INTERACTION BASED GAUSSIAN WEIGHTING CLUSTERING FOR FEDERATED LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Federated learning emerged as a decentralized paradigm to train models while securing privacy. However, conventional FL faces data heterogeneity and class imbalance challenges, affecting model performance. In response to these issues, Personalized FL has been developed as an innovative methodology that relies on fine-tuning the distinct local models based on individual training datasets. In this work, we propose a novel PFL method, FedGWC (Federated Gaussian Weighting), which groups clients based on their data distribution, allowing training of a more robust and personalized model on the identified clusters. FedGWC identifies homogeneous clusters by transforming individual empirical losses to model client interactions with a Gaussian reward mechanism. Additionally, we introduce a new clustering metric for FL to evaluate cluster cohesion with respect to the individual class distribution. Our experiments on benchmark datasets show that FedGWC outperforms existing FL algorithms in cluster quality and classification accuracy, validating the efficacy of our approach.

023

021

004

006

008

009

010

011

012

013

014

015

016

017

018

019

1 INTRODUCTION

025 Federated Learning (FL) has emerged as a promising paradigm for training models on decentralized 026 data while preserving privacy and improving communication efficiency. Unlike traditional machine 027 learning frameworks, FL enables collaborative training between multiple clients without requiring raw data transfer, making it particularly attractive in privacy-sensitive domains (McMahan et al., 2017; Bonawitz et al., 2019). FL was introduced primarily to address two major challenges in de-029 centralized scenarios: ensuring privacy (Kairouz et al., 2021) and reducing communication overhead (Hamer et al., 2020; Asad et al., 2020). In particular, FL algorithms must guarantee *communication* 031 efficiency to reduce the burden associated with the exchange of model updates between clients and the central server, while maintaining strong *privacy guarantees* is also essential, as clients should 033 not expose their private data during the training process. 034

A fundamental challenge in FL is given by data heterogeneity (Li et al., 2020), specifically in the form of class imbalance within individual clients and non-IID distributions throughout the federation. In this scenario, a single global model often fails to generalize due to clients contributing 037 updates from skewed distributions, leading to degraded performance (Zhao et al., 2018; Caldarola 038 et al., 2022) compared to the centralized counterpart. Furthermore, noisy or corrupted data from 039 some clients can further complicate the learning process (Cao et al., 2020; Zhang et al., 2022), while non-IID data often results in unstable convergence and conflicting gradient updates (Hsieh et al., 040 2020; Zhao et al., 2018). Despite the introduction of various techniques to mitigate these issues, 041 such as regularization methods (Li et al., 2020), momentum (Mendieta et al., 2022), and control 042 variates (Karimireddy et al., 2020b), statistical heterogeneity remains a critical unsolved problem. 043

Recently, Personalized FL (PFL) techniques have emerged to address the limitations of learning 044 a single global model for the entire federation, shifting the focus to a more flexible and cooperative strategy (Zhang et al., 2021). These approaches seek to maximize both the shared information 046 between clients and the specific data characteristics of individual clients, allowing the creation of 047 personalized models that adapt better to the unique distributions present in the data of each client 048 (T Dinh et al., 2020; Tan et al., 2022; Li et al., 2021; Sun et al., 2021). By doing so, these meth-049 ods mitigate the effects of data heterogeneity and class imbalance and improve the performance of 050 models on individual client tasks. Among PFL methods, clustering techniques have been employed 051 to group clients with similar data distributions, proposing an effective solution to face client statistical heterogeneity (Ghosh et al., 2020; Sattler et al., 2020). By partitioning clients based on data 052 distributions, clustering-based approaches reduce model divergence within each group, allowing for more homogeneous and targeted model updates.

054 In this work, we address the challenges of data heterogeneity and class imbalance through a novel 055 PFL clustering-based approach. We propose FedGWC (Federated Gaussian Weighting), a method 056 to group clients with similar data distributions into clusters, enabling the training of personalized 057 federated models for each group. Using a Gaussian reward mechanism to form homogeneous clus-058 ters, our method builds clusters based on an interaction matrix, which encodes the compatibility of all the couples of clients according to their data distribution. This is achieved by simply communicating the empirical losses to the server at each communication round. Each cluster benefits from 060 a personalized federated model that better captures the shared data characteristics within the group, 061 offering a more robust solution to the limitations of training a single global model over heteroge-062 neous data. Training on more homogeneous clusters allows us to exploit the advantages of FL, such 063 as aggregating knowledge from a larger pool of data, while avoiding the drawbacks of statistical 064 heterogeneity, such as client drift. The core idea of our approach is that client data distributions can 065 be inferred by a proper transformation of individual empirical loss functions, motivated by the fact 066 that training federated algorithms on homogeneous clusters leads to better performance compared 067 to training across the entire federation, as shown in (Sattler et al., 2020).

068 We present a comprehensive mathematical framework for the proposed algorithm and conduct a rig-069 orous analytical examination to establish the convergence properties of the Gaussian weights. Unlike 070 the majority of existing work on FL, we do not rely on model updates to cluster clients. Instead, inspired by (Cho et al., 2022), we extract valuable information from individual losses, hypothesizing 071 that clients with similar data distributions exhibit similar loss landscapes. Furthermore, FedGWC 072 can be integrated with any robust FL aggregation algorithm to furtherly handle data heterogeneity. 073 Additionally, we introduce a new clustering metric tailored for evaluating cluster cohesion in the 074 presence of class imbalance. Through experiments on benchmark datasets, we demonstrate that our 075 approach outperforms existing personalization and clustering algorithms in terms of accuracy and 076 clustering quality. 077

Contributions.

