{t+1})}{P(\tilde{w}_{t+1})}. \quad (8)$$

180 By utilizing the Bayesian filter in our predict-observe framework, an update mechanism is implemented  
 181 that first performs prediction and then observes the predicted model state, as described in the  
 182 following procedure:

$$f_{w_t}^+(w) \xrightarrow{\text{predict}} f_{\hat{w}_{t+1}}^-(w) = \int_{-\infty}^{+\infty} f_{Q_t}[w - f(v)] f_{w_t}^+(v) dv \quad (9)$$

$$\xrightarrow{\text{observe}} f_{w_{t+1}}^+(w) = \eta_t \cdot f_{R_t}[w_{t+1} - h(w)] \cdot f_{\hat{w}_{t+1}}^-(w),$$

183 where  $f_{w_t}^+(w)$  is the posterior probability of  $w_t$ ,  $f_{\hat{w}_{t+1}}^-(w)$  is the prior probability of  $w_{t+1}$ ,  $f_{Q_t}$  is the  
 184 PDF of period drift,  $f_{w_{t+1}}^+(w)$  is the posterior probability of  $w_{t+1}$ ,  $f_{R_t}$  is the PDF of client drift, and  
 185  $\eta_t = \left\{ \int_{-\infty}^{+\infty} f_{R_t}[\tilde{w}_{t+1} - h(\hat{w}_{t+1})] f_{\hat{w}_{t+1}}^-(w) dw \right\}^{-1}$ . By combining prediction and observation, the  
 186 fused model can be estimated by taking the expectation of the posterior probability as follow:

$$\hat{w}_{t+1} = E[f_{w_{t+1}}^+(w)] = \int_{-\infty}^{+\infty} w f_{w_{t+1}}^+(w) dw. \quad (10)$$

187 **Theorem 2.5.** Given assumption 2.2 and lemma 2.4, the composite model will exhibit a diminished  
 188 degree of variance in comparison to the individual variances of both period drift and client drift, and  
 189 the mean will be a linear combination that is weighted by the variances:

$$\mu_{\text{fused}} = \frac{\mu_1 \sigma_{R_t}^2 + \mu_2 \hat{\sigma}_{t+1}^2}{\sigma_{R_t}^2}, \quad \sigma_{\text{fused}}^2 = \frac{\hat{\sigma}_{t+1}^2 \sigma_{R_t}^2}{\hat{\sigma}_{t+1}^2 + \sigma_{R_t}^2}, \quad (11)$$

190 where  $\mu_1, \mu_2, \mu_{\text{fused}}$  is the mean of prediction, observation and fused model, and  $\hat{\sigma}_{t+1}^2, \sigma_{R_t}^2, \sigma_{\text{fused}}^2$   
 191 is the variance of prediction, client drift and fused model.

192 The complete proof of the bayesian filter  
 193 can be found in appendix D.5. The applica-  
 194 tion of Bayesian filtering allows for the  
 195 interaction of period drift and client drift  
 196 to generate a new model, which is charac-  
 197 terized by a reduced level of variance as  
 198 compared to the individual variances of pe-  
 199 riod drift and client drift, as depicted in  
 200 Figure 2(a). However, the computation of  
 201 the new model is challenging due to the  
 202 presence of infinite integrals in Equation  
 203 (10) and  $\eta_t$ , as it is a general framework for  
 204 any prediction and observation function. In the following section, we will propose a specialized  
 205 method to facilitate the convergence of FL.

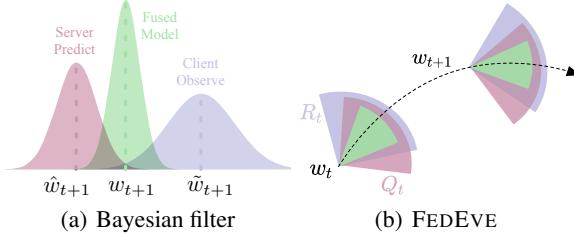


Figure 2: Illustrations of the framework and FEDeve

206 **2.4 The FEDeve Method**

207 The predict-observe framework has been proposed as a strategy for mitigating the challenges of  
 208 period drift and client drift. However, it also raises some questions regarding the effective method  
 209 of prediction and the variance associated with both period and client drift. In this section, we  
 210 demonstrate that the utilization of momentum as a server optimization [Hsu et al., 2019, Reddi et al.,  
 211 2020], can serve as an effective prediction method. Furthermore, we present a method for estimating  
 212 the variance of period drift and client drift. In the context of the predict-observe framework, we have  
 213 adapted it to a specific setting where Nesterov momentum is employed as the prediction function  $g(\cdot)$ ,  
 214 and the observation function  $h(\cdot)$  is the average of the models from the clients the same as FEDAVG .  
 215 We have reformulated FEDAVG in an incremental form as the starting point of our approach.

$$w_{t+1} = \sum_{k \in \mathcal{S}_t} p_k w_t^k = w_t - \sum_{k \in \mathcal{S}_t} p_k (w_t - w_t^k) \quad (12)$$

$$= w_t - \sum_{k \in \mathcal{S}_t} p_k \Delta w_t^k = w_t - \Delta w_t.$$

216 This formulation facilitates the accumulation of  $\Delta w_t$  as the momentum on the server, which serves  
 217 as a prediction of updates, as the empirical value of the hyperparameter  $\beta = 0.9$  suggests that the

218 direction of historical updates is likely to be maintained. By introducing the Nesterov momentum and  
 219 specialize  $g(w_t) = w_t - \eta_g M_t$  in Equation (6) as the prediction function. Additionally, we specialize  
 220  $h(\hat{w}_{t+1}) = \hat{w}_{t+1} - \eta_g \Delta \tilde{w}_t$  in Equation (7) as the observation function. Thus, the predict-observe  
 221 equation can be rewritten as follow:

$$\hat{w}_{t+1} = w_t - \eta_g M_t + Q_t, \quad (13)$$

$$\tilde{w}_{t+1} = \hat{w}_{t+1} - \eta_g \Delta \tilde{w}_t + R_t, \quad (14)$$

222 where  $M_t$  is the momentum (the accumulation of  $\Delta w_t$ ) at  $t$ -th round,  $\Delta \tilde{w}_t$  is the average of model  
 223 update in Equation (12) from clients at the states of  $\hat{w}_{t+1}$ , and  $\eta_g$  is the global learning rate. By  
 224 assuming a normal distribution for  $Q_t$  and  $R_t$  based on the equations (6), (7), and (8), the problem of  
 225 infinite integral in Equation (10) and  $\eta_t$  can be solved in a closed-form, as detailed in reference D.5.2.  
 226 Additionally, due to the normal distribution, the form of distribution like equations 15e, 10 is not  
 227 necessary, and only the mean and variance are used to depict the model update process. Since these  
 228 equations are linear in nature, the Bayesian filter can be specialized as the Kalman Filter (KF). The  
 229 process of model update can thus be summarized as the use of KF, as represented by the following  
 230 formulation:

$$\hat{w}_t + 1 = w_t - \eta_g M_t, \quad (15a)$$

$$\hat{\sigma}^2 t + 1 = \sigma_t^2 + \sigma_{Q_t}^2, \quad (15b)$$

$$G_{kal} = \frac{\hat{\sigma}_{t+1}^2}{\hat{\sigma}_{t+1}^2 + \sigma_{R_t}^2}, \quad (15c)$$

$$M_{t+1} = M_t + G_{kal}(\Delta \tilde{w}_t - M_t), \quad (15d)$$

$$w_{t+1} = w_t - \eta_g M_{t+1}, \quad (15e)$$

$$\sigma_{t+1}^2 = (1 - G_{kal})\hat{\sigma}_{t+1}^2. \quad (15f)$$

231 The six steps of model update for each communication round in our method are outlined in Equations  
 232 (15a)-(15f). Equation (15a) predicts the model states  $w_t$  using the momentum  $M_t$ . Equation (15b)  
 233 estimates the variance of the prediction model by summing the variance of  $w_t$  and the period drift  $Q_t$ .  
 234 To provide a clear representation, the variance of the prediction model is represented by  $\sigma^-$  and the  
 235 variance of the fused model is represented by  $\sigma^+$ . The core of our method is presented in Equation  
 236 (15c), where the Kalman gain  $G_{kal}$  is calculated based on the ratio of the variance of the prediction  
 237  $\hat{\sigma}_{t+1}^2$  and the observation (client drift)  $R_t$ . The value of  $G_{kal}$  determines the relative weight of  
 238 the prediction and observation when they are combined. Equation (15d) fuses the prediction and  
 239 observation in a linear fashion, weighted by the Kalman gain  $G_{kal}$  calculated in (15c). The fourth  
 240 line updates the global model with the fused  $M_{t+1}$  calculated in (15d). Equation (15e) estimates  
 241 the variance of the fused model  $w_{t+1}$  using  $G_{kal}$  in (15d) and  $\hat{\sigma}_{t+1}^2$  in (15b), which will be used in  
 242 the next communication round. It is worth noting that all these calculations are performed on the  
 243 server, thus our method retains the same level of communication cost as FEDAVG while also being  
 244 compatible with cross-device FL settings.

245 While Equations (15a)-(15f) provide an efficient and accurate method for model updates, the variance  
 246 of the period drift  $\sigma_{Q_t}^2$  in Equation (15b) and the client drift  $\sigma_{R_t}^2$  in Equation (15c) remains unresolved.  
 247 To address this issue, we propose an effective method for estimating the variance of the period drift  
 248 and client drift. The period drift, which is a measure of the deviation from the consistency of the  
 249 optimization objective at each communication round, can be quantified by analyzing the discrepancy  
 250 between the prediction and the observation. Specifically, this can be done by computing the variance  
 251 between the momentum  $M_t$  and the average of the model updates  $\Delta \tilde{w}_t$ . Similarly, the client drift,  
 252 which represents the inconsistency of the updates made by different clients, can be estimated by  
 253 computing the variance between the average of the model updates  $\Delta \tilde{w}_t$  and the updates made by  
 254 each individual client  $\Delta \tilde{w}_t^k$ . We formulate the estimation of the variance of period drift and client  
 255 drift as follows:

$$\begin{aligned} \sigma_{Q_t}^2 &:= \frac{\sum_{i=1}^d (M_t^i - \Delta \tilde{w}_t^i)^2}{|\mathcal{S}_t|d}, \\ \sigma_{R_t}^2 &:= \frac{\sum_{k \in \mathcal{S}_t} \sum_{i=1}^d (\Delta \tilde{w}_t^{k,i} - \Delta \tilde{w}_t^i)^2}{|\mathcal{S}_t|^2 d}, \end{aligned} \quad (16)$$

256 where the index of model parameters is represented by the uppercase  $i$  and the dimension of the  
 257 model is represented by  $d$ . With the estimation of the variance of period drift and client drift, the

Table 1: Results on FEMNIST and CIFAR-100. The best method is highlighted in **bold** fonts.

Dataset	FEMNIST				CIFAR-100		
	natural	$\alpha = 1$	$\alpha = 0.1$	$\alpha = 0.01$	$\alpha = 1$	$\alpha = 0.1$	$\alpha = 0.01$
FEDAVG	$82.37 \pm 0.18$	$83.60 \pm 0.11$	$82.02 \pm 0.23$	$73.23 \pm 1.36$	$47.04 \pm 0.21$	$43.93 \pm 0.36$	$30.11 \pm 0.53$
FEDAVGM	$82.53 \pm 0.43$	$83.67 \pm 0.10$	$82.30 \pm 0.49$	$74.96 \pm 2.34$	$48.22 \pm 0.19$	$44.74 \pm 0.40$	$31.59 \pm 0.98$
FEDPROX	$82.34 \pm 0.17$	$83.58 \pm 0.11$	$82.04 \pm 0.27$	$74.16 \pm 1.19$	$46.86 \pm 0.38$	$43.74 \pm 0.27$	$30.10 \pm 0.55$
SCAFFOLD	$81.66 \pm 0.28$	$83.06 \pm 0.14$	$79.82 \pm 0.42$	$5.13 \pm 0.00$	$47.26 \pm 1.49$	$36.36 \pm 4.98$	$1.00 \pm 0.00$
FEDOPT	$5.13 \pm 0.00$	$81.86 \pm 0.38$	$78.13 \pm 0.39$	$5.13 \pm 0.00$	$47.26 \pm 1.49$	$45.43 \pm 1.18$	$32.17 \pm 1.38$
<b>FEDEVE</b>	<b><math>82.68 \pm 0.19</math></b>	<b><math>83.81 \pm 0.09</math></b>	<b><math>82.69 \pm 0.31</math></b>	<b><math>75.99 \pm 1.61</math></b>	<b><math>48.38 \pm 0.24</math></b>	<b><math>45.68 \pm 0.16</math></b>	<b><math>32.68 \pm 0.62</math></b>

Table 2: Results on MovieLens-1M. The best method is highlighted in **bold** fonts.

	AUC	HR@5	HR@10	NGCG@5	NGCG@10
FEDAVG	$0.7633 \pm 0.0065$	$0.2774 \pm 0.0100$	$0.4294 \pm 0.0120$	$0.1835 \pm 0.0058$	$0.2324 \pm 0.0064$
FEDAVGM	$0.7555 \pm 0.0128$	$0.2705 \pm 0.0384$	$0.4290 \pm 0.0196$	$0.1771 \pm 0.0319$	$0.2280 \pm 0.0257$
FEDPROX	$0.7819 \pm 0.0033$	$0.2700 \pm 0.0129$	$0.4279 \pm 0.0083$	$0.1803 \pm 0.0078$	$0.2310 \pm 0.0065$
FEDOPT	$0.7751 \pm 0.0085$	$0.2868 \pm 0.0055$	$0.4392 \pm 0.0101$	$0.1886 \pm 0.0044$	$0.2377 \pm 0.0040$
<b>FEDEVE</b>	<b><math>0.7967 \pm 0.0016</math></b>	<b><math>0.2916 \pm 0.0077</math></b>	<b><math>0.4460 \pm 0.0088</math></b>	<b><math>0.1924 \pm 0.0039</math></b>	<b><math>0.2407 \pm 0.0037</math></b>

258 overall process of model update can be described in Algorithm 1. The fundamental principle of  
259 FEDEVE is to calculate the Kalman gain  $G_{kal}$  which is used to determine the relative weight of the  
260 prediction and observation when they are combined. The value of  $G_{kal}$  is calculated based on the  
261 ratio of the variance of the prediction  $\sigma_{Q_t}^2$  and the variance of the observation  $\sigma_{R_t}^2$ . This coefficient is  
262 used to adjust the update direction of the model. A small  $G_{kal}$  means that the observation is close to  
263 the prediction, hence the update direction will also be close to the prediction. A large  $G_{kal}$  means  
264 that the observation deviates significantly from the prediction, hence the update direction will deviate  
265 from the prediction and be closer to the observation. This allows the algorithm to adapt to different  
266 scenarios in which the observations may deviate more or less from the predictions.

### 267 3 Experiments

#### 268 3.1 Setup

269 **Datasets and models.** We evaluate FEDEVE on three computer vision (CV) and recommender  
270 system (RS) datasets under realistic cross-device FL settings. **For CV dataset**, we use FEMNIST<sup>2</sup>  
271 Caldas et al. [2018], consisting of 671,585 training examples and 77,483 test samples of 62 different  
272 classes including 10 digits, 26 lowercase and 26 uppercase images with 28x28 pixels, handwritten by  
273 3400 users. We also use CIFAR-10/100<sup>3</sup> Caldas et al. [2018], consisting of 50,000 training examples  
274 and 10,000 test samples of 10/100 different classes with 32x32 pixels. For FEMNIST dataset, we use  
275 the lightweight model LeNet5 LeCun et al. [1998] and for CIFAR-10/100 dataset, we use ResNet-18  
276 (replacing batch norm with group norm [Hsieh et al., 2020, Reddi et al., 2020]). **For RS dataset**,  
277 we use MovieLens 1M<sup>4</sup> Harper and Konstan [2015], including 1,000,209 ratings by unidentifiable  
278 6,040 users on 3,706 movies. It is a click-through rate (CTR) task, and we use the popular DIN Zhou  
279 et al. [2018] model. For performance evaluation, we follow a widely used leave-one-out protocol  
280 Muhammad et al. [2020]. For each user, we hold out their latest interaction as testset and use the  
281 remaining data as trainset, and binarize the user feedback where all ratings are converted to 1, and  
282 negative instances are sampled 4:1 for training and 99:1 for test times the number of positive ones.

283 **Federated learning settings.** It is important to note that the datasets FEMNIST and MovieLens 1M  
284 have a "natural" non-iid distribution, which means that the data is split by "user\_id". For example,  
285 in FEMNIST, images are handwritten by different users, and in MovieLens 1M, movies are rated  
286 by different users. Furthermore, we use the Dirichlet distribution, to simulate the label distribution  
287 skew setting for FEMNIST, as described in Hsu et al. [2019]. This distribution allows us to control  
288 the degree of heterogeneity by adjusting the hyperparameter  $\alpha$  (the smaller, the more non-iid). This  
289 allows us to test the robustness of the algorithm under different levels of heterogeneity, which is  
290 a common scenario in real-world FL settings. For the FL training, we set a total of  $T = 1500$   
291 communication rounds for the CV task and sample 10 clients per round with SGD optimizer. For the  
292 RS task, we set a total of  $T = 1000$  communication rounds and sample 20 clients per round with

<sup>2</sup><https://github.com/TalwalkarLab/leaf/tree/master/data/fmnist>

<sup>3</sup><https://www.cs.toronto.edu/~kriz/cifar.html>

<sup>4</sup><https://grouplens.org/datasets/movielens/>

Table 3: Results on FEMNIST with different  $\alpha$  and  $E$ . The best method is highlighted in **bold** fonts.