- We propose an efficient FL framework that clusters clients based on class imbalance and data heterogeneity, leading to personalized models within clusters.
- We introduce a new clustering metric specifically designed to evaluate clusters in the presence of class imbalance.
- We provide a rigorous mathematical framework to motivate the algorithm, showing its convergence properties.
- We compare FedGWC with clustered FL baselines, empirically showing that FedGWC provides the best clusters and improves performance in the most heterogeneous scenarios.
- We empirically show how FedGWC can be used with any FL aggregation algorithm.
- We investigate the behavior of our algorithm in class and domain imbalance scenarios, showing how our algorithm successfully clusters the clients according to different class distributions or domains.
- 089 090 091

880

079

081

083

084

085

2 RELATED WORK

092 FL with Heterogeneous Data. Handling data heterogeneity, especially class imbalance, remains a critical challenge in FL. FedProx(Li et al., 2020) was one of the first attempts to address het-094 erogeneity by introducing a proximal term that constrains local model updates close to the global 095 model. FedMD (Li & Wang, 2019) focuses on heterogeneity in model architectures, allowing collaborative training between clients with different neural network structures using model distillation. 096 Methods such as SCAFFOLD (Karimireddy et al., 2020b) and Mime (Karimireddy et al., 2020a) 097 have also been proposed to reduce client drift by using control variates during the optimization pro-098 cess, which helps mitigate the effects of non-IID data. Furthermore, strategies such as biased client selection (Cho et al., 2022) based on ranking local losses of clients and normalization of updates 100 in FedNova (Wang et al., 2021) have been developed to specifically address class imbalance in 101 federated networks, leading to more equitable global model performance. 102

Personalization in FL. Personalization in FL has been extensively explored to tackle the challenge of data heterogeneity across clients (Tan et al., 2022). Various techniques aim to adapt models to each client's specific needs while maintaining a collaborative training framework. A notable method is FedPer (Arivazhagan et al., 2019), which personalizes the last layers of the model while sharing the lower layers across all clients. Another important contribution is the pFedMe algorithm (T Dinh et al., 2020), which uses a Moreau envelope to decouple local and global model updates,

allowing for client-specific customization without losing the benefits of collaborative learning. Similarly, Per-FedAvg (Fallah et al., 2020) leverages a model-agnostic meta-learning (MAML) approach (Finn et al., 2017), optimizing a global model while fine-tuning each client locally, ensuring that the global model can generalize across clients but is also personalized. The Ditto framework (Li et al., 2021) further advances personalization by ensuring fairness and robustness, particularly in non-IID settings, through individualized model training. More recently, FedALA (Zhang et al., 2023) has been introduced, offering adaptive local aggregation for enhanced personalization.

115 **Clustering in FL.** Clustering has proven to be an effective strategy in FL for handling client het-116 erogeneity and improving personalization (Huang et al., 2022; Duan et al., 2021; Briggs et al., 2020; 117 Caldarola et al., 2021; Ye et al., 2023). Clustered FL (Sattler et al., 2020) is one of the first meth-118 ods proposed to group clients with similar data distributions to train specialized models rather than 119 relying on a single global one. Nevertheless, from a practical perspective, this method exhibits pro-120 nounced sensitivity to hyper-parameter tuning, especially concerning the gradient norms threshold, which is intricately linked to the dataset. This sensitivity can result in significant issues of either 121 excessive under-splitting or over-splitting. Additionally, as client sampling is independent of the 122 clustering, there may be privacy concerns due to the potential for updating cluster models with the 123 gradient of a single client. An extension of this is the efficient framework for clustered FL proposed 124 by (Ghosh et al., 2020), which strikes a balance between model accuracy and communication effi-125 ciency. Multi-Center FL (Long et al., 2023) builds on this concept by dynamically adjusting client 126 clusters to achieve better personalization, however a-priori knowledge on the number of clusters 127 is needed. Similarly, IFCA (Ghosh et al., 2020) addresses client heterogeneity by predefining a 128 fixed number of clusters and alternately estimating the cluster identities of the users by optimiz-129 ing model parameters for the user clusters via gradient descent. However, it imposes a significant 130 computational burden, as the server communicates all cluster models to each client, which must 131 evaluate every model locally to select the best fit based on loss minimization. This approach not only increases communication overhead but also introduces inefficiencies, as each client must test 132 all models, making it less scalable in larger networks. 133

Compared to previous approaches, the key advantage of the proposed algorithm, FedGWC, lies in its ability to effectively identify clusters of clients with similar levels of heterogeneity and class distribution through simple transformations of individual empirical losses. This is achieved without imposing significant communication overhead or requiring additional computational resources. Additionally, FedGWC can be seamlessly integrated with any aggregation method, enhancing its robustness and performance when dealing with heterogeneous scenarios.

139 140

3 FEDGWC: MATHEMATICAL FRAMEWORK

In this section, we introduce the mathematical framework of our algorithm, FEDERATED GAUSSIAN
 WEIGHTING CLUSTERING (FedGWC), a recursive clustered FL algorithm designed to group clients
 with similar data distributions, where each cluster develops its own model.