Method	Natural			$\alpha = 1$			$\alpha = 0.1$			$\alpha = 0.01$		
	$E = 1$	$E = 3$	$E = 5$	$E = 1$	$E = 3$	$E = 5$	$E = 1$	$E = 3$	$E = 5$	$E = 1$	$E = 3$	$E = 5$
FEDAVG	82.46 $\pm$ 0.18	81.95 $\pm$ 0.26	66.38 $\pm$ 30.63	83.64 $\pm$ 0.11	83.57 $\pm$ 0.12	83.38 $\pm$ 0.09	82.12 $\pm$ 0.23	81.91 $\pm$ 0.23	81.70 $\pm$ 0.26	74.51 $\pm$ 1.36	73.43 $\pm$ 1.73	72.67 $\pm$ 1.39
FEDAVGM	82.55 $\pm$ 0.43	81.99 $\pm$ 0.42	50.97 $\pm$ 37.43	83.67 $\pm$ 0.10	83.79 $\pm$ 0.11	83.65 $\pm$ 0.08	82.36 $\pm$ 0.49	82.23 $\pm$ 0.26	82.16 $\pm$ 0.34	75.18 $\pm$ 2.34	74.20 $\pm$ 2.81	73.79 $\pm$ 3.13
FEDPROX	82.43 $\pm$ 0.17	81.90 $\pm$ 0.27	51.07 $\pm$ 37.51	83.62 $\pm$ 0.11	83.52 $\pm$ 0.14	67.65 $\pm$ 31.26	82.17 $\pm$ 0.27	82.12 $\pm$ 0.19	81.94 $\pm$ 0.30	75.07 $\pm$ 1.19	74.38 $\pm$ 1.46	60.22 $\pm$ 27.56
SCAFOLD	81.66 $\pm$ 0.28	81.08 $\pm$ 0.36	80.76 $\pm$ 0.30	83.18 $\pm$ 0.14	82.75 $\pm$ 0.16	82.46 $\pm$ 0.16	79.82 $\pm$ 0.42	79.10 $\pm$ 0.65	78.44 $\pm$ 0.70	5.13 $\pm$ 0.00	5.13 $\pm$ 0.00	5.13 $\pm$ 0.00
FEDOPT	5.13 $\pm$ 0.00	5.13 $\pm$ 0.00	5.13 $\pm$ 0.00	81.86 $\pm$ 0.58	55.90 $\pm$ 37.69	55.66 $\pm$ 37.39	78.13 $\pm$ 0.39	5.13 $\pm$ 0.00				
FEDEVE	<b>82.66 <math>\pm</math> 0.19</b>	<b>82.20 <math>\pm</math> 0.20</b>	<b>81.93 <math>\pm</math> 0.16</b>	<b>83.81 <math>\pm</math> 0.09</b>	<b>83.88 <math>\pm</math> 0.08</b>	<b>83.72 <math>\pm</math> 0.05</b>	<b>82.69 <math>\pm</math> 0.31</b>	<b>82.66 <math>\pm</math> 0.18</b>	<b>82.52 <math>\pm</math> 0.19</b>	<b>75.99 <math>\pm</math> 1.61</b>	<b>75.00 <math>\pm</math> 2.24</b>	<b>74.56 <math>\pm</math> 2.05</b>

293 Adam optimizer Kingma and Ba [2014]. In all datasets, each client trains for  $E = 1$  epoch at the  
294 local update with a learning rate of  $\eta_l = 0.01$ . In our proposed FEDEVE , we set the global learning  
295 rate  $\eta_g = 1$  for all experiments.

296 **Baselines.** To evaluate the performance of FEDEVE , we compare it with several state-of-the-art  
297 FL methods: 1) The vanilla FL method FEDAVG McMahan et al. [2017], which is a widely used  
298 method for FL; 2) A client-side FL method FEDPROX Li et al. [2020b], which improves the  
299 model aggregation by adding a proximal term to the local update; 3) A server-side FL method  
300 FEDAVGM Hsu et al. [2019], which adapts the momentum in FL optimization; 4) A server-side  
301 FL method FEDOPT Reddi et al. [2020], which introduces adaptive optimization methods in FL.  
302 See more experimental details in the Appendix E.

### 303 3.2 Analysis

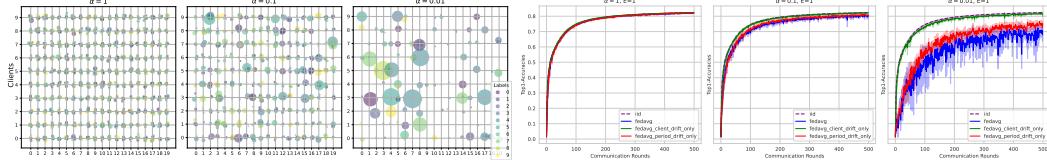


Figure 3: **Visualization of period drift and its impact on performance** (a) **Visualization of period drift** The color of the scatter points represents different classes, and the size denotes the number of samples of a given class on a particular client. (b) **Visualization of its impact on performance** It is revealed in cross-device FL when data is rather non-iid, period drift has a greater effect than client drift. Appendix E.2 for setting details.

304 **Visualizing the period drift and its impact.** Figure 3 (a) visualizes the data distributions of these  
305 sampled clients. Client drift arises due to the shift in label distribution among sampled clients **within**  
306 **a single round**, while period drift results from the shift in the data distribution of participating clients  
307 **across different rounds**. The scatter points’ size and distribution grow more diverse both within  
308 and across communication rounds as the value of  $\alpha$  decreases (indicating increasing non-iid). The  
309 implications of these drifts on the global model’s convergence are presented in Figure 3 (b). Utilizing  
310 the vanilla FEDAVG algorithm for illustration, we experimented with four settings: 1) FEDAVG with  
311 iid data; 2) FEDAVG experiencing only period drift; 3) FEDAVG subject to only client drift; and  
312 4) FEDAVG impacted by both drifts (See appendix E.2 for detailed settings). As heterogeneity  
313 intensifies, the effects of both drifts become evident. Specifically, in a highly non-iid environment  
314 ( $\alpha = 0.01$ ), FEDAVG affected only by client drift yields results akin to the iid setting. In contrast,  
315 FEDAVG influenced solely by period drift significantly disrupts the stability and convergence of the FL  
316 process. The combination of both drifts results in the poorest performance, underlining that in cross-  
317 device FL, period drift poses a more considerable challenge to model convergence than client drift.

318 **The performance of FEDEVE .** We evaluate our algorithm on real-world datasets and compare it with  
319 the relevant state-of-the-art methods in Tables 1 and 2. We conducted simulations on three datasets:  
320 FEMNIST, CIFAR-100, and MovieLens. The FEMNIST and MovieLens datasets have a naturally-  
321 arising client partitioning setting in real-world FL scenarios, making them highly representative. For  
322 FEMNIST and CIFAR-100 datasets, each of the datasets includes three non-iid settings, established  
323 through the Dirichlet distribution partition method [Hsu et al., 2019]. *Generally, the results show*  
324 *that our proposed algorithm, FEDEVE , consistently outperforms the baselines, and the performance*  
325 *gains are more dominant in more non-iid settings ( $\alpha = 0.01$ )*. We also conduct experiments with  
326 different local epochs ( $E$ ), please refer to the Figure 3.2 for the setting details. Also, our FEDEVE has  
327 more leading advantages in RS experiments, indicating its large potential in real-world industrial  
328 applications. Our method can better utilize the server-side adaptation through the Bayesian filter’s  
329 predict-observer framework. Besides, it is important to note that our method does not introduce other

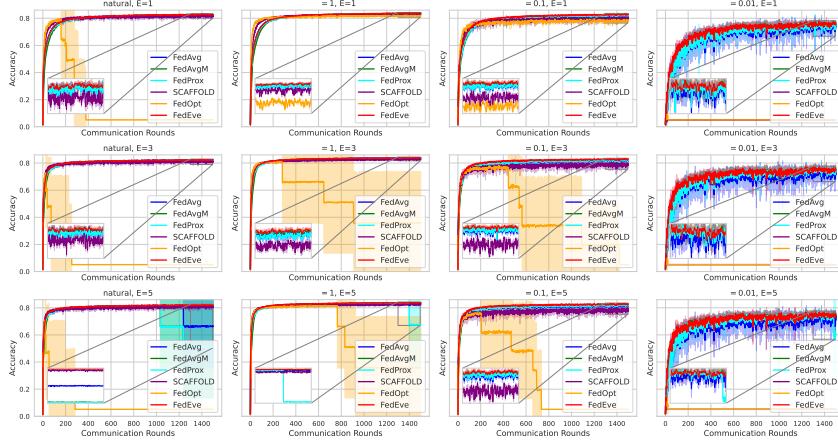


Figure 4: Accuracies with different  $\alpha$  and  $E$

hyperparameters while these baselines have multiple hyperparameters to tune, which means that our FEDEVE is more flexible and advantageous in real-world practices.

**Performance of FEDEVE on FEMNIST with different local epochs.** Table 3 showcases the results on the FEMNIST dataset across various methods, specifically: FEDAVG , FEDAVGM , FEDPROX , SCAFFOLD , FEOPT , and FEDEVE , with different settings of parameters  $\alpha$  and  $E$ . FEDEVE method stands out consistently as the superior approach across all configurations. This consistent performance signifies that FEDEVE is a potent and reliable method for the FEMNIST dataset across the tested configurations. The performance drop of SCAFFOLD in specific experiments may be attributed to two primary reasons: Staleness of control variate: SCAFFOLD mandates that each client maintain a control variant. However, given the large number of clients and the fact that only a limited subset is chosen for training during each communication round, most control variants remain outdated. As a result, they fail to effectively correct the bias in local updates. This point was also reported in FedOpt [Reddi et al., 2020]. Excessiveness of correction: Upon detailed inspection of our experiments, we discerned that the training of SCAFFOLD tends to fail when there exists a client with more substantial data than others. This stems from the fact that the fixed batch size and training epoch will result in more local updates in the clients with more data, but it will be corrected by the same control variant in SCAFFOLD. Excessive corrections drive the model further from the optimal point, resulting in the divergence of the model. We reckon that the poor performances of FedOpt in some settings primarily result from period drift. Period drift impedes FedOpt’s adaptivity across rounds. FedOpt tailors the learning rates of individual weights by accumulating past gradients’ squares. However, with the ever-shifting optimization objectives in each communication round (period drift), these rate adjustments become misaligned for subsequent updates, thereby skewing model training. It is validated in the experiments that FedOpt fails on FEMNIST with natural and heterogeneity, where period drift is more dominant (more client number, more non-i.i.d. data).

## 4 Conclusion

In this work, we conducted a comprehensive exploration of the impact of client drift and period drift on the performance of cross-device FL, discovering that period drift can be particularly harmful as data heterogeneity increases. To solve this challenge, we introduced a novel predict-observe framework and a method, FedEve, that views these drifts as noise associated with prediction and observation. By integrating these two sources in a principled way through a Kalman filter-based approach, we provided a better estimation of model update steps, reducing variance and improving the stability and convergence speed of FL. Our theoretical analysis showed that FedEve can effectively bound the variance of model updates, while extensive empirical evaluations demonstrated that it significantly outperforms alternative methods across different tasks and datasets. We believe our framework opens up new possibilities for developing more robust and efficient FL algorithms.

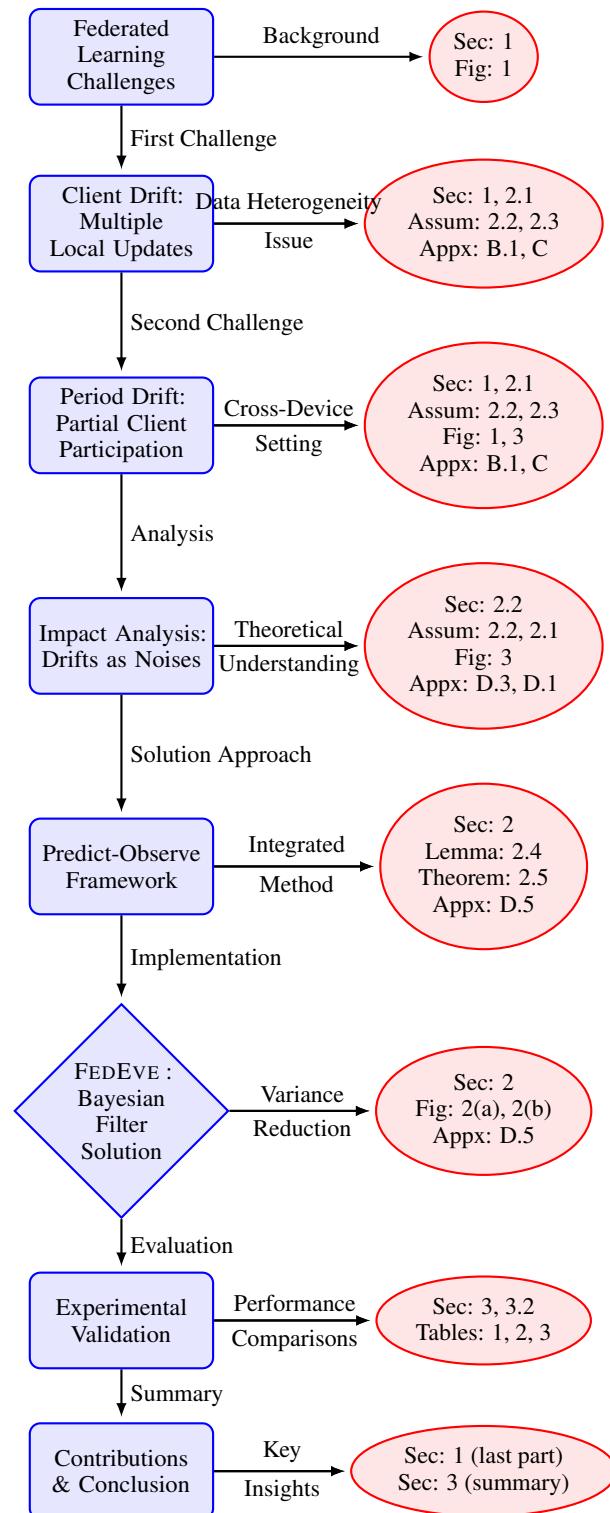
365 **References**

- 366 Barak Battash and Ofir Lindenbaum. Revisiting the noise model of stochastic gradient descent. *arXiv  
367 preprint arXiv:2303.02749*, 2023.
- 368 Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H Brendan McMahan,  
369 Virginia Smith, and Ameet Talwalkar. Leaf: A benchmark for federated settings. *arXiv  
370 preprint arXiv:1812.01097*, 2018.
- 371 Hong-You Chen and Wei-Lun Chao. On bridging generic and personalized federated learning for  
372 image classification. In *International Conference on Learning Representations*, 2021.
- 373 Yae Jee Cho, Jianyu Wang, and Gauri Joshi. Towards understanding biased client selection in  
374 federated learning. In *International Conference on Artificial Intelligence and Statistics*, pages  
375 10351–10375. PMLR, 2022.
- 376 Yann Fraboni, Richard Vidal, Laetitia Kameni, and Marco Lorenzi. A general theory for client  
377 sampling in federated learning. In *Trustworthy Federated Learning: First International Workshop,  
378 FL 2022, Held in Conjunction with IJCAI 2022, Vienna, Austria, July 23, 2022, Revised Selected  
379 Papers*, pages 46–58. Springer, 2023.
- 380 Yongxin Guo, Tao Lin, and Xiaoying Tang. Towards federated learning on time-evolving heteroge-  
381 neous data. *arXiv preprint arXiv:2112.13246*, 2021.
- 382 F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *Acm  
383 transactions on interactive intelligent systems (tiis)*, 5(4):1–19, 2015.
- 384 Kevin Hsieh, Amar Phanishayee, Onur Mutlu, and Phillip Gibbons. The non-iid data quagmire  
385 of decentralized machine learning. In *International Conference on Machine Learning*, pages  
386 4387–4398. PMLR, 2020.
- 387 Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data  
388 distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.
- 389 Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin  
390 Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Ad-  
391 vances and open problems in federated learning. *Foundations and Trends® in Machine Learning*,  
392 14(1–2):1–210, 2021.
- 393 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J. Reddi, Sebastian U. Stich, and  
394 Ananda Theertha Suresh. SCAFFOLD: Stochastic Controlled Averaging for Federated Learning.  
395 *arXiv:1910.06378 [cs, math, stat]*, 2021.
- 396 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint  
397 arXiv:1412.6980*, 2014.
- 398 Hendrik Anthony Kramers. Brownian motion in a field of force and the diffusion model of chemical  
399 reactions. *Physica*, 7(4):284–304, 1940.
- 400 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to  
401 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 402 Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges,  
403 methods, and future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020a.
- 404 Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith.  
405 Federated Optimization in Heterogeneous Networks. *arXiv:1812.06127 [cs, stat]*, 2020b.
- 406 Tao Lin, Sebastian U Stich, Kumar Kshitij Patel, and Martin Jaggi. Don't use large mini-batches, use  
407 local sgd. *arXiv preprint arXiv:1808.07217*, 2018.
- 408 Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust model  
409 fusion in federated learning. *Advances in Neural Information Processing Systems*, 33:2351–2363,  
410 2020.

- 411 Jarl Waldemar Lindeberg. Eine neue herleitung des exponentialgesetzes in der wahrscheinlichkeit-  
 412 srechnung. *Mathematische Zeitschrift*, 15(1):211–225, 1922.
- 413 Stephan Mandt, Matthew D Hoffman, and David M Blei. Stochastic gradient descent as approximate  
 414 bayesian inference. *Journal of Machine Learning Research*, 18:1–35, 2017.
- 415 H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas.  
 416 Communication-Efficient Learning of Deep Networks from Decentralized Data. *arXiv:1602.05629*  
 417 [cs], 2017.
- 418 Khalil Muhammad, Qin Qin Wang, Diarmuid O’Reilly-Morgan, Elias Tragos, Barry Smyth, Neil  
 419 Hurley, James Geraci, and Aonghus Lawlor. Fedfast: Going beyond average for faster training  
 420 of federated recommender systems. In *Proceedings of the 26th ACM SIGKDD International  
 Conference on Knowledge Discovery & Data Mining*, pages 1234–1242, 2020.
- 422 Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný,  
 423 Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. *arXiv preprint  
 arXiv:2003.00295*, 2020.
- 425 Umut Simsekli, Levent Sagun, and Mert Gurbuzbalaban. A tail-index analysis of stochastic gradient  
 426 noise in deep neural networks. In *International Conference on Machine Learning*, pages 5827–  
 427 5837. PMLR, 2019.
- 428 Zhenyu Sun, Xiaochun Niu, and Ermin Wei. Understanding generalization of federated learning via  
 429 stability: Heterogeneity matters. *arXiv preprint arXiv:2306.03824*, 2023.
- 430 Jingfeng Wu, Wenqing Hu, Haoyi Xiong, Jun Huan, Vladimir Braverman, and Zhanxing Zhu. On the  
 431 noisy gradient descent that generalizes as sgd. In *International Conference on Machine Learning*,  
 432 pages 10367–10376. PMLR, 2020.
- 433 Yuhuai Wu, Mengye Ren, Renjie Liao, and Roger Grosse. Understanding short-horizon bias in  
 434 stochastic meta-optimization. *arXiv preprint arXiv:1803.02021*, 2018.
- 435 Zeke Xie, Issei Sato, and Masashi Sugiyama. A diffusion theory for deep learning dynamics:  
 436 Stochastic gradient descent exponentially favors flat minima. In *9th International Conference on  
 437 Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net,  
 438 2021. URL [https://openreview.net/forum?id=wXgk\\_iCiYGo](https://openreview.net/forum?id=wXgk_iCiYGo).
- 439 Xin Yao, Tianchi Huang, Rui-Xiao Zhang, Ruiyu Li, and Lifeng Sun. Federated Learning with  
 440 Unbiased Gradient Aggregation and Controllable Meta Updating. *CoRR*, abs/1910.08234, 2019.
- 441 Michael Zhang, James Lucas, Jimmy Ba, and Geoffrey E Hinton. Lookahead optimizer: k steps  
 442 forward, 1 step back. *Advances in neural information processing systems*, 32, 2019.
- 443 Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated  
 444 Learning with Non-IID Data. *arXiv:1806.00582* [cs, stat], 2018.
- 445 Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin,  
 446 Han Li, and Kun Gai. Deep interest network for click-through rate prediction. In *Proceedings of  
 447 the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pages  
 448 1059–1068, 2018.
- 449 Chen Zhu, Zheng Xu, Mingqing Chen, Jakub Konečný, Andrew Hard, and Tom Goldstein. Diurnal  
 450 or nocturnal? federated learning of multi-branch networks from periodically shifting distributions.  
 451 In *International Conference on Learning Representations*, 2022.
- 452 Zhanxing Zhu, Jingfeng Wu, Bing Yu, Lei Wu, and Jinwen Ma. The anisotropic noise in stochastic  
 453 gradient descent: Its behavior of escaping from sharp minima and regularization effects. In  
 454 *International Conference on Machine Learning*, pages 7654–7663. PMLR, 2019.
- 455 Liu Ziyin, Kangqiao Liu, Takashi Mori, and Masahito Ueda. Strength of minibatch noise in sgd.  
 456 *arXiv preprint arXiv:2102.05375*, 2021.

457 **A Appendix**

458 **Reading Guide**



460 This reading guide provides a comprehensive roadmap to navigate through our paper on addressing  
461 the challenges in Federated Learning (FL).  
462 We begin by introducing the fundamental challenges of FL in Section 1, illustrated in Figure 1. Here,  
463 we identify two critical issues: Client Drift and Period Drift.  
464 Client Drift occurs due to multiple local updates and data heterogeneity, detailed in Sections 1 and  
465 2.1. This concept is formalized in Assumptions 2.2 and 2.3, with further explanations in Appendices  
466 B.1 and C.  
467 Period Drift, a less studied but equally important challenge, arises from partial client participation  
468 in cross-device settings. This phenomenon is thoroughly discussed in Sections 1 and 2.1, visually  
469 represented in Figures 1 and 3, and conceptually explained in Appendices B.1 and C.  
470 The impact of these drifts on FL performance is analyzed in Section 2.2, where we formally define  
471 them in Definition 2.1 and model them as noise under Assumption 2.2. The detrimental effects are  
472 visualized in Figure 3, with additional theoretical analysis in Appendix D.3 and justification of the  
473 noise assumption in Appendix D.1.  
474 To address these challenges, we propose a Predict-Observe Framework in Section 2. This framework  
475 leverages Lemma 2.4 and Theorem 2.5 to establish a principled approach for integrating predictions  
476 and observations, with detailed Bayesian derivations in Appendix D.5.  
477 Our specific solution, FEDEVE , implements this framework using a Bayesian filter. The methodology  
478 is presented in Section 2 with illustrations in Figures 2(a) and 2(b), demonstrating how it achieves  
479 variance reduction through the fusion of predictions and observations.  
480 We validate our approach through extensive experiments in Sections 3 and 3.2. The results, presented  
481 in Tables 1, 2, and 3, demonstrate that FEDEVE consistently outperforms state-of-the-art methods,  
482 particularly in highly heterogeneous settings.  
483 The paper concludes by summarizing our key contributions in the final part of Section 1 and providing  
484 a comprehensive overview of our findings in Section 3.  
485 This guide is designed to facilitate a clear understanding of our research, highlighting the logical  
486 progression from identifying challenges to proposing and validating solutions in the context of  
487 Federated Learning.

488 **B Related works**

489 There are many works that have attempted to address the non-iid problem in federated learning.  
 490 FEDAVG , first presented by McMahan et al. [2017], has been demonstrated to have issues with  
 491 convergence when working with non-iid data. Zhao et al. [2018] depict the non-iid trap as weight  
 492 divergence, and it can be reduced by sharing a small set of data. However, in traditional federal setting,  
 493 data sharing violates the principle of data privacy. Karimireddy et al. [2021] highlight the phenomenon  
 494 of “client drift” that occurs when data is heterogeneous (non-iid), and uses control variates to address  
 495 this problem. However, using Scaffold in cross-device FL may not be effective, as it requires clients  
 496 to maintain the control variates, which may become outdated and negatively impact performance.  
 497 Li et al. [2020b] propose FedProx that utilizes a proximal term to deal with heterogeneity.

498 In addition to these works, some research has noticed the presence of period drift, but have not  
 499 specifically addressed it in their analysis. For example, Cho et al. [2022], Fraboni et al. [2023]  
 500 investigate the problem of biased client sampling and proposes an sampling strategy that selects clients  
 501 with large loss. However, active client sampling can potentially alter the overall data distribution  
 502 by having unrandom clients participation, which can raise concerns about fairness. Similarly, Yao  
 503 et al. [2019] propose a meta-learning based method for unbiased aggregation, but it requires training  
 504 the global model on a proxy dataset, which may not be feasible in certain scenarios where such a  
 505 dataset is not available. Zhu et al. [2022] observe that the data on clients have periodically shifting  
 506 distributions that changed with the time of day, and model it using a mixture of distributions that  
 507 gradually shifted between daytime and nighttime modes. Guo et al. [2021] study the impact of  
 508 time-evolving heterogeneous data in real-world scenarios, and solve it in a framework of continual  
 509 learning. Although these two papers define similar terms, they focus on the case of client data  
 510 changing over time. However, in this paper, we find that even if the distribution of client data remains  
 511 unchanged, period drift can seriously affect the convergence of FL.

512 **B.1 The concepts of *Drifts***

513 To distinguish between the concepts of period drift and client drift, we conduct a thorough analysis of  
 514 the entire FL process, as shown in the upper part of Fig 5. We utilize a chain of sampling (from left  
 515 to right) to identify the errors introduced at each stage of FL. We aim to optimize an ideal objective  
 516  $f(x)$  to obtain the ideal optima  $x^*$  under the assumption of infinite data.. However, in the real world,  
 517 we can only optimize  $f_{\mathcal{N}}(x)$  over a finite training set  $\mathcal{D}_{\mathcal{N}}$  to obtain the optima  $x_{\mathcal{N}}^*$ . Sampling finite  
 518 data from infinite data introduces the first variance  $\sigma_{\mathcal{N}}^2$ , also known as the generalization error, which  
 519 is small because we usually assume that the training data and the overall data are IID in the context  
 520 of machine learning.

521 In FL, directly optimize  $f_{\mathcal{N}}(x)$  is not feasible due to the non-IID distribution of data across different  
 522 clients, requiring distributed training instead. In cross-device FL, only a small subset of clients is  
 523 sampled each round for training. However, due to data heterogeneity, the data distribution  $\mathcal{D}_S$  of  
 524 the sampled clients differs from the full client set  $\mathcal{D}_{\mathcal{N}}$ , leading to different optimization objectives  
 525  $f_S(x)$  and  $f_{\mathcal{N}}(x)$ , as well as a significant deviation between the optima  $x_S^*$  and  $x_{\mathcal{N}}^*$ . This difference  
 526 introduces the second variance  $\sigma_S^2$  referred to as “period drift”, which can be substantial in scenarios  
 527 with strong data heterogeneity and significantly slow down model convergence, as depicted in Fig 3  
 528 in the main paper. Additionally, these sampled clients independently optimize their local objectives  
 529  $f_k(x)$  on their respective datasets  $\mathcal{D}_k$ . Due to the presence of data heterogeneity, the optima  $x_k^*$  and  
 530  $x_S^*$  obtained by different clients also differ, introducing the third variance  $\sigma_k^2$ , which is “client drift”.  
 531 On this chain of sampling, each sampling introduces an error, each affecting the model’s convergence  
 532 and performance.

533 To clarify the concepts of “period drift” and “client drift” in FL, we refine the commonly cited  
 534 inequality  $\frac{1}{S} \sum_{i \in S} x_i^* \neq x^*$ . We propose a more nuanced formulation:  $x^* \neq x_{\mathcal{N}}^* \neq x_S^* \neq \bar{x} =$   
 535  $\frac{1}{|S|} \sum x_k^* \neq x_k^*$ . In FL, the goal is to estimate the global optimum  $x_{\mathcal{N}}^*$ —or ideally,  $x^*$ —through the  
 536 weighted average  $\bar{x} = \frac{1}{|S|} \sum x_k^*$ . However, this average  $\bar{x}$  often only approximates  $x_S^*$  and falls short  
 537 of estimating the true global optimum  $x_{\mathcal{N}}^*$  due to significant distributional discrepancies between  $\mathcal{D}_S$   
 538 and  $\mathcal{D}_{\mathcal{N}}$ . This variance is highlighted in the lower part of Figure 5 and illustrates the challenges in  
 539 drawing reliable inferences about the global data distribution from locally sampled subsets.

540 Given that we only have data from current samples and lack comprehensive information about  
 541 the broader sampling framework, it’s clear that attributing the inequality  $\frac{1}{S} \sum_{i \in S} x_i^* \neq x^*$  solely  
 542 to client drift is misleading. This inequality encompasses effects beyond local updates, including

Figure 1: A Chain of Sampling

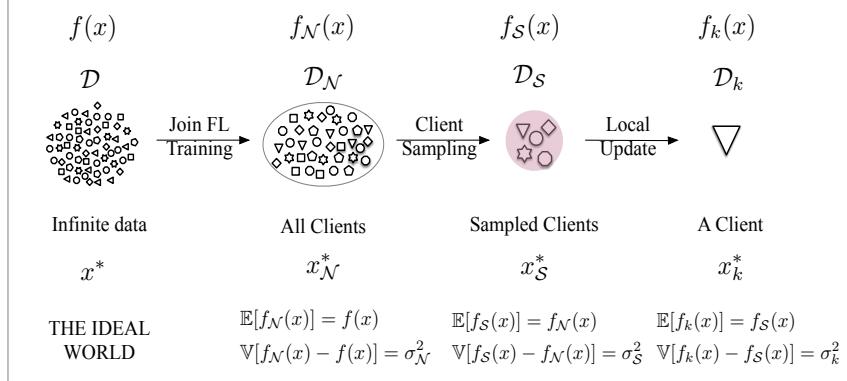


Figure 2: A Gap from  $x_S^*$  to  $x_N^*$

$$x^* \neq x_N^* \neq x_S^* \neq \bar{x} = \frac{1}{|\mathcal{S}|} \sum x_k^* \neq x_S^*$$

Client Drift

Period Drift

Figure 5: **Illustration of period drift and client drift in the FL process.**

543 generalization error  $\sigma_N^2$  and client sampling  $\sigma_S^2$ . A precise representation of client drift would  
 544 therefore be  $\bar{x} = \frac{1}{|\mathcal{S}|} \sum x_k^* \neq x_S^*$ . An illustrative example of this is when each client performs a  
 545 single gradient update step; there is effectively no client drift, as the average of these updates reflects  
 546 the gradient of the current batch. Yet, the disparity  $x^* \neq x_N^* \neq x_S^* = \bar{x} = \frac{1}{|\mathcal{S}|} \sum x_k^*$  still persists,  
 547 and the influence of period drift remains. Hence, period drift is fundamentally distinct from client  
 548 drift and operates independently.

## 549 C The difference between *Drift* in FL and the *Noise* of centralized SGD

550 Federated learning possesses unique characteristics compared to traditional centralized optimization,  
 551 such as client sampling, multiple local epochs, and non-iid data distribution. In this context, drifts  
 552 in federated learning can be viewed as noises to the training dynamics. More specifically, period  
 553 drift, originating from non-iid data and partial participation (only a subset of clients participate  
 554 in each round), can be likened to the implementation of a mini-batch technique in the full-batch  
 555 gradient descent of centralized training [Ziyin et al., 2021]. Here, the distinction is that in centralized  
 556 optimization issues, each batch is iid, and each batch's data distribution closely mirrors the overall  
 557 distribution, albeit with a noise component. This noise becomes remarkably pronounced in federated  
 558 learning, given that client data is non-iid. Client drift, arising from non-iid data and multiple local  
 559 updates (where each client runs local SGD with multiple steps), is a well-structured noise [Lin et al.,  
 560 2018]. Due to the combined impact of client drift and period drift, the situation can be perceived as  
 561 adding a noise term to the original model (or gradient). The research outlined in this paper based on  
 562 the difference between the drifts in federated learning and the noises in centralized SGD.

### 563 C.1 Visualization of difference between Period Drift and noise in SGD

564 In this section, we provide a visualization to illustrate the fundamental difference between period drift  
 565 in cross-device federated learning and the noise introduced by stochastic gradient descent (SGD).  
 566 While both phenomena introduce stochasticity into the optimization process, their nature and impact  
 567 on model convergence differ significantly.

568 Figure 6 depicts this distinction by comparing the gradient directions across different communication  
 569 rounds in federated learning. Period drift arises from the systematic shifts in client data distributions  
 570 between communication rounds, creating consistent biases in gradient direction for extended periods.  
 571 In contrast, SGD noise exhibits more random fluctuations that tend to average out over time.

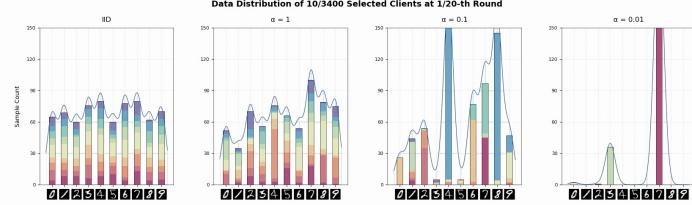


Figure 6: Visualization of the difference between period drift and SGD noise

572 This visualization highlights why period drift poses a more significant challenge to federated learning  
 573 than standard SGD noise. While SGD noise is often assumed to have zero mean and can be mitigated  
 574 through increased batch sizes or adaptive optimization methods, period drift introduces different kind  
 575 of noise that depends on data heterogeneity that require specialized techniques like our proposed  
 576 FEDIVE algorithm to effectively address. The systematic nature of period drift makes it particularly  
 577 problematic in highly heterogeneous cross-device settings, as demonstrated in our experimental  
 578 results in Section 3.

## 579 C.2 A simple analysis of data heterogeneity and period drift

580 Let us analyze how data heterogeneity impacts the data distribution difference between sampled  
 581 clients and the overall population. We begin by examining the expectation of the sampled distribution  
 582  $P_S$ . Since clients are randomly selected, the expectation equals the overall data distribution:

$$\mathbb{E}[P_S] = \mathbb{E}\left[\frac{1}{|S|} \sum_{i \in S} D_i\right] = \frac{1}{|S|} \sum_{i=1}^N D_i = P \quad (17)$$

583 Next, we calculate the variance of  $P_S$ . Using the linearity of variance and independence:

$$\begin{aligned} \text{Var}(P_S) &= \text{Var}\left(\frac{1}{|S|} \sum_{i \in S} D_i\right) \\ &= \frac{1}{|S|^2} \sum_{i \in S} \text{Var}(D_i) \\ &= \frac{1}{|S|^2} \sum_{i \in S} \frac{1}{M} \sum_{j=1}^M (p_{ij} - p_j)^2 \end{aligned} \quad (18)$$

584 Since each client in  $S$  is independently and identically distributed, taking the expectation yields:

$$\mathbb{E}[\text{Var}(P_S)] = \frac{1}{|S|} \cdot \frac{1}{N} \sum_{i=1}^N \text{Var}(D_i) = \frac{H}{|S|} \quad (19)$$

585 Finally, we derive the expectation of the data distribution difference between selected clients:

$$\begin{aligned} \mathbb{E}[D_S] &= \mathbb{E}[\text{Var}(P_S, P)] \\ &= \mathbb{E}\left[\frac{1}{M} \sum_{j=1}^M \left(\frac{1}{|S|} \sum_{i \in S} p_{ij} - p_j\right)^2\right] \\ &= \frac{1}{M} \sum_{j=1}^M \mathbb{E}\left[\left(\frac{1}{|S|} \sum_{i \in S} p_{ij} - p_j\right)^2\right] \\ &= \text{Var}(P_S) \\ &= \frac{H}{|S|} \end{aligned} \quad (20)$$

586 This result demonstrates that the expected data distribution difference among selected clients is  
 587 directly proportional to the system's data heterogeneity  $H$  and inversely proportional to the number  
 588 of selected clients  $|S|$ . This finding has two important implications:

- 589 1. As data heterogeneity  $H$  increases, the expected data distribution difference among selected  
 590 clients also increases. In the IID case where  $H = 0$ , we have  $\mathbb{E}[D_S] = 0$ , indicating that  
 591 any client selection will match the overall distribution.
- 592 2. As the number of selected clients  $|S|$  increases, the expected data distribution difference  
 593 decreases, suggesting that increasing participation can help mitigate the impact of data  
 594 heterogeneity.

## 595 D The relationship between period drift and client drift

596 The relationship between period drift and client drift can be understood through their distinct yet  
 597 interrelated effects on the FL optimization process. First, there is a clear sequential nature to these  
 598 drifts - period drift occurs first in the optimization process, as it is determined by client selection at  
 599 the beginning of each round. Client drift then follows as a consequence of local updates performed by  
 600 the selected clients. This sequential relationship means that period drift can influence the magnitude  
 601 and direction of subsequent client drift.

602 When data heterogeneity is high, period drift leads to biased optimization objectives for each round.  
 603 This bias can be amplified by client drift during local updates, as clients further optimize towards  
 604 their local objectives. The combination of these drifts results in increased variance in model updates,  
 605 slower convergence rates, and potential divergence from the global optimum. Despite their potentially  
 606 harmful effects, these drifts can sometimes compensate for each other. Period drift may help prevent  
 607 over-fitting to the data distribution of currently selected clients, while client drift can help explore the  
 608 local optimization landscape more thoroughly. Additionally, the interaction between these drifts can  
 609 create a natural regularization effect. The magnitude of these drifts operates at different scales, which  
 610 can be expressed mathematically as:

$$\begin{aligned} \text{Period Drift} &\propto \frac{H}{|S|} \\ \text{Client Drift} &\propto K \cdot H \end{aligned} \tag{21}$$

611 where  $H$  is the heterogeneity measure,  $|S|$  is the number of selected clients, and  $K$  is the number of  
 612 local steps. Understanding this relationship is crucial for designing effective FL algorithms that can  
 613 handle both types of drift simultaneously, as addressing one type of drift in isolation may inadvertently  
 614 exacerbate the other.

### 615 D.1 The justification of Gaussian-like noise assumption 2.2

616 **Assumption D.1** (2.2). *The aggregated model parameters on the server  $w_{server}$ , can be represented  
 617 as the sum of the optimal parameters  $w^*$  and a drift (noise) that follows a normal distribution  
 618  $w_{drift} \sim \mathcal{N}(0, \sigma_{drift}^2)$ :*

$$w_{server} = w^* + w_{drift} \leftarrow \text{noise}, \tag{22}$$

619 where  $w^*$  represents the optimal parameters obtained through the use of stochastic gradient descent  
 620 (SGD),  $w_{drift}$  represents the noise term caused by factors such as client sampling, multiple local  
 621 epochs, and non-iid data distribution that we assume a normal distribution, and  $w_{server}$  represents  
 622 the aggregated model parameters also follows a normal distribution  $w_{server} \sim \mathcal{N}(w^*, \sigma_{drift}^2)$ ,  
 623 with the expectation of the aggregate model parameters being equal to the optimal parameters, i.e.  
 624  $\mathbb{E}[w_{server}] = w^*$ .