144 First, in Section 3.1 we introduce the mathematical formulation and notations of the FL problem 145 we aim to solve. Then, in Sections 3.2 and 3.3, we present the recursive step of our algorithm or, in other words, how FedGWC decides if a cluster needs to be split further into smaller but more 146 homogeneous clusters, or if the current cluster is already homogeneous enough. Finally, Section 3.4 147 introduces the full algorithmic notation, including cluster indices, and presents the overall recursive 148 structure. To improve clarity, we initially omit cluster indices, focusing instead on the internal 149 mechanics of the recursive step. In particular, without loss of generality, we consider a cluster of 150 clients as the totality of the clients in the federation in Sections 3.2 and 3.3, as the final objective 151 here is to explain how FedGWC eventually partitions a set of clients into some smaller clusters of 152 clients. Then, we reintroduce the cluster indices in Section 3.4.

In addition, we propose a novel metric to evaluate the quality of clustered FL algorithms in Section
 3.5. This metric, which derives from the Wasserstein distance (Kantorovich, 1942), quantifies the cohesion of client groups based on their class distribution similarities.

157 3.1 PROBLEM FORMULATION

We adopt the general structure of most FL methods where, during each round of training, a subset of clients downloads the model from the server, trains it locally using its data, and sends back the updated models for future rounds. Let K be the total number of clients and T the total number of communication rounds. At each communication round $t \in [T]$, a subset \mathcal{P}_t of participating clients are selected for training. Let S be the total number of training iterations that are performed at each com-



Figure 1: The principle underlying Gaussian reward mechanisms is that if the client's loss lies within the 172 confidence region, the algorithm assigns a higher Gaussian reward. Conversely, if the loss is further from the 173 confidence region, it assigns a lower reward. Left: Empirical loss processes, dashed lines, over S = 8 local iterations for 10 sampled clients *i.e.* $l_k^{t,s}$ for $k \in \mathcal{P}_t$ and $s = 1, \ldots, S$, the average loss is represented in black, 174 *i.e.* $\hat{\mu}^{t,s}$ for s = 1, ..., S and the light blue region delineating the confidence within the standard deviation *i.e.* $\hat{\sigma}^{t,s}$ for s = 1, ..., S. Right: The same representation for the identical clients in the same round is provided 175 176 with violin plots instead of intervals. Blacked dashed line is the average process of the empirical loss, *i.e.* $\hat{\mu}^{t,s}$, 177 with the associated standard errors, *i.e.* $\hat{\mu}^{t,s} \pm \hat{\sigma}^{t,s}$ for $s = 1, \dots, S$. 178

munication round during the training. In general, FL aims to solve an optimization problem that can 179 be stated in the following form (McMahan et al., 2017): $\min_{\theta \in \Theta} \mathcal{L}(\theta) = \min_{\theta \in \Theta} \sum_{k=1}^{K} \frac{n_k}{n} \mathcal{L}_k(\theta)$ 181 where $\mathcal{L}_k(\cdot)$ is the loss function of client k with n_k training samples, and $n = \sum_k n_k^{\prime \prime}$. When 182 k is sampled during round t, its local parameters are updated at every iteration s with a stochas-183 tic optimizer. For instance, Stochastic Gradient Descent (SGD) has the following update rule: $\theta_k^{t,s+1} = \theta_k^{t,s} - \eta \nabla \mathcal{L}_k(\theta_k^{t,s})$, where $\theta_k^{t,s}$ is the vector of the model parameters, η the learning rate. The stochastic process identified by the optimization algorithm suggests that the evolution of the 184 185 186 empirical loss can be modeled using random variables and tools from probability theory. 187

In the following sections, we will use capital letters (e.g., X) to denote random variables and lower-188 case letters (e.g., x) to represent their specific observations. 189

3.2 GAUSSIAN REWARDS WEIGHTS 190

191 To assess how closely each client's local data aligns with the global distribution, we introduce the 192 Gaussian Weights γ_k , statistical estimators that capture the closeness of each clients' class distribu-193 tion to the main distribution of the cluster. A weight near zero suggests that the client's distribution is far from the main distribution, meaning that it should probably belong to a separate cluster. We 194 pictorially represent the idea of the Gaussian rewards in Figure 1. 195

The core idea of FedGWC is to group clients based on the similarity of their empirical losses, which 196 are used to compute the *rewards*, continuous random variables in (0, 1). A high reward indicates that 197 a client's loss is close to the cluster's mean loss, while a lower reward reflects greater divergence. Gaussian weights estimate the expected value of these rewards, quantifying the closeness between 199 each client's distribution and the central distribution. 200

At each training round t and local training iteration s, we define the loss random variable $L_k^{t,s}$, of which we observe its samples $l_k^{t,s} = \mathcal{L}_k(\theta_k^{t,s})$. These values are naturally computed in the clients 201 202 during the local training. For each client k, it is possible to define a family of random *rewards* that 203 are stationary by construction, *i.e.*, their moments do not depend on the iteration and are obtained 204 with a Gaussian transformation of the loss 205

206 207

208

 $R_k^{t,s} = \exp\left(-\frac{(L_k^{t,s} - \hat{\mu}^{t,s})^2}{2(\hat{\sigma}^{t,s})^2}\right),\,$ (1)

209 210

where $\hat{\mu}^{t,s} = 1/|\mathcal{P}_t| \sum_{k \in \mathcal{P}_t} L_k^{t,s}$ is the sample mean across clients, and $(\hat{\sigma}^{t,s})^2 = 1/(|\mathcal{P}_t| - 1) \sum_{k \in \mathcal{P}_t} (L_k^{t,s} - \hat{\mu}^{t,s})^2$ is the sample variance. The rewards $R_k^{t,s} \to 1$ as the distance between $L_k^{t,s}$ and the average $\hat{\mu}^{t,s}$ decreases, while $R_k^{t,s} \to 0$ as the distance increases. Therefore, the ex-211 212 pected reward $\mathbb{E}[R_k^{t,s}]$ for out-of-distribution clients is lower than that of the in-distribution clients. 213 In Eq. 1, during any local iteration s, a Gaussian kernel with mean loss and sample variance as-214 sesses a client's proximity to the confidence interval's center, indicating their probability of sharing 215 the same learning process distribution.