625 In this paper, we conceptualize the aggregated model parameters on the server as the summation of  
 626 optimal parameters and a certain drift (or noise), represented as:  $w_{server} = w^* + w_{drift}$ . We also  
 627 assume that  $w_{drift}$  is subject to a normal (Gaussian-like) distribution, and justify this assumption by  
 628 demonstrating its prevalence, and explaining it in FL.

629 From a historical standpoint, modeling noise in dynamic systems as a Gaussian-like distribution is a  
 630 widely accepted practice. This dates back to [Kramers, 1940], and many studies analyzing Stochastic  
 631 Gradient Descent (SGD) optimization have emphasized the Gaussian nature of noise on gradients  
 632 or parameters [Mandt et al., 2017, Zhu et al., 2019, Simsekli et al., 2019, Ziyin et al., 2021]. This  
 633 assumption of Gaussianity for SGD noise is justified by Wu et al. [2020], which guarantees the

634 SGD noise's convergence to a specific infinite divisible distribution. This falls under the Gaussian  
 635 class provided the noise's second moment is finite (as per Lindeberg's condition). While it has been  
 636 proposed that the noise in SGD might be better represented by *SaS* noise [Simsekli et al., 2019], this  
 637 idea has been challenged and redirected back to the earlier proposed Gaussian noise model [Xie et al.,  
 638 2021, Battash and Lindenbaum, 2023].

639 We further elucidate the occurrence of Gaussian noise in the context of FL. The Lindeberg-Feller  
 640 Central Limit Theorem (CLT) [Lindeberg, 1922] plays a key role in explaining the prevalence of  
 641 the normal distribution. It posits that the sum (or average) of random variables gravitates towards a  
 642 normal distribution (no need to assume iid of these random variables themselves), irrespective of the  
 643 individual distributions of these variables. In FL, the emergence of noise can be attributed to partial  
 644 participation (referred to as period drift) and multiple local updates (referred to as client drift). The  
 645 drifted model that we observe,  $w_{server}$ , is typically the result of the combination of these factors,  
 646 making the normal distribution an apt model for characterizing the noise.

**The impact of noise on generalization** In order to investigate the effect of the deviation on performance in FL, we utilize a regression optimization objective as in previous studies, such as [Zhang et al., 2019] and [Wu et al., 2018]:

$$\hat{\mathcal{L}}(w) = \frac{1}{2}(w + w_{drift})^T A(w + w_{drift}),$$

647 where  $w_{drift} \sim \mathcal{N}(0, \sigma^2)$  is the drift caused by the characteristics of FL. Therefore, the generalization  
 648 error can be formulated as:

$$\begin{aligned}\mathcal{L}(w^t) &= \mathbb{E}[\hat{\mathcal{L}}(w^t)] = \frac{1}{2}\mathbb{E}\left[\sum_i a_i (w_i^{t^2} + \sigma_i^2)\right] \\ &= \frac{1}{2} \sum_i a_i (\mathbb{E}[w_i^t]^2 + \mathbb{V}[w_i^t] + \sigma_i^2),\end{aligned}$$

649 where  $A$  is the matrix of quadratic form of the MSE loss function, and  $a_i$  is the elements of  $A$ . As  
 650 results, the generalization error can also be decomposed into three components: bias, variance, and  
 651 noise. The noise component in FL context is further influenced by factors such as client sampling,  
 652 multiple local epochs, and non-iid data distribution, leading to a much larger overall generalization  
 653 error compared to centralized SGD. This formulation reveals the reason of why FL usually performs  
 654 worse than centralized training. Thus, our goal is to reduce the variance of drift  $\sigma^2$  in order to improve  
 655 both the convergence and performance of the model.

## 656 D.2 The justification of independence assumption 2.3 of client drift and period drift

657 **Assumption D.2** (2.3). *The initialization model parameters are independent of all period drifts  
 658  $Q_t$  and client drifts  $R_t$  at each communication round, that is  $w_0 \perp Q_0, Q_1, \dots, Q_t$  and  $w_0 \perp R_0, R_1, \dots, R_t$ .*

660 This assumption can be justified from various perspectives, demonstrating its reasonableness:

- 661 • **Independence of the initial model parameters from other noise variables:** This assumption  
 662 suggests that there is no direct relationship between the initial model parameters ( $w_0$ ) and period  
 663 drifts or client drifts. This is a reasonable assumption since initial model parameters are typically  
 664 determined prior to training and hence are not influenced by any noise processes.
- 665 • **Independence of client drifts between each communication round:** According to this assump-  
 666 tion, client drifts ( $R_0, R_1, \dots, R_t$ ) in different communication rounds are independent of each  
 667 other. The client drift is influenced by data heterogeneity and multiple local updates. A higher  
 668 degree of data heterogeneity and an increased number of local updates can result in greater client  
 669 drift. Each client has its own fixed client drift[Guo et al., 2021], and the client drift in each  
 670 communication round doesn't impact other rounds, the assumption of client drift independence is  
 671 reasonable.
- 672 • **Independence of period drifts between each communication round:** This assumption contends  
 673 that the period drifts ( $Q_0, Q_1, \dots, Q_t$ ) across different communication rounds are independent.  
 674 Period drift is influenced by data heterogeneity and client sampling. Though period drift may be  
 675 caused by biased client sampling due to factors like time and geographic locations, leading to a  
 676 dependence between period drifts in different communication rounds, this paper considers the case  
 677 where random client sampling occurs in each communication round. Here, one round of client  
 678 sampling doesn't affect others, thus rendering the independence of period drift reasonable.

- 679 • **Independence between period drifts and client drifts at each communication round:** This  
 680 assumption argues that the client drifts ( $R_0, R_1, \dots, R_t$ ) and the period drifts ( $Q_0, Q_1, \dots, Q_t$ ) at  
 681 each communication round are independent. While both client drift and period drift are influenced  
 682 by data heterogeneity, they are conditionally independent given the heterogeneous data. We offer a  
 683 causal graph to depict their relationships:  
 684 multiply local updates  $\rightarrow$  client drift  $\leftarrow$  data heterogeneity  $\rightarrow$  period drift  $\leftarrow$  client sampling  
 685 This graph indicates that data heterogeneity is a common cause of client drift and period drift, and  
 686 varying levels of data heterogeneity result in different magnitudes of client drift and period drift.  
 687 However, given that the heterogeneous data is constant across clients (i.e., given D), we can express  
 688  $P(\text{client drift, period drift}|D) = P(\text{client drift}|D) * P(\text{period drift}|D)$ . Indeed, conditioning on a given  
 689 heterogeneous data set is a fundamental assumption in Federated Learning.

690 **Empirical Analysis** To verify that both period drift and client drift adhere to the Gaussian assump-  
 691 tion, we design this experiment by means of randomly selecting a model parameter and tracing its  
 692 period drift and client drift across the subsequent 1500 rounds of training. During this process, we  
 693 collected the global optima  $x_{\mathcal{N}}^*$  of global objective on the whole training dataset for each communica-  
 694 tion round, as well as the optima  $x_S^*$  on the dataset of sampled clients for each communication round,  
 695 and the local optima  $x_k^*$  of the dataset of a single client, respectively. The period drift is represented  
 696 by  $x_S^* - x_{\mathcal{N}}^*$ , while the client drift is depicted by  $x_k^* - x_S^*$ . The sampling results are visualized  
 697 using histograms. Experiments were conducted with alpha values in the scope  $\alpha = [1, 0.1, 0.01]$ ,  
 698 as illustrated in Fig 7. We also employ the ‘normaltest’ function from ‘scipy.stats’, utilizing the  
 699 \*D’Agostino-Pearson Test\* to determine if a sample deviates from a normal distribution. Here,  
 700 a p-value  $> 0.05$  indicates conformity to a normal distribution. The experimental results showed  
 701 that both period drift and client drift follow a normal distribution, confirming the reliability of the  
 702 Gaussian distribution assumption.

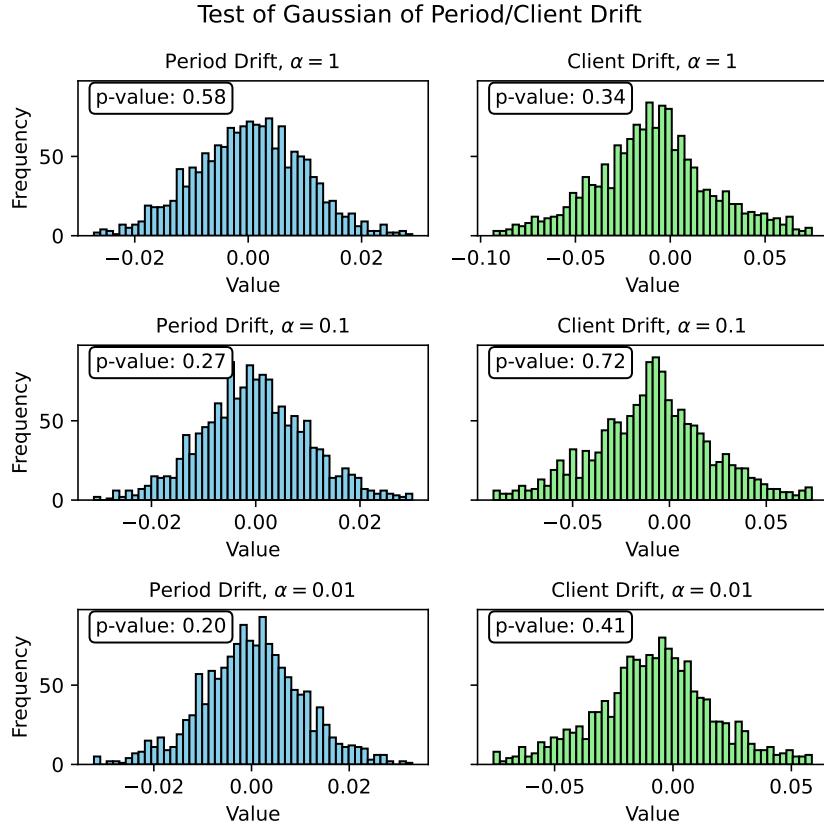


Figure 7: **Empirical justification of Gaussian distribution assumption.**

703 **D.3 The Independence of Noise**

704 **Lemma D.3** (2.4). (*Independence of Noise*). *the noise present in the prediction and observation at each communication round is independent of the current state of the model, specifically,  $w_t \perp Q_t$  and  $w_t \perp R_t$ .*

707 To prove this, we will first need to understand the relationship between the variables  $w_{t+1}$ ,  $Q_t$ , and  
 708  $R_t$ . In the given context,  $w_{t+1}$  represents the state of the model at a particular communication round  
 709 (say, round  $t + 1$ ). It is affected by the values of the period drift ( $Q_t$ ) and client drift ( $R_t$ ) at the  
 710 previous round. This is represented by the state transfer function  $w_{t+1} = T_t(w_t, Q_t, R_t)$ . However,  
 711 this relationship does not imply that  $w_{t+1}$  is dependent on  $Q_t$  or  $R_t$ . To see why, let's look at how  $w_t$   
 712 is formed. Using the state transfer function, we can express  $w_t$  as:

$$\begin{aligned} w_t &= T_{t-1}(w_{t-1}, Q_{t-1}, R_{t-1}), \\ w_{t-1} &= T_{t-2}(w_{t-2}, Q_{t-2}, R_{t-2}), \\ &\vdots \\ w_2 &= T_1(w_1, Q_1, R_1), \\ w_1 &= T_0(w_0, Q_0, R_0), \end{aligned} \tag{23}$$

713 From this chain of equations, it is evident that  $w_t$  is not only a function of the current round's  
 714 period drift  $Q_t$  and client drift  $R_t$ , but also of their past values and the past values of  $w_t$  itself. In  
 715 a more generalized form, we can write this as:  $w_t = T(w_0, Q_0, Q_1, \dots, Q_{t-1}, R_0, R_1, \dots, R_{t-1})$ .  
 716 Assumption 2.3 states that the period drift and client drift are independent of each other at each  
 717 communication round and also independent of the initial model parameters. That is,  $w_0 \perp Q_0 \perp$   
 718  $Q_1 \perp \dots \perp Q_t \perp R_0 \perp R_1 \perp \dots \perp R_t$ . This assumption implies that the previous states of  $w_t$  ( $w_{t-1}$ ,  
 719  $w_{t-2}$ , and so on) do not have any influence on the current values of  $Q_t$  and  $R_t$ . In other words, the  
 720 noise present at each round is independent of the model's current state. Thus,  $w_t$  is independent of  
 721  $Q_t$  and  $R_t$ , denoted as  $w_t \perp Q_t \perp R_t$ . Therefore, it can be concluded that the noise present in the  
 722 prediction and observation at each communication round is indeed independent of the current state of  
 723 the model, thereby confirming the independence of noise. This statement about the independence of  
 724 noise is significant because it confirms that the noise encountered during each communication round  
 725 does not affect the model's state. This means that the model is robust and not affected by random  
 726 perturbations, which is a desirable property in any machine learning model.

727 **D.4 Analysis of Kalman Gain in FEDEVE**

728 We conducted an in-depth analysis of the Kalman Gain  $K$  of FedEve under various experimental  
 729 settings, incorporating four levels of data heterogeneity and various local epochs, as shown in Figure  
 730 4. We observed that as data heterogeneity increases (i.e., as the value of  $\alpha$  decreases), the Kalman  
 731 Gain  $K$  progressively enlarges. With the rise of data heterogeneity, the period drift starts to play a  
 732 more dominant role. In this context, the primary role of Kalman Gain  $K$  is to adjust the weights  
 733 between global and local updates, as depicted by Equation (15b) and (15c). Further, according to  
 734 Equation (15d), the model update tends to trust local updates more, stabilizing the optimization  
 735 process. For varied counts of local updates, the relative change in Kalman Gain  $K$  is marginal. This  
 736 is primarily because, in cross-device FL, the client drift is not a pivotal or dominant factor, which  
 737 aligns with our prior analysis in Figure 3.

738 **D.5 Model Update with Bayesian filter**

739 **D.5.1 Bayesian filter**

740 This section begins by describing the main idea behind the approach: the combination of prediction  
 741 and observation models using a Bayesian filter, which is a statistical tool for estimating an unknown  
 742 probability density function (PDF) based on observations. The prediction process is explained using  
 743 the concept of cumulative distribution functions (CDFs), which are functions giving the probability  
 744 that a random variable is less than or equal to a certain value. The prediction process is characterized

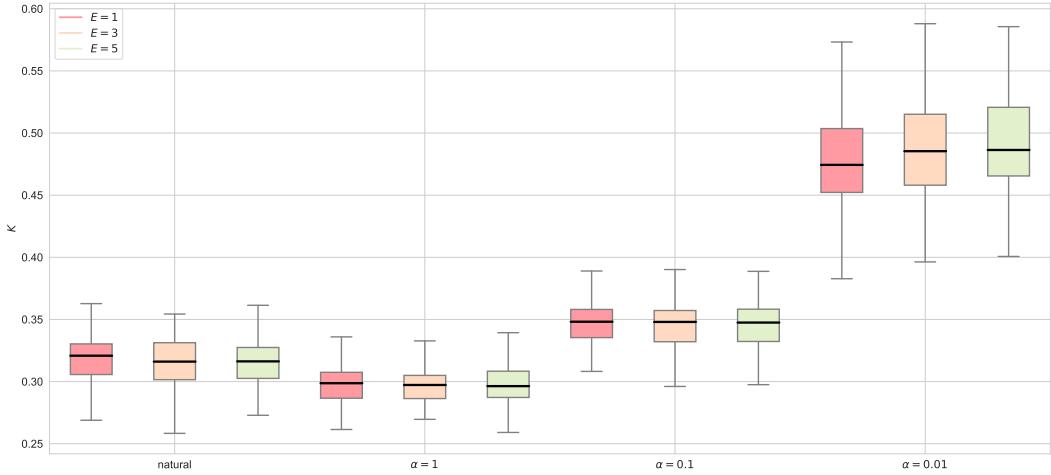


Figure 8: **Boxplots for Kalman Factors**

745 by the distribution:

$$\begin{aligned}
 F_{\hat{w}_{t+1}}^-(w) &= P(\hat{w}_{t+1} \leq w) \\
 &= \sum_{u=-\infty}^w P(\hat{w}_{t+1} = u) \\
 &= \sum_{u=-\infty}^w \sum_{v=-\infty}^{+\infty} P(\hat{w}_{t+1} = u \mid w_t = v) P(w_t = v) \\
 &= \sum_{u=-\infty}^w \sum_{v=-\infty}^{+\infty} P[\hat{w}_{t+1} - f(w_t) = u - f(v) \mid w_t = v] P(w_t = v) \\
 &= \sum_{u=-\infty}^w \sum_{v=-\infty}^{+\infty} P[Q_t = u - f(v) \mid w_t = v] P(w_t = v) \quad \Rightarrow \text{Prediction Equation} \\
 &= \sum_{u=-\infty}^w \sum_{v=-\infty}^{+\infty} P[Q_t = u - f(v)] P(w_t = v) \quad \Rightarrow \text{Lemma (2.4)} \\
 &= \sum_{u=-\infty}^w \left\{ \lim_{\epsilon \rightarrow 0} \sum_{v=-\infty}^{+\infty} f_{Q_t}[u - f(v)] \cdot \epsilon \cdot f_{w_t}^+(v) \cdot \epsilon \right\} \\
 &= \sum_{u=-\infty}^w \left\{ \lim_{\epsilon \rightarrow 0} \int_{-\infty}^{+\infty} f_{Q_t}[u - f(v)] f_{w_t}^+(v) dv \cdot \epsilon \right\} \\
 &= \int_{-\infty}^w \int_{-\infty}^{+\infty} f_{Q_t}[u - f(v)] f_{w_t}^+(v) dv du \\
 &= \int_{-\infty}^w \int_{-\infty}^{+\infty} f_{Q_t}[w - f(v)] f_{w_t}^+(v) dv dw
 \end{aligned} \tag{24}$$

746 The first three steps apply the definition of the cumulative distribution function (CDF), which is the  
 747 sum of probabilities up to a certain point. In the fourth step, the change of variables is used to switch  
 748 from  $u$  to  $v$ . The fifth and sixth steps apply Lemma 2.4, which states that the drift is independent of  
 749 the weights  $w_t$ . The last three steps show how to convert the sum to an integral, which is a common  
 750 method in probability theory for dealing with continuous random variables. Finally, the PDF of the

751 prediction is obtained by taking the derivative of the CDF, which can be expressed as:

$$f_{\hat{w}_{t+1}}^-(w) = \frac{dF_{\hat{w}_{t+1}}^-(w)}{dw} = \int_{-\infty}^{+\infty} f_{Q_t}[w - f(v)] f_{w_t}^+(v) dv \quad (25)$$