To estimate the expected reward $\mathbb{E}[R_k^{t,s}]$, we introduce the random variable Ω_k^t , which is the average of the rewards across the S local iterations, *i.e.*, $\Omega_k^t = 1/S \sum_{s \in [S]} R_k^{t,s}$. We introduce Ω_k^t to 216 217 218 reduce noise in the estimation caused by stochastic fluctuations in the loss. Instead of using a single 219 sample, such as the last value of the loss as done in Cho et al. (2022), we opted for averaging across 220 iterations to provide a more stable estimate. Due to the linearity of the expectation operator, the 221 expected reward $\mathbb{E}[R_k^{t,s}]$ for the k-th client at round t, local iteration s equals the expected Gaussian 222 reward $\mathbb{E}[\Omega_k^t]$ that, to simplify the notation, we denote by μ_k . μ_k is the theoretical value that we 223 aim to estimate by designing our Gaussian weights Γ_k^t appropriately, as it encodes the ideal reward 224 to quantify the closeness of the distribution of each client k to the main distribution. Note that the 225 process is stationary by construction. Therefore, it does not depend on t but differs between clients, as it reaches a higher value for in-distribution clients and a lower for out-of-distribution clients. 226

Practically, we calculate the observed Gaussian weights γ_k^t . We initialize γ_k^t to zero, as we do not want to bias the estimation of the expectation of the reward. During each round t, a group of participating clients \mathcal{P}_t is sampled. Each client $k \in \mathcal{P}_t$ updates its model parameters θ_k^{t+1} , stores the observed loss process $l_k^t = (l_k^{t,1}, \ldots, l_k^{t,S}) \in \mathbb{R}^S$ and communicates it to the server, along with the update model parameters θ_k^{t+1} . Then, the server computes the observed reward process $t = (t_k^{t,1}, \ldots, t_k^{t,S}) \in \mathbb{R}^S$ and $t = (l_k^{t,1}, \ldots, l_k^{t,S}) \in \mathbb{R}^S$. 227 228 229 230 231 $r_k^t = (r_k^{t,1}, \ldots, r_k^{t,S}) \in [0,1]^S$ according to Eq. 1, and $\omega_k^t = 1/S \sum_{s \in [S]} r_k^{t,s}$, *i.e.*, the realization of the random variable Ω_k^t , averaging the entries of r_k^t . Finally the server can update the weight γ_k^t for 232 233 234 each selected client $k \in \mathcal{P}_t$ as 235

$$k^{t+1} = (1 - \alpha_t)\gamma_k^t + \alpha_t \omega_k^t \tag{2}$$

236 for a sequence of update coefficients $\{\alpha_t\}_{t=0}^{\infty}$ such that $0 < \alpha_t < 1 \ \forall t$. The weight definition in Eq.2 237 is closely related to the Robbins-Monro stochastic approximation method (Robbins & Monro, 1951). 238 If a client is not participating in the training, its weight is not updated. FedGWC mitigates biases in 239 the estimation of rewards by employing two mechanisms: (1) uniform random sampling method for 240 clients, with a dynamic adjustment process to prioritize clients that are infrequently sampled, thus 241 ensuring equitable participation across time periods; and (2) when a client is not sampled in a round, its weight and contribution to the reward estimate remain unchanged. 242

To rigorously motivate the construction of our algorithm and the reliability of the weights, we intro-243 duce the following theoretical results. Theorems 3.1 and 3.2 demonstrate that the weights converge 244 to a finite value and, more importantly, that this limit serves as an unbiased estimator of the theoret-245 ical reward μ_k . The first theorem provides a strong convergence result, showing that, with suitable 246 choices of the sequence $\{\alpha_t\}_t$, the expectation of the Gaussian weights Γ_k^t converges to μ_k in L^2 247 and almost surely. In addition, Theorem 3.2 extends this to the case of constant α_t , proving that the 248 weights still converge and remain unbiased estimators of the rewards as $t \to \infty$.

249 **Theorem 3.1.** Let $\{\alpha_t\}_{t=1}^{\infty}$ be a sequence of positive real values, and $\{\Gamma_k^t\}_{t=1}^{\infty}$ the sequence of Gaussian weights. If $\{\alpha_t\}_{t=1}^{\infty} \in l^2(\mathbb{N})/l^1(\mathbb{N})$, then Γ_k^t converges in L^2 . Furthermore, for $t \to \infty$, 250 251

$$\Gamma_k^t \longrightarrow \mu_k \ a.s.$$
 (3)

Theorem 3.2. Let $\alpha \in (0,1)$ be a fixed constant, then in the limit $t \to \infty$, the expectation of the weights converges to the individual theoretical reward μ_k , for each client $k = 1, \ldots, K$, i.e.,

$$\mathbb{E}[\Gamma_k^t] \longrightarrow \mu_k \ t \to \infty \,. \tag{4}$$

257 Finally, Proposition 3.1 shows that Gaussian weights reduce the variance of the estimate, thus de-258 creasing the error and enabling the construction of a confidence interval for μ_k .