752 The observation process is also characterized by a PDF. Similar steps are used to manipulate and  
 753 simplify the expression for the PDF of the observed value  $w_{t+1}$ , given the predicted value  $\hat{w}_{t+1}$ .  
 754 Specifically:

$$\begin{aligned} f_{\tilde{w}_{t+1}|\hat{w}_{t+1}}(w_{t+1} | w) &= \lim_{\epsilon \rightarrow 0} \frac{F_{\tilde{w}_{t+1}|\hat{w}_{t+1}}(w_{t+1} + \epsilon | w) - F_{\tilde{w}_{t+1}|\hat{w}_{t+1}}(w_{t+1} | w)}{\epsilon} \\ &= \lim_{\epsilon \rightarrow 0} \frac{P(w_{t+1} \leq \tilde{w}_{t+1} \leq w_{t+1} + \epsilon | \hat{w}_{t+1} = w)}{\epsilon} \\ &= \lim_{\epsilon \rightarrow 0} \frac{P[w_{t+1} - h(w) \leq \tilde{w}_{t+1} - h(\hat{w}_{t+1}) \leq w_{t+1} - h(w) + \epsilon | \hat{w}_{t+1} = w]}{\epsilon} \\ &= \lim_{\epsilon \rightarrow 0} \frac{P[w_{t+1} - h(w) \leq R_t \leq w_{t+1} - h(w) + \epsilon | \hat{w}_{t+1} = w]}{\epsilon} \Rightarrow \text{Lemma (2.4)} \\ &= \lim_{\epsilon \rightarrow 0} \frac{F_{R_t}[w_{t+1} - h(w) + \epsilon] - F_{R_t}[w_{t+1} - h(w)]}{\epsilon} \\ &= f_{R_t}[w_{t+1} - h(w)]. \end{aligned} \quad (26)$$

755 The first two steps apply the definition of the probability density function (PDF), which is the  
 756 derivative of the cumulative distribution function (CDF). The third and fourth steps use a change of  
 757 variables to express the PDF in terms of the difference between the observed and predicted values.  
 758 The fifth step applies the independence Lemma (2.4) to simplify the conditional probability. The  
 759 sixth step calculates the limit to arrive at the PDF of the observation process. As a consequence of  
 760 combining the prediction and observation distributions, a fused model can be obtained, which can be  
 761 described by its own probability density function (PDF) as:

$$f_{w_{t+1}}^+(w) = \eta_t \cdot f_{\tilde{w}_{t+1}|\hat{w}_{t+1}}(\tilde{w}_{t+1} | \hat{w}_{t+1}) \cdot f_{\hat{w}_{t+1}}^-(w) = \eta_t \cdot f_{R_t}[\tilde{w}_{t+1} - h(\hat{w}_{t+1})] \cdot f_{\hat{w}_{t+1}}^-(w), \quad (27)$$

762 where

$$\eta_t = \left[ \int_{-\infty}^{+\infty} f_{\tilde{w}_{t+1}|\hat{w}_{t+1}}(\tilde{w}_{t+1} | \hat{w}_{t+1}) f_{\hat{w}_{t+1}}^-(w) dw \right]^{-1} = \left\{ \int_{-\infty}^{+\infty} f_{R_t}[\tilde{w}_{t+1} - h(\hat{w}_{t+1})] f_{\hat{w}_{t+1}}^-(w) dw \right\}^{-1}. \quad (28)$$

763 The PDF of the fused model is obtained by multiplying the PDFs of the prediction and observation  
 764 processes, normalized by a factor  $\eta_t$ . The process of updating the fused model, also known as the  
 765 Bayesian filter, can be summarized in the following steps:

$$\begin{aligned} f_{w_t}^+(w) &\xrightarrow{\text{predict}} f_{\hat{w}_{t+1}}^-(w) = \int_{-\infty}^{+\infty} f_{Q_t}[w - f(v)] f_{w_t}^+(v) dv \\ &\xrightarrow{\text{observe}} f_{w_{t+1}}^+(w) = \eta_t \cdot f_{R_t}[w_{t+1} - h(w)] \cdot f_{\hat{w}_{t+1}}^-(w), \end{aligned} \quad (29)$$

766 where  $\eta_t = \left\{ \int_{-\infty}^{+\infty} f_{R_t}[\tilde{w}_{t+1} - h(\hat{w}_{t+1})] f_{\hat{w}_{t+1}}^-(w) dw \right\}^{-1}$ . The fused model combines the predi-  
 767 cition and observation distributions, and it describes the sequential steps of the Bayesian filter: starting  
 768 with the PDF at time  $t$ , a prediction is made for the PDF at time  $t + 1$ , and then this prediction  
 769 is updated based on the observation to obtain the PDF at time  $t + 1$ . The final estimation of the  
 770 parameter can be obtained as a result of these update steps and can be represented as:

$$\hat{w}_{t+1} = E[f_{w_{t+1}}^+(w)] = \int_{-\infty}^{+\infty} w f_{w_{t+1}}^+(w) dw, \quad (30)$$

771 which is done by calculating the expected value of the PDF at time  $t + 1$ . This is performed by  
 772 multiplying the parameter  $w$  by its probability density and integrating over all possible values of  $w$ .  
 773 The integral provides a single, average value for  $w$  weighted by its probability density, which serves  
 774 as the final estimate of the parameter.

775 **D.5.2 FEDEVE**

776 In this section, we delve into the derivation of the FEDEVE algorithm using the Bayesian filter. The  
 777 model update process within the FEDEVE algorithm is outlined as follows: Firstly, we calculate the  
 778 predictive value of the model's parameters ( $\hat{w}_{t+1}$ ) at the next time step using the current parameters  
 779 ( $w_t$ ) and the step-size scaled momentum ( $\eta_g M_t$ ):

$$\hat{w}_{t+1} = w_t - \eta_g M_t, \quad (31)$$

780 The inverse of the variance at time  $t + 1$  ( $\hat{\sigma}_{t+1}^2$ ) is determined as the sum of the predicted variance at  
 781 time  $t$  ( $\sigma_t^2$ ) and the squared variance associated with the process noise ( $\sigma_{Q_t}^2$ ):

$$\hat{\sigma}_{t+1}^2 = \sigma_t^2 + \sigma_{Q_t}^2, \quad (32)$$

782 The Kalman Gain ( $K$ ) is computed as the ratio of the inverse of the variance at time  $t + 1$  to the sum  
 783 of the inverse of the variance at time  $t + 1$  and the variance of the observation noise ( $\sigma_{R_t}^2$ ):

$$K = \frac{\hat{\sigma}_{t+1}^2}{\hat{\sigma}_{t+1}^2 + \sigma_{R_t}^2}, \quad (33)$$

784 The momentum for the next time step ( $M_{t+1}$ ) is obtained by adjusting the current momentum ( $M_t$ )  
 785 based on the difference between the observed value ( $\Delta \tilde{w}_t$ ) and the current momentum:

$$M_{t+1} = M_t + K(\Delta \tilde{w}_t - M_t), \quad (34)$$

786 The parameters ( $w_{t+1}$ ) for the next time step are calculated by subtracting the step-size scaled updated  
 787 momentum from the current parameters:

$$w_{t+1} = w_t - \eta_g M_{t+1}, \quad (35)$$

788 Finally, the predicted variance for the next time step ( $\sigma_t^2$ ) is computed by scaling the inverse of the  
 789 variance at time  $t + 1$  by  $(1 - K)$ :

$$\sigma_{t+1}^2 = (1 - K)\hat{\sigma}_{t+1}^2. \quad (36)$$

790 A key assumption for this derivation is that the two random variables,  $A$  and  $B$ , follow a normal  
 791 distribution. Specifically,  $A$  is assumed to be distributed as  $\mathcal{N}(\mu_A, \sigma_A^2)$ , and  $B$  is assumed to be  
 792 distributed as  $\mathcal{N}(\mu_B, \sigma_B^2)$ . Given these assumptions, it can be mathematically proven that the sum  
 793 and the product of  $A$  and  $B$  also follow a normal distribution. In particular, we have:

$$A + B \sim \mathcal{N}(\mu_A + \mu_B, \sigma_A^2 + \sigma_B^2), \quad (37)$$

$$A \times B \sim \mathcal{N}\left(\frac{\mu_A \sigma_B^2 + \mu_B \sigma_A^2}{\sigma_A^2 + \sigma_B^2}, \frac{\sigma_A^2 \sigma_B^2}{\sigma_A^2 + \sigma_B^2}\right), \quad (38)$$

794 In the context of the predict-observe framework and the Bayesian filter, we make a few key assumptions:  
 795 firstly,  $w_t$  is distributed as  $\mathcal{N}(w_t, \sigma_t^2)$ ; secondly,  $Q_t$  is distributed as  $\mathcal{N}(0, \sigma_{Q_t}^2)$ ; and finally,  
 796 the momentum  $M_t$  is considered a non-random variable. Specializing the prediction function as a  
 797 linear function, as shown in Equation (13), leads us to the following result:

$$\hat{w}_{t+1} \sim \mathcal{N}(w_t - \eta_g M_t, \sigma_t^2 + \sigma_{Q_t}^2), \quad (39)$$

798 This is essentially equivalent to the equations (15a) and (15b) in the model update process. Moreover,  
 799 we set  $\sigma_{t+1}^{-2} = \sigma_t^2 + \sigma_{Q_t}^2$ . The distribution of  $w_{t+1}$  is considered as the posterior, which is calculated  
 800 by applying the Bayesian formula and combining the product of the likelihood and the prior. Here, the  
 801 likelihood corresponds to the observation, and the prior corresponds to the prediction. The observed  
 802  $\tilde{w}_{t+1}$  is distributed as follows:

$$\tilde{w}_{t+1} \sim \mathcal{N}(\hat{w}_{t+1} - \eta_g \Delta \tilde{w}_t, \sigma_{R_t}^2). \quad (40)$$

803 To calculate the product of the prior and the likelihood, we evaluate the proportionality coefficient  
 804  $K = \frac{\hat{\sigma}_{t+1}^2}{\hat{\sigma}_{t+1}^2 + \sigma_{R_t}^2}$ . We can then assert that  $w_{t+1}$  also adheres to a normal distribution:

$$\tilde{w}_{t+1} \sim \mathcal{N}(w_t - \eta_g((1 - K)M_t + K\Delta \tilde{w}_t), (1 - K)\hat{\sigma}_{t+1}^2), \quad (41)$$

805 This result serves as a stepping stone for the subsequent round of computation. Consequently, the  
 806 variance of  $w_{t+1}$  is minimized as:

$$\sigma_{\text{fused}}^2 = \frac{\hat{\sigma}_{t+1}^2 \sigma_{R_t}^2}{\hat{\sigma}_{t+1}^2 + \sigma_{R_t}^2}. \quad (42)$$

807 In summary, the above proof demonstrates how the Bayesian filter can be used to derive the model  
 808 update process of the FEDEVE algorithm. The predict-observe framework and the feature of normal  
 809 distribution are key elements in this derivation.

810 **D.6 Pseudo-code**

811 The pseudo-code of FED-EVE is depicted in Algorithm 1.

---

**Algorithm 1 FED-EVE** The selected clients are indexed by  $k$ ;  $E$  is the number of local epochs, and  $\eta_l$  is the local learning rate.

---

**Server executes:**

```

initialize  $w_0$ 
for each round  $t = 1, 2, \dots, T$  do
     $\hat{w}_{t+1} \leftarrow w_t - \eta_g M_t$  as in Equation (12) in the main paper // predict
     $\mathcal{S}_t \leftarrow$  randomly select  $|\mathcal{S}_t|$  clients
    for each client  $k \in \mathcal{S}_t$  in parallel do
         $w_t^k \leftarrow \text{ClientUpdate}(k, \hat{w}_{t+1})$ 
    end for
     $\Delta\tilde{w}_t = \sum_{k \in \mathcal{S}_t} p_k (\hat{w}_{t+1} - w_t^k)$  // observe
    Model update: executes Equations (13)-(17) in the main paper
end for

```

```

ClientUpdate( $k, w$ ): // run on client  $k$ 
 $\mathcal{B} \leftarrow$  (split local data into batches of size)
for each local epoch  $i$  from 1 to  $E$  do
    for batch  $b \in \mathcal{B}$  do
         $w \leftarrow w - \eta_l F_k(w; b)$ 
    end for
end for
return  $w$  to server

```

---

812 **E Experimental Details**

813 **Implementation.** All the experiments are implemented using PyTorch, a popular machine learning  
 814 library. We simulate the federated learning environment, including clients, and run all experiments  
 815 on a deep learning server equipped with an NVIDIA GTX 2080 ti GPU.

816 **E.1 Evaluation metric.**

817 For the RS task, the model performance is evaluated using the following metrics: area under curve  
 818 (AUC), Hit Ratio (HR) and Normalized Discounted Cumulative Gain (NDCG). For the CV task, the  
 819 model performance is measured by the widely used Top-1 accuracy metric. In the experiments, for  
 820 the CTR (Click-Through Rate) task, the model performance is evaluated using the following metrics:  
 821 area under curve (AUC), Hit Ratio (HR) and Normalized Discounted Cumulative Gain (NDCG).

$$\text{AUC} = \frac{\sum_{x_0 \in D_T} \sum_{x_1 \in D_F} \mathbf{1}[f(x_1) < f(x_0)]}{|D_T| |D_F|},$$

$$\text{HitRate@K} = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \mathbf{1}(R_{u,g_u} \leq K),$$

$$\text{NDCG@K} = \sum_{u \in \mathcal{U}} \frac{1}{|\mathcal{U}|} \frac{2^{\mathbf{1}(R_{u,g_u} \leq K)} - 1}{\log_2(\mathbf{1}(R_{u,g_u} \leq K) + 1)},$$

822 where  $\mathcal{U}$  is the set of users,  $\mathbf{1}$  is the indicator function,  $R_{u,g_u}$  is the rank generated by the model  
 823 for the ground truth item  $g_u$ ,  $f$  is the model being evaluated, and  $D_T$  and  $D_F$  are the positive and  
 824 negative sample sets in the testing data, respectively. For the image classification task, the model  
 825 performance is measured by the widely used Top-1 accuracy metric.

826 **E.2 Detailed Settings of Figure 3 in the Main Paper.**

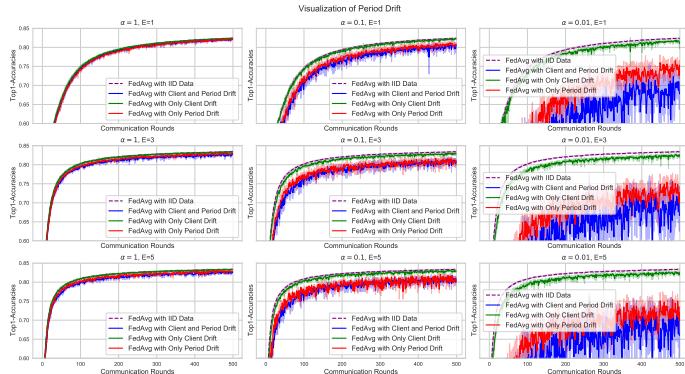


Figure 9: Visualization of period drift with different  $\alpha$  and  $E$

- 827 • **Figure 3 (a).** To provide a clearer illustration, we displayed the label distribution of the  
 828 10-digit classes, rather than the complete 62 classes in the original FEMNIST dataset, for 20  
 829 communication rounds.
- 830 • **Figure 3 (b).** 1) FEDAVG with iid data with no period or client drift: FEDAVG with iid data,  
 831 where training data is randomly partitioned among all clients, resulting in no period or client  
 832 drift; 2) FEDAVG with only period drift: non-iid data is initially partitioned, but the training  
 833 data of the sampled clients is randomly reshuffled and iid-distributed evenly among clients  
 834 in each round, resulting in period drift but no client drift; 3) FEDAVG with only client drift:  
 835 iid data is initially partitioned, but the training data of the sampled clients in each round is  
 836 re-partitioned in non-iid setting, resulting in client drift but no period drift; 4) FEDAVG with  
 837 both period and client drift: FEDAVG with non-iid data, where data is partitioned in non-iid  
 838 setting, resulting in both period and client drift.

839 To make the results more convincing, we conducted more experiments on FEMNIST. Specifically,  
 840 we add experiments with various local epochs and data heterogeneity. Each experiment is repeated 5  
 841 times, and the results are shown as follow:

842 **F Convergence Analysis**

843 This section provides a comprehensive theoretical analysis of the FedEve algorithm. We establish  
 844 formal guarantees on the convergence of our proposed method by first stating necessary assumptions,  
 845 then presenting the main convergence theorem, and finally providing supporting lemmas and detailed  
 846 proofs.

847 **F.1 Assumptions**

848 For our theoretical analysis, we rely on the following standard assumptions in federated optimization  
 849 literature:

850 **Assumption F.1** (Lipschitz Continuity). The loss function  $l(\cdot, z)$  is  $L$ -Lipschitz continuous, that is,  
 851  $|l(w; z) - l(w'; z)| \leq L\|w - w'\|$ , and is  $L$ -smooth for any  $z$ , that is,  $\|\nabla l(w; z) - \nabla l(w'; z)\| \leq$   
 852  $L\|w - w'\|$  for any  $z, w, w'$ .

853 The Lipschitz continuity assumption ensures that the loss function doesn't change too rapidly as the  
 854 model parameters change, which is crucial for establishing the stability of our optimization procedure.  
 855 The smoothness property ensures that the gradients of the loss function are well-behaved, which is  
 856 necessary for proving convergence rates.

857 **Assumption F.2** (Bounded Variance). The function  $f_i$  has  $\sigma$ -bounded variance, i.e.,  $\mathbb{E}\|f_i(w) -$   
 858  $\nabla f_i(w)\| \leq \sigma$  for all  $w \in \mathbb{R}^d$  and  $i \in [N]$ .

859 This assumption limits the variance of the stochastic gradients at each client, which is essential  
 860 in federated learning where we only have access to gradient estimates from a subset of clients.  
 861 This bounded variance allows us to control the error introduced by client sampling and stochastic  
 862 optimization.

863 **F.2 Theorems**

864 Our main theoretical result establishes the convergence rate of the FedEve algorithm:

865 **Theorem F.3** (Convergence of FedEve). *Suppose Assumptions F.1,F.2 hold and  $\eta_g \leq \frac{1}{L}$ . Then, for  
 866 the FedEve algorithm, we have:*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(w_t)\|^2] \leq \mathcal{O} \left( \frac{L(f(w_0) - f^*)}{T} + \frac{G_{kal}^2 \sigma^2}{S} \left( 1 - \frac{S}{N} \right) \right).$$

867 This theorem provides a bound on the average squared norm of gradients, which is a standard measure  
 868 of convergence in non-convex optimization. The bound consists of two terms:

- 869 • The first term  $\mathcal{O} \left( \frac{L(f(w_0) - f^*)}{T} \right)$  decreases with the number of iterations  $T$ , indicating that  
 870 the algorithm converges as  $T$  increases.
- 871 • The second term  $\mathcal{O} \left( \frac{G_{kal}^2 \sigma^2}{S} \left( 1 - \frac{S}{N} \right) \right)$  represents the irreducible error due to client sampling  
 872 and stochastic gradients, which decreases as the number of sampled clients  $S$  increases and  
 873 vanishes when all clients participate ( $S = N$ ).

874 This convergence bound demonstrates that FedEve achieves the optimal convergence rate for non-  
 875 convex optimization while effectively handling the challenges of federated learning settings.

876 **F.3 Lemmas**

877 To prove our main theorem, we establish the following key lemmas that characterize the behavior of  
 878 client sampling and the Kalman filter-based momentum updates:

879 **Lemma F.4** (Variance of Local Gradients). *Given Assumptions F.1,F.2, the variance of the local  
 880 gradients can be bounded as:*

$$\mathbb{E}[\|\nabla f_S(w) - \nabla f(w)\|^2] \leq \frac{\sigma^2}{S} \left( 1 - \frac{S}{N} \right).$$

881 This lemma quantifies the error introduced by client sampling in federated learning. It shows that the  
 882 variance of the gradients from a subset of clients  $S$  compared to the full gradient decreases as we  
 883 sample more clients (larger  $S$ ), and vanishes when all clients participate ( $S = N$ ).

884 **Lemma F.5** (Kalman Filter Update Variance). *Given Assumptions F.1,F.2, the variance of the  
885 momentum update using the Kalman filter can be bounded as:*

$$\mathbb{E}[\|M_{t+1} - M_t\|^2] \leq G_{kal}^2 \left( \frac{\sigma^2}{S} \left( 1 - \frac{S}{N} \right) + \mathbb{E}[\|\nabla f(w_t) - M_t\|^2] \right).$$

886 This lemma characterizes the stability of the momentum updates in our FedEve algorithm. It shows  
887 that the variance of the momentum changes is controlled by both the client sampling variance and  
888 the current error in the momentum estimate. The Kalman gain parameter  $G_{kal}$  allows us to balance  
889 between adaptivity and stability in the momentum updates.

## 890 F.4 Proofs

891 We now provide detailed proofs for our lemmas and main theorem. Each proof is structured to  
892 clearly show the logical progression from assumptions to conclusions, with detailed explanations of  
893 intermediate steps.

### 894 F.4.1 Proof of Lemma F.4

895 *Proof.* Our goal is to bound the variance between the gradient estimated from a subset of clients and  
896 the full gradient across all clients. We start by analyzing the definition of the sampled gradient:

$$\nabla f_{\mathcal{S}}(w) = \frac{1}{S} \sum_{i \in \mathcal{S}} \nabla f_i(w).$$

897 Given that  $\mathcal{S}$  is a randomly selected subset of  $S$  clients from the total  $N$  clients without replacement,  
898 the variance we want to bound can be expressed as:

$$\mathbb{E}[\|\nabla f_{\mathcal{S}}(w) - \nabla f(w)\|^2] = \mathbb{E} \left[ \left\| \frac{1}{S} \sum_{i \in \mathcal{S}} \nabla f_i(w) - \frac{1}{N} \sum_{i=1}^N \nabla f_i(w) \right\|^2 \right].$$

899 To analyze this expression, we will decompose the difference by separating the selected clients  $\mathcal{S}$  and  
900 the non-selected clients  $\mathcal{S}^c$ . Note that  $\mathcal{S} \cup \mathcal{S}^c = \{1, 2, \dots, N\}$  and  $\mathcal{S} \cap \mathcal{S}^c = \emptyset$ . We can rewrite:

$$\begin{aligned} \frac{1}{S} \sum_{i \in \mathcal{S}} \nabla f_i(w) - \frac{1}{N} \sum_{i=1}^N \nabla f_i(w) &= \frac{1}{S} \sum_{i \in \mathcal{S}} \nabla f_i(w) - \frac{1}{N} \left( \sum_{i \in \mathcal{S}} \nabla f_i(w) + \sum_{i \in \mathcal{S}^c} \nabla f_i(w) \right) \\ &= \left( \frac{1}{S} - \frac{1}{N} \right) \sum_{i \in \mathcal{S}} \nabla f_i(w) - \frac{1}{N} \sum_{i \in \mathcal{S}^c} \nabla f_i(w). \end{aligned}$$

901 To compute the expected squared norm of this expression, we note that the two terms are independent  
902 due to the random sampling process. Therefore:

$$\begin{aligned} \mathbb{E} \left[ \left\| \left( \frac{1}{S} - \frac{1}{N} \right) \sum_{i \in \mathcal{S}} \nabla f_i(w) - \frac{1}{N} \sum_{i \in \mathcal{S}^c} \nabla f_i(w) \right\|^2 \right] \\ = \left( \frac{1}{S} - \frac{1}{N} \right)^2 \mathbb{E} \left[ \left\| \sum_{i \in \mathcal{S}} \nabla f_i(w) \right\|^2 \right] + \frac{1}{N^2} \mathbb{E} \left[ \left\| \sum_{i \in \mathcal{S}^c} \nabla f_i(w) \right\|^2 \right]. \end{aligned}$$

903 Now, we need to evaluate the expectations. Let's denote  $\mu = \frac{1}{N} \sum_{i=1}^N \nabla f_i(w)$  as the mean gradient  
904 across all clients. Under Assumption F.2, the client gradients have bounded variance  $\sigma^2$  around this  
905 mean. For a set of  $S$  randomly sampled clients, the variance of their sum is:

$$\mathbb{E} \left[ \left\| \sum_{i \in \mathcal{S}} \nabla f_i(w) - S\mu \right\|^2 \right] = S\sigma^2.$$

906 Since  $\mathbb{E}[\nabla f_i(w)] = \mu$  for all clients, we have:

$$\mathbb{E} \left[ \left\| \sum_{i \in \mathcal{S}} \nabla f_i(w) \right\|^2 \right] = S\sigma^2 + S^2\|\mu\|^2.$$

907 Similarly, for the non-selected clients:

$$\mathbb{E} \left[ \left\| \sum_{i \in \mathcal{S}^c} \nabla f_i(w) \right\|^2 \right] = (N - S)\sigma^2 + (N - S)^2 \|\mu\|^2.$$

908 Substituting these expressions into our variance formula:

$$\begin{aligned} & \left( \frac{1}{S} - \frac{1}{N} \right)^2 (S\sigma^2 + S^2 \|\mu\|^2) + \frac{1}{N^2} ((N - S)\sigma^2 + (N - S)^2 \|\mu\|^2) \\ &= \left( \frac{1}{S} - \frac{1}{N} \right)^2 S\sigma^2 + \left( \frac{1}{S} - \frac{1}{N} \right)^2 S^2 \|\mu\|^2 + \frac{(N - S)\sigma^2}{N^2} + \frac{(N - S)^2 \|\mu\|^2}{N^2}. \end{aligned}$$

909 After algebraic manipulation, the terms involving  $\|\mu\|^2$  cancel out (which can be verified by expanding  
910 the expressions), leaving us with:

$$\begin{aligned} \mathbb{E}[\|\nabla f_S(w) - \nabla f(w)\|^2] &= \left( \frac{1}{S} - \frac{1}{N} \right)^2 S\sigma^2 + \frac{(N - S)\sigma^2}{N^2} \\ &= \frac{\sigma^2}{S} \left( 1 - \frac{S}{N} \right). \end{aligned}$$

911 This final result elegantly quantifies how the variance scales with both the number of sampled clients  
912  $S$  and the total number of clients  $N$ . We observe that:

- 913 • As  $S$  increases, the variance decreases, showing the benefit of sampling more clients.
- 914 • When  $S = N$  (i.e., we use all clients), the variance becomes zero, as expected.
- 915 • The term  $1 - \frac{S}{N}$  represents the finite population correction factor from sampling theory,  
916 accounting for sampling without replacement.

917  $\square$

#### 918 F.4.2 Proof of Lemma F.5

919 *Proof.* This lemma characterizes the stability of momentum updates using the Kalman filter. We  
920 begin with the momentum update equation in the FedEve algorithm:

$$M_{t+1} = M_t + G_{kal}(\Delta \tilde{w}_t - M_t),$$

921 where  $\Delta \tilde{w}_t$  represents the average update from the selected clients at time  $t$ , and  $G_{kal}$  is the Kalman  
922 gain parameter that controls how much new information is incorporated into the momentum.

923 Rearranging the update equation, we have:

$$M_{t+1} - M_t = G_{kal}(\Delta \tilde{w}_t - M_t).$$

924 Our goal is to bound the expected squared norm of this difference, which measures how much the  
925 momentum changes between iterations:

$$\mathbb{E}[\|M_{t+1} - M_t\|^2] = \mathbb{E}[\|G_{kal}(\Delta \tilde{w}_t - M_t)\|^2].$$

926 Since  $G_{kal}$  is a scalar constant, we can factor it out:

$$\mathbb{E}[\|M_{t+1} - M_t\|^2] = G_{kal}^2 \mathbb{E}[\|\Delta \tilde{w}_t - M_t\|^2].$$

927 To analyze  $\mathbb{E}[\|\Delta \tilde{w}_t - M_t\|^2]$ , we introduce the true full gradient  $\nabla f(w_t)$  as an intermediate term:

$$\begin{aligned} \mathbb{E}[\|\Delta \tilde{w}_t - M_t\|^2] &= \mathbb{E}[\|\Delta \tilde{w}_t - \nabla f(w_t) + \nabla f(w_t) - M_t\|^2] \\ &= \mathbb{E}[\|(\Delta \tilde{w}_t - \nabla f(w_t)) + (\nabla f(w_t) - M_t)\|^2]. \end{aligned}$$

928 Now, we expand the squared norm of the sum. If we can show that the cross-term  $\mathbb{E}[\langle \Delta \tilde{w}_t -$   
929  $\nabla f(w_t), \nabla f(w_t) - M_t \rangle] = 0$ , then we can separate the expression. This is indeed the case because:

- 930 •  $M_t$  depends only on information up to time  $t - 1$

- 931     •  $\nabla f(w_t)$  is a deterministic function of  $w_t$   
 932     •  $\Delta \tilde{w}_t - \nabla f(w_t)$  is the sampling error at time  $t$ , which is independent of past information

933 Therefore, we can write:

$$\begin{aligned}\mathbb{E}[\|\Delta \tilde{w}_t - M_t\|^2] &= \mathbb{E}[\|\Delta \tilde{w}_t - \nabla f(w_t)\|^2] + \mathbb{E}[\|\nabla f(w_t) - M_t\|^2] \\ &\quad + 2\mathbb{E}[\langle \Delta \tilde{w}_t - \nabla f(w_t), \nabla f(w_t) - M_t \rangle] \\ &= \mathbb{E}[\|\Delta \tilde{w}_t - \nabla f(w_t)\|^2] + \mathbb{E}[\|\nabla f(w_t) - M_t\|^2].\end{aligned}$$

934 The first term,  $\mathbb{E}[\|\Delta \tilde{w}_t - \nabla f(w_t)\|^2]$ , represents the variance of the client updates. From Lemma  
 935 F.4, we know that:

$$\mathbb{E}[\|\Delta \tilde{w}_t - \nabla f(w_t)\|^2] = \frac{\sigma^2}{S} \left(1 - \frac{S}{N}\right).$$

936 The second term,  $\mathbb{E}[\|\nabla f(w_t) - M_t\|^2]$ , represents how far the current momentum estimate is from  
 937 the true gradient. This is the estimation error from previous iterations.

938 Substituting these results back into our original expression:

$$\mathbb{E}[\|M_{t+1} - M_t\|^2] = G_{kal}^2 \left( \frac{\sigma^2}{S} \left(1 - \frac{S}{N}\right) + \mathbb{E}[\|\nabla f(w_t) - M_t\|^2] \right).$$

939 This result provides valuable insights into the momentum dynamics:

- 940     • The variance of momentum changes is controlled by two factors: client sampling variance  
 941       and current momentum estimation error.  
 942     • The Kalman gain  $G_{kal}$  directly impacts the magnitude of momentum changes. A smaller  
 943        $G_{kal}$  leads to more stable but slower-adapting momentum, while a larger  $G_{kal}$  allows faster  
 944       adaptation but with potentially higher variance.  
 945     • When  $S = N$  (using all clients), the client sampling variance term disappears, but the  
 946       estimation error term remains, showing that momentum still provides a beneficial smoothing  
 947       effect even with full client participation.

948  $\square$

#### 949 F.4.3 Proof of Theorem F.3

950 *Proof.* Our convergence analysis follows the standard approach for analyzing first-order optimization  
 951 methods for non-convex functions. We start by using the  $L$ -smoothness property from Assumption  
 952 F.1 to bound the progress in the objective function between consecutive iterations.

953 For an  $L$ -smooth function  $f$ , we have the following inequality for any two points  $w$  and  $w'$ :

$$f(w') \leq f(w) + \langle \nabla f(w), w' - w \rangle + \frac{L}{2} \|w' - w\|^2.$$

954 Applying this to consecutive iterates  $w_t$  and  $w_{t+1}$  in our algorithm:

$$f(w_{t+1}) \leq f(w_t) + \langle \nabla f(w_t), w_{t+1} - w_t \rangle + \frac{L}{2} \|w_{t+1} - w_t\|^2.$$

955 In the FedEve algorithm, the update rule is  $w_{t+1} = w_t - \eta_g M_{t+1}$ , where  $M_{t+1}$  is the momentum  
 956 term updated using the Kalman filter and  $\eta_g$  is the global learning rate. Substituting this update rule:

$$\begin{aligned}f(w_{t+1}) &\leq f(w_t) + \langle \nabla f(w_t), -\eta_g M_{t+1} \rangle + \frac{L}{2} \|-\eta_g M_{t+1}\|^2 \\ &= f(w_t) - \eta_g \langle \nabla f(w_t), M_{t+1} \rangle + \frac{\eta_g^2 L}{2} \|M_{t+1}\|^2.\end{aligned}$$

957 Taking the expectation, we have:

$$\mathbb{E}[f(w_{t+1})] \leq f(w_t) - \eta_g \mathbb{E}[\langle \nabla f(w_t), M_{t+1} \rangle] + \frac{\eta_g^2 L}{2} \mathbb{E}[\|M_{t+1}\|^2].$$

958 Now we need to analyze the terms  $\mathbb{E}[\langle \nabla f(w_t), M_{t+1} \rangle]$  and  $\mathbb{E}[\|M_{t+1}\|^2]$ . Let's start with the inner  
959 product term.

960 The momentum update in FedEve is given by:

$$M_{t+1} = M_t + G_{kal}(\Delta \tilde{w}_t - M_t),$$

961 where  $\Delta \tilde{w}_t$  is the average gradient from the sampled clients. An important property of this update is  
962 that, under expectation, it is unbiased:

$$\mathbb{E}[\Delta \tilde{w}_t] = \nabla f(w_t).$$

963 Therefore:

$$\begin{aligned} \mathbb{E}[M_{t+1}] &= \mathbb{E}[M_t + G_{kal}(\Delta \tilde{w}_t - M_t)] \\ &= \mathbb{E}[M_t] + G_{kal}(\mathbb{E}[\Delta \tilde{w}_t] - \mathbb{E}[M_t]) \\ &= \mathbb{E}[M_t] + G_{kal}(\nabla f(w_t) - \mathbb{E}[M_t]). \end{aligned}$$

964 In the long run, this recursive relation converges to  $\mathbb{E}[M_t] = \nabla f(w_t)$ . For simplicity, we assume this  
965 has approximately been achieved, which gives us:

$$\mathbb{E}[\langle \nabla f(w_t), M_{t+1} \rangle] = \langle \nabla f(w_t), \mathbb{E}[M_{t+1}] \rangle = \langle \nabla f(w_t), \nabla f(w_t) \rangle = \|\nabla f(w_t)\|^2.$$

966 Next, we need to bound  $\mathbb{E}[\|M_{t+1}\|^2]$ . Using the update rule and the variance from Lemma F.5:

$$\begin{aligned} \mathbb{E}[\|M_{t+1}\|^2] &= \mathbb{E}[\|\mathbb{E}[M_{t+1}] + (M_{t+1} - \mathbb{E}[M_{t+1}])\|^2] \\ &= \|\mathbb{E}[M_{t+1}]\|^2 + \mathbb{E}[\|M_{t+1} - \mathbb{E}[M_{t+1}]\|^2] \\ &= \|\nabla f(w_t)\|^2 + \mathbb{E}[\|M_{t+1} - \nabla f(w_t)\|^2] \\ &\leq \|\nabla f(w_t)\|^2 + G_{kal}^2 \frac{\sigma^2}{S} \left(1 - \frac{S}{N}\right). \end{aligned}$$

967 The last inequality uses the bound on the variance of  $M_{t+1}$  derived from Lemma F.5.

968 Substituting these results back into our progress bound:

$$\begin{aligned} \mathbb{E}[f(w_{t+1})] &\leq f(w_t) - \eta_g \|\nabla f(w_t)\|^2 + \frac{\eta_g^2 L}{2} \left( \|\nabla f(w_t)\|^2 + G_{kal}^2 \frac{\sigma^2}{S} \left(1 - \frac{S}{N}\right) \right) \\ &= f(w_t) - \eta_g \|\nabla f(w_t)\|^2 + \frac{\eta_g^2 L}{2} \|\nabla f(w_t)\|^2 + \frac{\eta_g^2 L G_{kal}^2 \sigma^2}{2S} \left(1 - \frac{S}{N}\right) \\ &= f(w_t) - \eta_g \left(1 - \frac{\eta_g L}{2}\right) \|\nabla f(w_t)\|^2 + \frac{\eta_g^2 L G_{kal}^2 \sigma^2}{2S} \left(1 - \frac{S}{N}\right). \end{aligned}$$

969 The coefficient of  $\|\nabla f(w_t)\|^2$  is maximized when  $\eta_g = \frac{1}{L}$ . Setting this optimal learning rate, we get:

$$\begin{aligned} \mathbb{E}[f(w_{t+1})] &\leq f(w_t) - \frac{1}{L} \left(1 - \frac{1}{2}\right) \|\nabla f(w_t)\|^2 + \frac{1}{L^2} \frac{L G_{kal}^2 \sigma^2}{2S} \left(1 - \frac{S}{N}\right) \\ &= f(w_t) - \frac{1}{2L} \|\nabla f(w_t)\|^2 + \frac{G_{kal}^2 \sigma^2}{2LS} \left(1 - \frac{S}{N}\right). \end{aligned}$$

970 This inequality shows that in each iteration, we make progress in decreasing the objective function,  
971 but there's a constant error term due to the variance in client sampling.