259 **Proposition 3.1.** The variance of the weights Γ_k^t is smaller than the variance σ_k^2 of the theoretical 260 rewards $R_k^{t,s}$.

261 Complete proofs of Theorems 3.1, 3.2 and Proposition 3.1 are detailed in Appendix A. 262

263 3.3 INTERACTION MATRIX AND CLUSTERING

252 253

254

255

256

264 Interaction Matrix. In the previous section, we defined the Gaussian weights as a measure of 265 proximity between each client's data distribution and the main distribution of the cluster. Gaussian 266 weights are scalar quantities that offer an absolute measure of the alignment between a client's data distribution and the global distribution. Although these weights indicate the conformity of 267 each client's distribution individually, they do not consider the interrelations among the distributions 268 of different clients. Therefore, we propose to encode these interactions in an interaction matrix 269 $P^t \in \mathbb{R}^{K \times K}$ whose element P_{kj}^t estimates the similarity between the k-th and the j-th client data

distribution. The interaction matrix is initialized to the null matrix, *i.e.* $P_{kj}^0 = 0$ for every couple $k, j \in [K]$.

Specifically, we define the update rule for the matrix P^t as follows:

$$P_{kj}^{t+1} = \begin{cases} (1 - \alpha_t) P_{kj}^t + \alpha_t \omega_k^t &, \quad (k, j) \in \mathcal{P}_t \times \mathcal{P}_t \\ P_{kj}^t &, \quad (k, j) \notin \mathcal{P}_t \times \mathcal{P}_t \end{cases}$$
(5)

where $\{\alpha_t\}_t$ is the same sequence used to update the weights, and \mathcal{P}_t is the subset of clients sampled in round t.

Intuitively, in the long run, since ω_k^t measures the proximity of the loss process of client k to the average loss process of clients in \mathcal{P}_t at round t, we are estimating the *expected perception* of client k by client j with P_{kj}^t , *i.e.* a larger value indicates a higher degree of similarity between the loss profiles, whereas smaller values indicate a lower degree of similarity. For example, if P_{kj}^t is close to 1, it suggests that on average, whenever k and j have been simultaneously sampled prior to round t, ω_k^t was high, meaning that the two clients are well-represented within the same distribution.

Interestingly, it is possible to show that the diagonal values of the interaction matrix are exactly the Gaussian weights computed as in Eq. 2. Indeed, by looking at Eq. 5, if we take k = j, we obtain the same relation in Eq. 2; namely, for any k and any t, the equality $P_{kk}^t = \gamma_k^t$ holds. Moreover, the entries of the interaction matrix P^t are bounded, as shown in Proposition A.2. Furthermore, as a direct consequence of Theorem 3.2, there exists a matrix $P \in \mathbb{R}^{K \times K}$, such that, in the limit for $t \to \infty$, $\mathbb{E}[P_{kj}^t] \to P_{kj}$ entry-wise.

To effectively extract the information embedded in P, we introduce the concept of *unbiased perception vectors* (UPV). For any pair of clients $k, j \in [K]$, the UPV $v_k^j \in \mathbb{R}^{K-2}$ represents the k-th row of P, excluding the k-th and j-th entries. Recalling the construction of P^t , where each row indicates how a client is perceived to share the same distribution as other clients in the federation, the UPV v_k^j captures the collective perception of client k by all other clients, excluding both itself and client j. This exclusion is why we refer to v_k^j as *unbiased*.

The UPVs encode information about the relationships between clients, which can be exploited for clustering. However, the UPVs cannot be directly used as their entries are only aligned when considered in pairs. Instead, we construct the *affinity matrix* W by transforming the information encoded by the UPVs through an RBF kernel, as this choice allows to effectively model the affinity between clients: two clients are considered *affine* if similarly perceived by others. This relation is encoded by the entries of $W \in \mathbb{R}^{K \times K}$, which we define as:

$$W_{kj} = \mathcal{K}(v_k^j, v_j^k) = \exp\left(-\beta \left\|v_k^j - v_j^k\right\|^2\right).$$
(6)

The spread of the RBF kernel is controlled by a single hyper-parameter β : changes in this value provide different clustering outcomes, as shown in the sensitivity analysis in Appendix H.

Clustering. The affinity matrix W, designed to be symmetric, highlights features that capture 307 similarities between clients' distributions. Clustering is performed by the server using the rows of 308 W as feature vectors, as they contain the relevant information. We apply the spectral clustering algorithm (Ng et al., 2001) to W due to its effectiveness in detecting non-convex relationships em-310 bedded within the client affinities. Symmetrizing the interaction matrix P into the affinity matrix 311 W is fundamental for spectral clustering as it refines inter-client relationship representation. It mod-312 els interactions, reducing biases, and emphasizing reliable similarities. This improves robustness 313 to noise, allowing spectral clustering to effectively detect the distributional structure underlying the 314 clients' network (Von Luxburg, 2007). During the iterative training process, the server determines 315 whether to perform clustering by checking the convergence of the matrix P^t . Convergence is nu-316 merically verified when the mean squared error (MSE) between consecutive updates is less than a small threshold $\epsilon > 0$. To reduce the computational cost of computing the MSE at each round, we 317 employ a running average update. Specifically, the MSE is initialized to 1 in order to ensure stability 318 and avoid erratic updates in early iterations, and it is updated as follows: 319

$$MSE^{t+1} = (1 - \alpha_t)MSE^t + \alpha_t m^t \quad \text{with} \quad m^t = \frac{1}{|\mathcal{P}_t|} \sum_{k,j \in \mathcal{P}_t} |P_{kj}^t - P_{kj}^{t+1}|^2.$$
(7)