972 To derive the final convergence rate, we sum this inequality over all iterations from  $t = 0$  to  $t = T - 1$ :

$$\begin{aligned} \sum_{t=0}^{T-1} \mathbb{E}[f(w_{t+1})] &\leq \sum_{t=0}^{T-1} \left( f(w_t) - \frac{1}{2L} \|\nabla f(w_t)\|^2 + \frac{G_{kal}^2 \sigma^2}{2LS} \left(1 - \frac{S}{N}\right) \right) \\ &= \sum_{t=0}^{T-1} f(w_t) - \frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(w_t)\|^2 + \frac{T G_{kal}^2 \sigma^2}{2LS} \left(1 - \frac{S}{N}\right). \end{aligned}$$

973 Note that  $\sum_{t=0}^{T-1} \mathbb{E}[f(w_{t+1})] = \sum_{t=1}^T \mathbb{E}[f(w_t)]$ . Rearranging the terms:

$$\begin{aligned}\sum_{t=1}^T \mathbb{E}[f(w_t)] &\leq \sum_{t=0}^{T-1} f(w_t) - \frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(w_t)\|^2 + \frac{TG_{kal}^2\sigma^2}{2LS} \left(1 - \frac{S}{N}\right) \\ \Rightarrow \mathbb{E}[f(w_T)] &\leq f(w_0) - \frac{1}{2L} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(w_t)\|^2] + \frac{TG_{kal}^2\sigma^2}{2LS} \left(1 - \frac{S}{N}\right).\end{aligned}$$

974 This is because the sum telescopes:  $f(w_0) + (f(w_1) - f(w_0)) + \dots + (f(w_{T-1}) - f(w_{T-1})) -$   
975  $\mathbb{E}[f(w_T)] = f(w_0) - \mathbb{E}[f(w_T)]$ .

976 Rearranging to isolate the gradient norm terms:

$$\frac{1}{2L} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(w_t)\|^2] \leq f(w_0) - \mathbb{E}[f(w_T)] + \frac{TG_{kal}^2\sigma^2}{2LS} \left(1 - \frac{S}{N}\right).$$

977 Since  $f$  is bounded below by some value  $f^*$  (which could be the global minimum), we have  
978  $\mathbb{E}[f(w_T)] \geq f^*$ . Thus:

$$\frac{1}{2L} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(w_t)\|^2] \leq f(w_0) - f^* + \frac{TG_{kal}^2\sigma^2}{2LS} \left(1 - \frac{S}{N}\right).$$

979 Dividing both sides by  $T$ :

$$\begin{aligned}\frac{1}{2LT} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(w_t)\|^2] &\leq \frac{f(w_0) - f^*}{T} + \frac{G_{kal}^2\sigma^2}{2LS} \left(1 - \frac{S}{N}\right) \\ \Rightarrow \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(w_t)\|^2] &\leq \frac{2L(f(w_0) - f^*)}{T} + \frac{G_{kal}^2\sigma^2}{S} \left(1 - \frac{S}{N}\right).\end{aligned}$$

980 This final bound can be expressed in big-O notation as:

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(w_t)\|^2] \leq \mathcal{O}\left(\frac{L(f(w_0) - f^*)}{T} + \frac{G_{kal}^2\sigma^2}{S} \left(1 - \frac{S}{N}\right)\right).$$

981 This result demonstrates several important properties of the FedEve algorithm:

- 982 • The first term  $\mathcal{O}\left(\frac{L(f(w_0) - f^*)}{T}\right)$  shows that the convergence rate is  $\mathcal{O}(1/T)$ , which is  
983 optimal for first-order methods on non-convex functions.
- 984 • The second term  $\mathcal{O}\left(\frac{G_{kal}^2\sigma^2}{S} \left(1 - \frac{S}{N}\right)\right)$  represents the irreducible error due to client sampling  
985 and stochastic gradients.
- 986 • This error decreases as we increase the number of sampled clients  $S$ , and vanishes completely  
987 when  $S = N$  (i.e., when we use all clients).
- 988 • The Kalman gain parameter  $G_{kal}$  appears quadratically in the error term, showing that  
989 smaller values of  $G_{kal}$  can reduce the impact of stochasticity, but at the cost of potentially  
990 slower adaptation to changes in the gradient.

991 In practical implementations, the parameters  $\eta_g$  and  $G_{kal}$  can be tuned to balance the trade-off  
992 between convergence speed and stability based on the specific characteristics of the federated learning  
993 task.  $\square$

994 **G Generalization bound**

995 This Generalization bound is inspired by Sun et al. [2023].

996 **G.1 Theorems**

997 **Theorem G.1** (On-average Algorithmic Stability). *Suppose a federated learning algorithm  $\mathcal{A}$  is  
998  $\epsilon$ -on-averagely stable. Then,*

$$\epsilon_{gen} \leq \mathbb{E}_{\mathcal{A}, \mathcal{S}} \left[ |f(\mathcal{A}(\mathcal{S})) - \hat{f}_{\mathcal{S}}(\mathcal{A}(\mathcal{S}))| \right] \leq \epsilon.$$

999 **Theorem G.2** (Generalization Bound). *Suppose Assumptions G.3-G.5 hold and  $\eta_l \leq \frac{1}{L_m K(t+1)}$ .  
1000 Then,*

$$\epsilon_{gen} \leq \mathcal{O} \left( \frac{L_p}{n L_m} \right) \left[ T(\sigma_l + \sigma_g) + \sqrt{T L_m \Delta_0} + T \sqrt{\frac{G_{kal}^2 \sigma^2}{S}} \right].$$

1001 **G.2 Assumptions**

1002 **Assumption G.3** (Lipschitz Continuity). The loss function  $l(\cdot, z)$  is  $L_p$ -Lipschitz continuous, that is,  
1003  $|l(w; z) - l(w'; z)| \leq L_p \|w - w'\|$ , and is  $L_m$ -smooth for any  $z$ , that is,  $\|\nabla l(w; z) - \nabla l(w'; z)\| \leq  
1004 L_m \|w - w'\|$  for any  $z, w, w'$ .

1005 **Assumption G.4** (Bounded Variance). The function  $F_i$  have  $\sigma_l$ -bounded (local) variance i.e.,  
1006  $\mathbb{E}\|f_i(w) - \nabla f_i(w)\| \leq \sigma_l$  for all  $w \in \mathcal{U}^d$ ,  $j \in [d]$  and  $i \in [m]$ . Furthermore, we assume the  
1007 (global) variance is bounded,  $\mathbb{E}\|f_i(w) - \nabla f(w)\| \leq \sigma_g$  for all  $x \in \mathcal{U}^d$  and  $j \in [d]$ .

1008 **Assumption G.5** (Bounded Gradients). The function  $f_i(x, z)$  have  $G$ -bounded gradients i.e., for any  
1009  $i \in [m]$ ,  $x \in \mathcal{U}^d$  and  $z \in \mathcal{Z}$  we have  $|\nabla f_i(x, z)| \leq G$  for all  $j \in [d]$ .

1010 **G.3 Lemmas**

1011 **Lemma G.6** (Bounded Local Updates). *Suppose Assumptions G.3-G.5 hold. For any step-size, we  
1012 can bound the local updates as*

$$\mathbb{E}\|w_{i,k} - w_t\| \leq \frac{(1 + \eta_l L_m)^k - 1}{L_m} (\mathbb{E}\|\nabla f(w_t)\| + \sigma_l + \sigma_g).$$

1013 where  $w_{i,k}$  is the model parameters of client  $i$  at  $k$ -th local updates.

1014 **Lemma G.7** (Bounded Local Gradients). *Given Assumptions G.3-G.5. For any step-size, we can  
1015 bound the local gradients as  $\mathbb{E}\|f_i(w_{i,k})\| \leq (1 + \eta_l L_m)^k (\mathbb{E}\|\nabla f(w_t)\| + \sigma_l + \sigma_g)$ , where  $f_i(\cdot)$  is  
1016 the sampled gradient of client  $i$ .*

1017 **Lemma G.8** (Bounded Global Model with Sample Perturbation). *Given Assumptions G.3-  
1018 G.5. For any step-size, we can bound the local gradients as  $\mathbb{E}\|w_T - w'_T\| \leq  
1019 \sum_{t=0}^{T-1} \frac{2e^{\eta_l K(t+1)L_m}}{n L_m} (\mathbb{E}\|\nabla f(w_t)\| + \sigma_l + \sigma_g)$ , where  $w'_T$  is the model parameters with sample per-  
1020 turbation at  $T$ -th communication rounds.*

1021 **G.4 Proofs**

1022 **G.4.1 Proof of Lemma G.6**

1023 *Proof.* Bounding Local Updates:

$$\begin{aligned} & \mathbb{E} \|w_{i,k+1} - w_t\| \\ &= \mathbb{E} \|w_{i,k} - \eta_l f_i(w_{i,k}) - w_t\| \\ &\leq \mathbb{E} \|w_{i,k} - w_t - \eta_l (f_i(w_{i,k}) - f_i(w_t))\| + \eta_l \mathbb{E} \|f_i(w_t)\| \\ &\leq (1 + \eta_l L_m) \mathbb{E} \|w_{i,k} - w_t\| + \eta_l \mathbb{E} \|f_i(w_t)\| \\ &\leq (1 + \eta_l L_m) \mathbb{E} \|w_{i,k} - w_t\| + \eta_l (\mathbb{E} \|f_i(w_t) - \nabla f_i(w_t)\| + \mathbb{E} \|\nabla f_i(w_t) - \nabla f(w_t)\| + \mathbb{E} \|\nabla f(w_t)\|) \\ &\leq (1 + \eta_l L_m) \mathbb{E} \|w_{i,k} - w_t\| + \eta_l (\sigma_l + \sigma_g + \mathbb{E} \|\nabla f(w_t)\|), \end{aligned}$$

1024 unrolling the above and noting  $w_{i,0} = w_t$  yields

$$\mathbb{E}\|w_{i,k} - w_t\| \leq \frac{(1 + \eta_l L_m)^k - 1}{L_m} (\mathbb{E}\|\nabla f(w_t)\| + \sigma_l + \sigma_g).$$

1025  $\square$

1026 **G.4.2 Proof of Lemma G.7**

1027 *Proof.* Bounding Local Gradients:

$$\begin{aligned}\mathbb{E}\|f_i(w_{i,k})\| &= \mathbb{E}\|f_i(w_{i,k}) - \nabla f_i(w_{i,k}) + \nabla f_i(w_{i,k}) - \nabla f(w_t) + \nabla f(w_t)\| \\ &\leq \mathbb{E}\|f_i(w_{i,k}) - \nabla f_i(w_{i,k})\| + \mathbb{E}\|\nabla f_i(w_{i,k}) - \nabla f(w_t)\| + \mathbb{E}\|\nabla f(w_t)\| \\ &\leq \sigma_l + L_m \mathbb{E}\|w_{i,k} - w_t\| + \mathbb{E}\|\nabla f(w_t)\|,\end{aligned}$$

based on Lemma G.6, we obtain:

$$\begin{aligned}&\leq \sigma_l + \mathbb{E}\|\nabla f(w_t)\| \\ &+ ((1 + \eta_l L_m)^k - 1) (\mathbb{E}\|\nabla f(w_t)\| + \sigma_l + \sigma_g) \\ &\leq (1 + \eta_l L_m)^k (\mathbb{E}\|\nabla f(w_t)\| + \sigma_l + \sigma_g).\end{aligned}$$

1028  $\square$

1029 **G.4.3 Proof of Lemma G.8**

1030 *Proof.* Given time index  $t$  and for client  $j$  with  $j \neq i$ , we have

$$\begin{aligned}\mathbb{E}\|w_{j,k+1} - w'_{j,k+1}\| &= \mathbb{E}\|w_{j,k} - w'_{j,k} - \eta_l(g_j(w_{j,k}) - g_j(w'_{j,k}))\| \\ &\leq (1 + \eta_l L_m) \mathbb{E}\|w_{j,k} - w'_{j,k}\|.\end{aligned}$$

1031 And unrolling it gives

$$\mathbb{E}\|w_{j,K} - w'_{j,K}\| \leq e^{\eta_l K L_m} \mathbb{E}\|w_t - w'_t\|, \quad \forall j \neq i,$$

1032 since  $1 + x < e^x$ . For client  $i$ , there are two cases to consider. In the first case, SGD selects  
1033 non-perturbed samples in  $\mathcal{S}$  and  $\mathcal{S}^{(i)}$ , which happens with probability  $1 - 1/n_i$ . Then, we have  
1034  $\|w_{i,k+1} - w'_{i,k+1}\| \leq (1 + \eta_l L_m) \|w_{i,k} - w'_{i,k}\|$ . In the second case, SGD encounters the perturbed  
1035 sample at time step  $k$ , which happens with probability  $1/n_i$ . Then, we have

$$\begin{aligned}\|w_{i,k+1} - w'_{i,k+1}\| &= \|w_{i,k} - w'_{i,k} - \eta_l(f_i(w_{i,k}) - g'_i(w'_{i,k}))\| \\ &\leq \|w_{i,k} - w'_{i,k} - \eta_l(f_i(w_{i,k}) - f_i(w'_{i,k}))\| + \eta_l \|f_i(w'_{i,k}) - g'_i(w'_{i,k})\| \\ &\leq (1 + \eta_l L_m) \|w_{i,k} - w'_{i,k}\| + \eta_l \|f_i(w'_{i,k}) - g'_i(w'_{i,k})\|.\end{aligned}$$

1036 Combining these two cases for client  $i$  we have

$$\begin{aligned}\mathbb{E}\|w_{i,k+1} - w'_{i,k+1}\| &\leq (1 + \eta_l L_m) \mathbb{E}\|w_{i,k} - w'_{i,k}\| + \frac{\eta_l}{n_i} \mathbb{E}\|f_i(w'_{i,k}) - g'_i(w'_{i,k})\| \\ &\leq (1 + \eta_l L_m) \mathbb{E}\|w_{i,k} - w'_{i,k}\| + \frac{2\eta_l}{n_i} \mathbb{E}\|f_i(w_{i,k})\|,\end{aligned}$$

based on Lemma G.7, we obtain:

$$\begin{aligned}&\leq (1 + \eta_l L_m) \mathbb{E}\|w_{i,k} - w'_{i,k}\| \\ &+ \frac{2\eta_l}{n_i} e^{\eta_l k L_m} (\mathbb{E}\|\nabla f(w_t)\| + \sigma_l + \sigma_g),\end{aligned}$$

1037 then unrolling it gives

$$\begin{aligned}&\mathbb{E}\|w_{i,K} - w'_{i,K}\| \\ &\leq e^{\eta_l K L_m} \mathbb{E}\|w_t - w'_t\| \\ &+ \frac{2e^{\eta_l K L_m}}{n_i L_m} (\mathbb{E}\|\nabla f(w_t)\| + \sigma_l + \sigma_g) \quad \forall j = i.\end{aligned}\tag{43}$$

1038 Combines 43 and 43 we have

$$\begin{aligned}\mathbb{E}\|w_{t+1} - w'_{t+1}\| &\leq \sum_{i=1}^m p_i \mathbb{E}\|w_{i,K} - w'_{i,K}\| \\ &\leq e^{\eta_l K L_m} \mathbb{E}\|w_t - w'_t\| + \frac{2e^{\eta_l K L_m}}{n L_m} (\mathbb{E}\|\nabla f(w_t)\| + \sigma_l + \sigma_g)\end{aligned}$$

1039 where we also use  $p_i = n_i/n$  in the last step. Further, unrolling the above over  $t$  and noting  $w_0 = w'_0$ ,  
1040 we obtain

$$\mathbb{E}\|w_T - w'_T\| \leq \sum_{t=0}^{T-1} \frac{2e^{\eta_l K(t+1)L_m}}{n L_m} (\mathbb{E}\|\nabla f(w_t)\| + \sigma_l + \sigma_g).$$

1041  $\square$

1042 **G.4.4 Proof of Theorem G.2**

1043 *Proof.* According to the fact that:

$$\left( \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(w_t)\| \right)^2 \leq T \sum_{t=0}^{T-1} (\mathbb{E} \|\nabla f(w_t)\|)^2 \leq T \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(w_t)\|^2,$$

1044 where the second inequality follows Jensen's inequality, and the convergence analysis of FedEve:

$$\frac{1}{T} \sum_{r=0}^{T-1} \mathbb{E} \|\nabla f(w_t)\|^2 \leq \mathcal{O} \left( \frac{L_m \Delta_0}{T} + \frac{G_{kal}^2 \sigma^2}{S} \left( 1 - \frac{S}{N} \right) \right),$$

1045 where  $\Delta_0 := \mathbb{E}[f(w_0) - f(w^*)]$ . The generalization bound is

$$\begin{aligned} \epsilon_{gen} &\leq L_p \mathbb{E} \|w_T - w'_T\| \\ &\leq L_p \sum_{t=0}^{T-1} \frac{2e^{\eta_l K(t+1)L_m}}{nL_m} (\mathbb{E} \|\nabla f(w_t)\| + \sigma_l + \sigma_g), \end{aligned}$$

when  $\eta_l < \frac{1}{K(t+1)L_m}$ , we obtain:

$$\begin{aligned} &\leq L_p \sum_{t=0}^{T-1} \frac{2e^{\eta_l K(t+1)L_m}}{nL_m} \left( \sqrt{T \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(w_t)\|^2} + T(\sigma_l + \sigma_g) \right), \\ &\leq L_p \sum_{t=0}^{T-1} \frac{2e^{\eta_l K(t+1)L_m}}{nL_m} \left( \sqrt{T \mathcal{O} \left( \frac{L_m \Delta_0}{T} + \frac{G_{kal}^2 \sigma^2}{S} \left( 1 - \frac{S}{N} \right) \right)} + T(\sigma_l + \sigma_g) \right), \\ &\leq L_p \sum_{t=0}^{T-1} \frac{2e^{\eta_l K(t+1)L_m}}{nL_m} \left( \sqrt{L_m \Delta_0 + T \frac{G_{kal}^2 \sigma^2}{S} \left( 1 - \frac{S}{N} \right)} + T(\sigma_l + \sigma_g) \right), \\ &\leq L_p \sum_{t=0}^{T-1} \frac{2e^{\eta_l K(t+1)L_m}}{nL_m} \left( \sqrt{L_m \Delta_0} + \sqrt{T \frac{G_{kal}^2 \sigma^2}{S} \left( 1 - \frac{S}{N} \right)} + T(\sigma_l + \sigma_g) \right), \\ &\leq \mathcal{O} \left( \frac{L_p}{nL_m} \right) \left[ T(\sigma_l + \sigma_g) + \sqrt{TL_m \Delta_0} + T \sqrt{\frac{G_{kal}^2 \sigma^2}{S} \left( 1 - \frac{S}{N} \right)} \right]. \end{aligned}$$

1046

□

1047 **NeurIPS Paper Checklist**

1048 The checklist is designed to encourage best practices for responsible machine learning research,  
1049 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
1050 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should  
1051 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
1052 towards the page limit.