321 322

320

302

303

274 275 276

Algorithm 1 summarizes the clustering procedure. If the MSE is below ϵ , the server computes the matrix W^t and performs spectral clustering over W^t with a number of clusters $n \in \{2, ..., n_{max}\}$.

For each clustering outcome, the Davies-Bouldin (DB) score (Davies & Bouldin, 1979) is computed: DB larger than one means that clusters are not well separated, while if it is smaller than one, the clusters are well separated.

We denote by n_{cl} the optimal number of clusters detected by FedGWC. If $\min_{n=2,...,n_{max}} DB_n > 1$, 327 328 we do not split the current cluster. Hence, the optimal number of clusters is n_{cl} is one. In the other case, the optimal number of clusters is $n_{cl} \in \arg \min_{n=2,\dots,n_{max}} DB_n$. This requirement ensures proper control over the over-splitting phenomenon, a common issue in hierarchical clustering al-330 gorithms in FL. Over-splitting can undermine key principles of FL by creating degenerate clusters 331 with very few clients. Moreover, since each cluster requires a server for aggregation, two options 332 arise: either one client from each cluster acts as the server, or the central server manages the clusters 333 separately. In the first case, this assumption is unrealistic in cross-device settings, as clients typi-334 cally have limited resources. In the second case, the server would face an excessive workload, along 335 with a significant communication overhead from managing multiple models and performing sepa-336 rate aggregations. For these reasons, over-splitting is particularly problematic in cross-device FL. 337 Finally, on each cluster $\mathcal{C}^{(1)}, \ldots, \mathcal{C}^{(n_{cl})}$ an FL aggregation algorithm is trained separately, resulting 338 in models $\theta_{(1)}, \ldots, \theta_{(n_{cl})}$ personalized for each cluster. 339

340 3.4 FEDGWC ALGORITHM

In the previous sections, we have detailed the recursive step within the individual clusters. In this section, we present the full FedGWC recursive procedure (Algorithm 2), introducing the complete notation with indices for the distinct clusters. We denote the clustering index as n, and the total number of clusters N_{cl} .

The interaction matrix $P_{(1)}^0$ is initialized to the null matrix $0_{K \times K}$, and the *total number of clusters* 345 346 N_{cl}^0 , as no clusters have been formed yet, and $MSE_{(1)}^0$ are initialized to 1, in order to ensure stability 347 in early updates, allowing a gradual decrease. At each communication round t, and for cluster $C^{(n)}$, 348 where $n = 1, ..., N_{cl}^t$, the cluster server independently samples the participating clients $\mathcal{P}_t^{(n)} \subseteq$ 349 $\mathcal{C}^{(n)}$. Each client $k \in \mathcal{P}_t^{(n)}$ receives the current cluster model $\theta_{(n)}^t$. After performing local updates, 350 each client sends its updated model θ_k^{t+1} and empirical loss l_k^t back to the cluster server. The server 351 aggregates these updates to form the new cluster model $\theta_{(n)}^{t+1}$, computes the Gaussian rewards ω_k^t for the sampled clients, and updates the interaction matrix $P_{(n)}^{t+1}$ and $MSE_{(n)}^{t+1}$ according to Eqs. 5 and 352 353 354 7. If $MSE_{(n)}^{t+1}$ is lower than a threshold ϵ , the server of the cluster performs clustering to determine 355 whether to split cluster $\mathcal{C}^{(n)}$ into n_{cl} sub-groups, as outlined in Algorithm 1. The matrix $P_{(n)}^{t+1}$ is 356 then partitioned into sub-matrices by filtering its columns and rows according to the newly formed 357 clusters, with the MSE for these sub-matrices reinitialized to 1. This process results in a distinct 358 model $\theta_{(n)}$ for each cluster $\mathcal{C}^{(n)}$. When the final iteration T is reached we are left with N_{cl}^T clusters 359 with personalized models $\theta_{(n)}$ for $n = 1, \ldots, N_{cl}^T$. 360

Thanks to the Gaussian Weights, and the recursive spectral clustering on the affinity matrices, our algorithm, FedGWC, is able to detect groups of clients that display similar levels of heterogeneity. The clusters formed are more uniform, *i.e.* the class distributions within each group are more similar. These results are supported by experimental evaluations, discussed in Section 4.2.

365

3.5 A NEW METRIC TO EVALUATE CLUSTERING IN FL

In the previous section we observed that when clustering clients according to different heterogeneity
 levels, the outcome must be evaluated using a metric that assesses the cohesion of individual distributions. In this paragraph, we introduce a novel metric to evaluate the performance of clustering
 algorithms in FL. This metric, derived the Wasserstein distance (Kantorovich, 1942), quantifies the
 cohesion of client groups based on their class distribution similarities.

We propose a general method for adapting clustering metrics to account for class imbalances. This adjustment is particularly relevant when the underlying class distributions across clients are skewed.
The formal derivation and mathematical details of the proposed metric are provided in Appendix B. In this section, we provide a high-level overview of our new metric.

Consider a generic clustering metric *s*, e.g. Davies-Bouldin score (Davies & Bouldin, 1979) or the Silhouette score (Rousseeuw, 1987). Let *C* denote the total number of classes, and x_i^k the empirical frequency of the *i*-th class in the *k*-th client's local training set. Following theoretical reasonings, as shown in Appendix B, the empirical frequency vector for client *k*, denoted by $x_{(i)}^k$, is ordered according to the rank statistic of the class frequencies, *i.e.* $x_{(i)}^k \ge x_{(i+1)}^k$ for any i = 1, ..., C-1. The class-adjusted clustering metric \tilde{s} is defined as the standard clustering metric s computed on the ranked frequency vectors $x_{(i)}^k$. Specifically, the distance between two clients j and k results in

$$\frac{1}{C} \left(\sum_{i=1}^{C} \left| x_{(i)}^{k} - x_{(i)}^{j} \right|^{2} \right)^{1/2} .$$
(8)

This modification ensures that the clustering evaluation is sensitive to the distributional characteristics of the class imbalance. As we show in Appendix B, this adjustment is mathematically equivalent to assessing the dispersion between the empirical class probability distributions of different clients using the Wasserstein distance, also known as the Kantorovich-Rubenstein metric (Kantorovich, 1942). This equivalence highlights the theoretical soundness of using ranked class frequencies to better capture the variation in class distributions when evaluating clustering outcomes in FL.

³⁹¹ 4 EXPERIMENTS

382

384

In this section, we present the experimental results on widely used FL benchmark datasets (Caldas et al., 2018), comparing the performance of FedGWC with other baselines from the literature, including standard FL algorithms, personalized FL approaches, and clustering methods. The details on the dataset we use and the implementation choices are shown in Appendix G.

396 In Section 4.1, we first evaluate our method, FedGWC, against various clustering algorithms, in-397 cluding CFL (Sattler et al., 2020), FeSEM (Long et al., 2023), and IFCA (Ghosh et al., 2020), 398 to demonstrate that FedGWC yields superior and more consistent clustering outcomes in heteroge-399 neous scenarios than the baselines. Then, we investigate the benefits of integrating FedGWC with established FL baselines. Specifically, we test our approach with FedAvg (McMahan et al., 2017), 400 FedAvgM (Asad et al., 2020) and FedProx (Li et al., 2020), showing how our approach is orthog-401 onal to conventional FL aggregation methods. Finally, we compare FedGWC with FeSEM when 402 paired with personalized FL algorithms such as pFedMe (T Dinh et al., 2020) and Per-FedAvq 403 (Fallah et al., 2020), showing that our proposed solution exhibits greater robustness than existing 404 clustering techniques and can be effectively combined with pure personalization FL techniques to 405 achieve improved performance over these baselines. Finally, in Section 4.2, we propose analyses on 406 class and domain imbalance, showing that our algorithm successfully detects clients belonging to 407 separate distributions. Further experiments are presented in Appendix J. 408

409 4.1 FEDGWC IN HETEROGENEOUS SETTINGS

410 FedGWC vs. clustering baselines. We assess the performance of FedGWC in comparison to sev-411 eral FL clustering baselines (IFCA, FeSEM, and CFL) utilizing FedAvg aggregation. For the al-412 gorithms that require knowing the number of clusters in advance, *i.e.* FeSEM and IFCA, we show 413 only the best result among 2, 3, 4, and 5 clusters. The complete tuning is shown in Appendix I. 414 While IFCA showcases competitive outcomes, it incurs significant communication overhead, as 415 each client is required to evaluate models from every cluster in each round, rendering this approach impractical in cross-device scenarios and positioning it as an upper bound in our analysis. Although 416 FeSEM is less costly, it is constrained by a predefined number of clusters, reducing flexibility. On 417 the other hand, CFL demands extensive hyperparameter tuning - in contrast, FedGWC requires only 418 one hyperparameter – and produces overly granular clustering, resulting in an unrealistic number of 419 models for cross-device scenarios, or no clusters at all. 420

- In Table 1, we present a comparative analysis of these algorithms with respect to balanced accuracy, adjusted silhouette score (AS), and adjusted Davies-Bouldin index (ADB), employing FedAvg as the aggregation strategy. Recall that higher the value of AS the better the clustering outcome, as, for ADB, a lower value suggests a better cohesion between clusters.
- Notably, both FedGWC and CFL autonomously determine the optimal number of clusters based 425 on data heterogeneity, thereby offering a more scalable solution for large-scale cross-device FL. In 426 contrast to CFL, FedGWC consistently produced a reasonable number of clusters, even when using 427 the optimal hyperparameters for CFL, which resulted in no splits, thereby achieving performance 428 equivalent to FedAvg. Interestingly, as illustrated in Figure 2, FedGWC exhibits a significant im-429 provement in accuracy on Cifar100 precisely at the rounds where clustering occurs. Furthermore, as shown in Table 1, FedGWC consistently outperformed other methods on both Cifar10 and Cifar100 430 when evaluated using adjusted clustering metrics, indicating its superior ability to partition clients 431 into homogeneous clusters.

Table 1: Clustered FL baselines: FedGWC consistently outperforms the baselines in the most heterogeneous settings (Cifar100 and Femnist), even surpassing the IFCA upper-bound.