1053 Please read the checklist guidelines carefully for information on how to answer these questions. For  
1054 each question in the checklist:

- 1055 • You should answer **[Yes]**, **[No]**, or **[NA]**.
- 1056 • **[NA]** means either that the question is Not Applicable for that particular paper or the  
1057 relevant information is Not Available.
- 1058 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

1059 **The checklist answers are an integral part of your paper submission.** They are visible to the  
1060 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it  
1061 (after eventual revisions) with the final version of your paper, and its final version will be published  
1062 with the paper.

1063 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
1064 While "**[Yes]**" is generally preferable to "**[No]**", it is perfectly acceptable to answer "**[No]**" provided a  
1065 proper justification is given (e.g., "error bars are not reported because it would be too computationally  
1066 expensive" or "we were unable to find the license for the dataset we used"). In general, answering  
1067 "**[No]**" or "**[NA]**" is not grounds for rejection. While the questions are phrased in a binary way, we  
1068 acknowledge that the true answer is often more nuanced, so please just use your best judgment and  
1069 write a justification to elaborate. All supporting evidence can appear either in the main paper or the  
1070 supplemental material, provided in appendix. If you answer **[Yes]** to a question, in the justification  
1071 please point to the section(s) where related material for the question can be found.

1072 **IMPORTANT**, please:

- 1073 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”.**
- 1074 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 1075 • **Do not modify the questions and only use the provided macros for your answers.**

1076 **1. Claims**

1077 Question: Do the main claims made in the abstract and introduction accurately reflect the  
1078 paper’s contributions and scope?

1079 Answer: **[Yes]**

1080 Justification: The abstract and introduction clearly state the paper’s contributions, including:  
1081 identifying and defining period drift, proposing a predict-observe framework, developing  
1082 FedEve as an implementation, and providing theoretical and empirical evidence of its  
1083 effectiveness.

1084 Guidelines:

- 1085 • The answer NA means that the abstract and introduction do not include the claims  
1086 made in the paper.
- 1087 • The abstract and/or introduction should clearly state the claims made, including the  
1088 contributions made in the paper and important assumptions and limitations. A No or  
1089 NA answer to this question will not be perceived well by the reviewers.
- 1090 • The claims made should match theoretical and experimental results, and reflect how  
1091 much the results can be expected to generalize to other settings.
- 1092 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
1093 are not attained by the paper.

1094 **2. Limitations**

1095 Question: Does the paper discuss the limitations of the work performed by the authors?

1096 Answer: **[Yes]**

1097 Justification: The paper discusses limitations in Appendix, particularly with respect to  
1098 assumptions of independence between model parameters and drift, and provides justifications  
1099 for these assumptions.

1100 Guidelines:

- 1101 • The answer NA means that the paper has no limitation while the answer No means that  
1102 the paper has limitations, but those are not discussed in the paper.
- 1103 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1104 • The paper should point out any strong assumptions and how robust the results are to  
1105 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
1106 model well-specification, asymptotic approximations only holding locally). The authors  
1107 should reflect on how these assumptions might be violated in practice and what the  
1108 implications would be.
- 1109 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
1110 only tested on a few datasets or with a few runs. In general, empirical results often  
1111 depend on implicit assumptions, which should be articulated.
- 1112 • The authors should reflect on the factors that influence the performance of the approach.  
1113 For example, a facial recognition algorithm may perform poorly when image resolution  
1114 is low or images are taken in low lighting. Or a speech-to-text system might not be  
1115 used reliably to provide closed captions for online lectures because it fails to handle  
1116 technical jargon.
- 1117 • The authors should discuss the computational efficiency of the proposed algorithms  
1118 and how they scale with dataset size.
- 1119 • If applicable, the authors should discuss possible limitations of their approach to  
1120 address problems of privacy and fairness.
- 1121 • While the authors might fear that complete honesty about limitations might be used by  
1122 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
1123 limitations that aren't acknowledged in the paper. The authors should use their best  
1124 judgment and recognize that individual actions in favor of transparency play an impor-  
1125 tant role in developing norms that preserve the integrity of the community. Reviewers  
1126 will be specifically instructed to not penalize honesty concerning limitations.

### 1127 3. Theory assumptions and proofs

1128 Question: For each theoretical result, does the paper provide the full set of assumptions and  
1129 a complete (and correct) proof?

1130 Answer: [Yes]

1131 Justification: The paper provides clear assumptions (Assumption 1 and 2 in Section 3.2  
1132 and 3.3) for the theoretical results, with complete proofs in the Appendix (e.g., Appendix  
1133 sections on noise independence and Bayesian filter proofs).

1134 Guidelines:

- 1135 • The answer NA means that the paper does not include theoretical results.
- 1136 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
1137 referenced.
- 1138 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 1139 • The proofs can either appear in the main paper or the supplemental material, but if  
1140 they appear in the supplemental material, the authors are encouraged to provide a short  
1141 proof sketch to provide intuition.
- 1142 • Inversely, any informal proof provided in the core of the paper should be complemented  
1143 by formal proofs provided in appendix or supplemental material.
- 1144 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 1145 4. Experimental result reproducibility

1146 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
1147 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
1148 of the paper (regardless of whether the code and data are provided or not)?

1149 Answer: [Yes]

1150 Justification: The paper provides comprehensive details on datasets (FEMNIST, CIFAR-  
1151 10/100, MovieLens), models (LeNet5, ResNet-18, DIN), training settings (rounds, clients  
1152 per round, epochs), and algorithm parameters in Section 4.1 and Appendix.  
1153

1154 Guidelines:

- 1155 • The answer NA means that the paper does not include experiments.
- 1156 • If the paper includes experiments, a No answer to this question will not be perceived  
1157 well by the reviewers: Making the paper reproducible is important, regardless of  
whether the code and data are provided or not.
- 1158 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
1159 to make their results reproducible or verifiable.
- 1160 • Depending on the contribution, reproducibility can be accomplished in various ways.  
1161 For example, if the contribution is a novel architecture, describing the architecture fully  
1162 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
1163 be necessary to either make it possible for others to replicate the model with the same  
1164 dataset, or provide access to the model. In general, releasing code and data is often  
1165 one good way to accomplish this, but reproducibility can also be provided via detailed  
1166 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
1167 of a large language model), releasing of a model checkpoint, or other means that are  
1168 appropriate to the research performed.
- 1169 • While NeurIPS does not require releasing code, the conference does require all submissions  
1170 to provide some reasonable avenue for reproducibility, which may depend on the  
1171 nature of the contribution. For example
  - 1172 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
1173 to reproduce that algorithm.
  - 1174 (b) If the contribution is primarily a new model architecture, the paper should describe  
1175 the architecture clearly and fully.
  - 1176 (c) If the contribution is a new model (e.g., a large language model), then there should  
1177 either be a way to access this model for reproducing the results or a way to reproduce  
1178 the model (e.g., with an open-source dataset or instructions for how to construct  
1179 the dataset).
  - 1180 (d) We recognize that reproducibility may be tricky in some cases, in which case  
1181 authors are welcome to describe the particular way they provide for reproducibility.  
1182 In the case of closed-source models, it may be that access to the model is limited in  
1183 some way (e.g., to registered users), but it should be possible for other researchers  
1184 to have some path to reproducing or verifying the results.

1185 **5. Open access to data and code**

1186 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
1187 tions to faithfully reproduce the main experimental results, as described in supplemental  
1188 material?

1189 Answer: [No]

1190 Justification: While the paper uses publicly available datasets (FEMNIST, CIFAR-10/100,  
1191 MovieLens) with appropriate citations and links, there is no explicit mention of code  
1192 availability for the proposed FedEve algorithm.

1193 Guidelines:

- 1194 • The answer NA means that paper does not include experiments requiring code.
- 1195 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
1196 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1197 • While we encourage the release of code and data, we understand that this might not be  
1198 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
1199 including code, unless this is central to the contribution (e.g., for a new open-source  
1200 benchmark).
- 1201 • The instructions should contain the exact command and environment needed to run to  
1202 reproduce the results. See the NeurIPS code and data submission guidelines ([https://nips.cc/  
1203 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.

- 1204           • The authors should provide instructions on data access and preparation, including how  
 1205            to access the raw data, preprocessed data, intermediate data, and generated data, etc.  
 1206           • The authors should provide scripts to reproduce all experimental results for the new  
 1207            proposed method and baselines. If only a subset of experiments are reproducible, they  
 1208            should state which ones are omitted from the script and why.  
 1209           • At submission time, to preserve anonymity, the authors should release anonymized  
 1210            versions (if applicable).  
 1211           • Providing as much information as possible in supplemental material (appended to the  
 1212            paper) is recommended, but including URLs to data and code is permitted.

1213       **6. Experimental setting/details**

1214       Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
 1215           parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
 1216           results?

1217       Answer: [Yes]

1218       Justification: The paper details data splits (e.g., Dirichlet distribution for non-iid settings),  
 1219           hyperparameters (learning rates, epochs, communication rounds), optimizers (SGD for CV,  
 1220           Adam for RS), and evaluation protocols in Section 4.1 and Appendix.

1221       Guidelines:

- 1222           • The answer NA means that the paper does not include experiments.
- 1223           • The experimental setting should be presented in the core of the paper to a level of detail  
 1224            that is necessary to appreciate the results and make sense of them.
- 1225           • The full details can be provided either with the code, in appendix, or as supplemental  
 1226            material.

1227       **7. Experiment statistical significance**

1228       Question: Does the paper report error bars suitably and correctly defined or other appropriate  
 1229           information about the statistical significance of the experiments?

1230       Answer: [Yes]

1231       Justification: The paper reports standard deviations for all experimental results as shown in  
 1232           Tables 1, 2, and 3, indicating statistical significance across multiple runs of each experiment.

1233       Guidelines:

- 1234           • The answer NA means that the paper does not include experiments.
- 1235           • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
 1236            dence intervals, or statistical significance tests, at least for the experiments that support  
 1237            the main claims of the paper.
- 1238           • The factors of variability that the error bars are capturing should be clearly stated (for  
 1239            example, train/test split, initialization, random drawing of some parameter, or overall  
 1240            run with given experimental conditions).
- 1241           • The method for calculating the error bars should be explained (closed form formula,  
 1242            call to a library function, bootstrap, etc.)
- 1243           • The assumptions made should be given (e.g., Normally distributed errors).
- 1244           • It should be clear whether the error bar is the standard deviation or the standard error  
 1245            of the mean.
- 1246           • It is OK to report 1-sigma error bars, but one should state it. The authors should  
 1247            preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
 1248            of Normality of errors is not verified.
- 1249           • For asymmetric distributions, the authors should be careful not to show in tables or  
 1250            figures symmetric error bars that would yield results that are out of range (e.g. negative  
 1251            error rates).
- 1252           • If error bars are reported in tables or plots, The authors should explain in the text how  
 1253            they were calculated and reference the corresponding figures or tables in the text.

1254       **8. Experiments compute resources**

1255 Question: For each experiment, does the paper provide sufficient information on the com-  
1256 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
1257 the experiments?

1258 Answer: [No]

1259 Justification: The paper does not explicitly mention the hardware specifications, memory  
1260 requirements, or running time for the experiments.

1261 Guidelines:

- 1262 • The answer NA means that the paper does not include experiments.  
1263 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
1264 or cloud provider, including relevant memory and storage.  
1265 • The paper should provide the amount of compute required for each of the individual  
1266 experimental runs as well as estimate the total compute.  
1267 • The paper should disclose whether the full research project required more compute  
1268 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
1269 didn't make it into the paper).

## 9. Code of ethics

1271 Question: Does the research conducted in the paper conform, in every respect, with the  
1272 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

1273 Answer: [Yes]

1274 Justification: The research focuses on improving federated learning algorithms which  
1275 inherently promote privacy protection. It uses standard benchmark datasets and does not  
1276 present any apparent ethical concerns.

1277 Guidelines:

- 1278 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.  
1279 • If the authors answer No, they should explain the special circumstances that require a  
1280 deviation from the Code of Ethics.  
1281 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
1282 eration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

1284 Question: Does the paper discuss both potential positive societal impacts and negative  
1285 societal impacts of the work performed?

1286 Answer: [No]

1287 Justification: While the paper focuses on improving federated learning, which has positive  
1288 privacy implications, it does not explicitly discuss broader societal impacts, either positive  
1289 or negative.

1290 Guidelines:

- 1291 • The answer NA means that there is no societal impact of the work performed.  
1292 • If the authors answer NA or No, they should explain why their work has no societal  
1293 impact or why the paper does not address societal impact.  
1294 • Examples of negative societal impacts include potential malicious or unintended uses  
1295 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
1296 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
1297 groups), privacy considerations, and security considerations.  
1298 • The conference expects that many papers will be foundational research and not tied  
1299 to particular applications, let alone deployments. However, if there is a direct path to  
1300 any negative applications, the authors should point it out. For example, it is legitimate  
1301 to point out that an improvement in the quality of generative models could be used to  
1302 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
1303 that a generic algorithm for optimizing neural networks could enable people to train  
1304 models that generate Deepfakes faster.

- 1305           • The authors should consider possible harms that could arise when the technology is  
 1306           being used as intended and functioning correctly, harms that could arise when the  
 1307           technology is being used as intended but gives incorrect results, and harms following  
 1308           from (intentional or unintentional) misuse of the technology.  
 1309           • If there are negative societal impacts, the authors could also discuss possible mitigation  
 1310           strategies (e.g., gated release of models, providing defenses in addition to attacks,  
 1311           mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
 1312           feedback over time, improving the efficiency and accessibility of ML).

1313           **11. Safeguards**

1314           Question: Does the paper describe safeguards that have been put in place for responsible  
 1315           release of data or models that have a high risk for misuse (e.g., pretrained language models,  
 1316           image generators, or scraped datasets)?

1317           Answer: [NA]

1318           Justification: The paper does not release any high-risk models or datasets. It uses standard  
 1319           benchmark datasets and proposes an optimization algorithm for federated learning which  
 1320           does not pose risks of misuse.

1321           Guidelines:

- 1322           • The answer NA means that the paper poses no such risks.  
 1323           • Released models that have a high risk for misuse or dual-use should be released with  
 1324           necessary safeguards to allow for controlled use of the model, for example by requiring  
 1325           that users adhere to usage guidelines or restrictions to access the model or implementing  
 1326           safety filters.  
 1327           • Datasets that have been scraped from the Internet could pose safety risks. The authors  
 1328           should describe how they avoided releasing unsafe images.  
 1329           • We recognize that providing effective safeguards is challenging, and many papers do  
 1330           not require this, but we encourage authors to take this into account and make a best  
 1331           faith effort.

1332           **12. Licenses for existing assets**

1333           Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
 1334           the paper, properly credited and are the license and terms of use explicitly mentioned and  
 1335           properly respected?

1336           Answer: [Yes]

1337           Justification: The paper properly cites and links to all datasets used (FEMNIST, CIFAR-  
 1338           10/100, MovieLens) and acknowledges the original creators through standard citations.

1339           Guidelines:

- 1340           • The answer NA means that the paper does not use existing assets.  
 1341           • The authors should cite the original paper that produced the code package or dataset.  
 1342           • The authors should state which version of the asset is used and, if possible, include a  
 1343           URL.  
 1344           • The name of the license (e.g., CC-BY 4.0) should be included for each asset.  
 1345           • For scraped data from a particular source (e.g., website), the copyright and terms of  
 1346           service of that source should be provided.  
 1347           • If assets are released, the license, copyright information, and terms of use in the  
 1348           package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
 1349           has curated licenses for some datasets. Their licensing guide can help determine the  
 1350           license of a dataset.  
 1351           • For existing datasets that are re-packaged, both the original license and the license of  
 1352           the derived asset (if it has changed) should be provided.  
 1353           • If this information is not available online, the authors are encouraged to reach out to  
 1354           the asset's creators.

1355           **13. New assets**

1356           Question: Are new assets introduced in the paper well documented and is the documentation  
 1357           provided alongside the assets?

1358                  Answer: [NA]

1359                  Justification: The paper does not introduce new datasets or models, but rather proposes  
1360                  a new algorithm (FedEve) based on existing FL paradigms and evaluates it on standard  
1361                  datasets.

1362                  Guidelines:

- 1363                  • The answer NA means that the paper does not release new assets.
- 1364                  • Researchers should communicate the details of the dataset/code/model as part of their  
1365                  submissions via structured templates. This includes details about training, license,  
1366                  limitations, etc.
- 1367                  • The paper should discuss whether and how consent was obtained from people whose  
1368                  asset is used.
- 1369                  • At submission time, remember to anonymize your assets (if applicable). You can either  
1370                  create an anonymized URL or include an anonymized zip file.

#### 1371                  14. Crowdsourcing and research with human subjects

1372                  Question: For crowdsourcing experiments and research with human subjects, does the paper  
1373                  include the full text of instructions given to participants and screenshots, if applicable, as  
1374                  well as details about compensation (if any)?

1375                  Answer: [NA]

1376                  Justification: The paper does not involve crowdsourcing or research with human subjects. It  
1377                  uses pre-existing datasets and simulation-based evaluations.

1378                  Guidelines:

- 1379                  • The answer NA means that the paper does not involve crowdsourcing nor research with  
1380                  human subjects.
- 1381                  • Including this information in the supplemental material is fine, but if the main contribu-  
1382                  tion of the paper involves human subjects, then as much detail as possible should be  
1383                  included in the main paper.
- 1384                  • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
1385                  or other labor should be paid at least the minimum wage in the country of the data  
1386                  collector.

#### 1387                  15. Institutional review board (IRB) approvals or equivalent for research with human 1388                  subjects

1389                  Question: Does the paper describe potential risks incurred by study participants, whether  
1390                  such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1391                  approvals (or an equivalent approval/review based on the requirements of your country or  
1392                  institution) were obtained?

1393                  Answer: [NA]

1394                  Justification: The research does not involve human subjects or study participants, so IRB  
1395                  approval was not required for this work.

1396                  Guidelines:

- 1397                  • The answer NA means that the paper does not involve crowdsourcing nor research with  
1398                  human subjects.
- 1399                  • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1400                  may be required for any human subjects research. If you obtained IRB approval, you  
1401                  should clearly state this in the paper.
- 1402                  • We recognize that the procedures for this may vary significantly between institutions  
1403                  and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1404                  guidelines for their institution.
- 1405                  • For initial submissions, do not include any information that would break anonymity (if  
1406                  applicable), such as the institution conducting the review.

#### 1407                  16. Declaration of LLM usage

1408 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
1409 non-standard component of the core methods in this research? Note that if the LLM is used  
1410 only for writing, editing, or formatting purposes and does not impact the core methodology,  
1411 scientific rigorosity, or originality of the research, declaration is not required.

1412 Answer: [NA]

1413 Justification: The paper does not use large language models as part of its core methodology  
1414 or experiments. The proposed FedEve method is based on mathematical formulations and  
1415 traditional ML techniques.

1416 Guidelines:

- 1417 • The answer NA means that the core method development in this research does not  
1418 involve LLMs as any important, original, or non-standard components.  
1419 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)  
1420 for what should or should not be described.