		Clustering method	С	Acc	AS	ADB
	0	IFCA	2	$\textbf{78.6} \pm \textbf{2.3}$	0.0 ± 0.0	17.1 ± 8.1
	arl	FeSem	3	71.5 ± 1.3	-0.1 ± 0.0	52.7 ± 29.9
	Cif	FedGWC	3	76.2 ± 0.9 75.8 ± 1.1	/ 0.1 + 0.0	2.6 + 0.0
taset	ar100	IFCA FeSem	5 5	47.5 ± 3.5 53.4 ± 1.8	-0.8 ± 0.2 -0.3 ± 0.1	5.2 ± 5.1 38.4 ± 13.0
Da	Cif	FedGWC	4	41.0 ± 1.3 53.4 ± 0.4	0.1 ± 0.0	2.4 ± 0.4
	ït	IFCA	5	$\textbf{76.7} \pm \textbf{0.6}$	$\textbf{0.3} \pm \textbf{0.1}$	$\textbf{0.5} \pm \textbf{0.1}$
	nis	FeSem	2	75.6 ± 0.2	0.0 ± 0.0	25.6 ± 7.8
	en	CFL	1	76.0 ± 0.1	/	/
	ц	FedGWC	4	76.1 ± 0.1	-0.2 ± 0.1	18.0 ± 6.2

Figure 2: Balanced accuracy on Cifar100 for FedGWC using FedAvg aggregation compared to the clustered FL baselines. FedGWC detects two splits, and demonstrates significant improvements in accuracy when clustering is performed. FedGWC has a faster and more stable convergence with respect to baseline algorithms.



Table 2: Comparison of classical FL algorithms using FedGWC vs. the same algorithms without FedGWC. Our empirical results suggest that employing our clustering algorithm provides benefits in most heterogeneous settings, such as Cifar100 and Femnist, while it is unnecessary in simpler settings like Cifar10.

FI mothod	Cifa	r10	Cifar	·100	Fem	nist
r L methou	no clusters	FedGWC	no clusters	FedGWC	no clusters	FedGWC
FedAvg	76.2 ± 0.9	75.8 ± 1.1	41.6 ± 1.3	53.4 ± 0.4	76.0 ± 0.1	$\textbf{76.1} \pm \textbf{0.1}$
FedAvgM	$\textbf{78.6} \pm \textbf{1.3}$	77.8 ± 2.0	41.5 ± 0.5	50.5 ± 0.3	$\textbf{83.3} \pm \textbf{0.3}$	83.2 ± 0.4
FedProx	$\textbf{76.1} \pm \textbf{0.1}$	75.6 ± 0.8	41.8 ± 1.0	49.1 ± 1.0	75.9 ± 0.2	$\textbf{76.3} \pm \textbf{0.2}$

453 Training FL methods with FedGWC. In this paragraph, we empirically show how FedGWC can 454 be used on top of any FL aggregation algorithm. Although FedGWC does not improve the results 455 on Cifar10 because of the simplicity of the task, our method consistently improved the performance of FL algorithms for the more heterogeneous settings of Cifar100 and Femnist. We show these re-456 sults in Table 2. FedGWC consistently boosted balanced accuracies across all methods. Notably, 457 on Cifar100, the scenario that most resembles cross-device FL, FedGWC substantially improves the 458 performance, leveraging the increased heterogeneity to construct more homogeneous clusters. This 459 resulted in an average increase of over 10% in balanced accuracy on Cifar100. While the improve-460 ment on Femnist was less pronounced, it remained beneficial in most cases, highlighting FedGWC's 461 adaptability to less heterogeneous environments without introducing overhead. These experiments 462 underscore the orthogonal nature of FedGWC to FL aggregation strategies, demonstrating its ability 463 to enhance the performance of classical FL algorithms.

464 FedGWC improves performances of PFL algorithms. To evaluate the benefits of combining pure 465 personalization with cluster-wise personalization, we conducted a comparative analysis of FedGWC 466 with FeSEM and the PFL baselines, pFedMe and Per-FedAvg. Given the cross-device nature of 467 our algorithm, we focused on FeSEM as a representative clustering baseline. While PFL methods 468 can be considered upper bounds, as they perform a more fine-grained personalization than clus-469 tered methods at the cost of more client resources, our experiments demonstrate that incorporating 470 FedGWC with PFL methods can further enhance their performance when personalization is feasible. 471 However, personalization entails higher computational overhead and diverges from the fundamental philosophy of FedGWC, which emphasizes grouping clients based on distribution similarity rather 472 than optimizing individual models of the clients. As illustrated in Table 3, FedGWC consistently 473 outperforms FeSEM and surpasses the pure personalization performance bound in all benchmark 474 datasets when combined with pFedME and Per-FedAvg, achieving a notable 4.5% improvement 475 on Cifar100. 476

477 4.2 ANALYSIS ON THE CLUSTERING DECISIONS OF FEDGWC

478 FedGWC detects different client class distributions. Here, we investigate the underlying mech-479 anisms behind FedGWC 's clustering decisions in heterogeneous scenarios. Specifically, we explore 480 how the algorithm identifies and groups clients based on the non-IID nature of their data distribu-481 tions, represented by the Dirichlet concentration parameter α . For the Cifar-10 dataset, we propose 482 three distinct client partitions: (1) 90 clients with $\alpha = 0$ and 10 clients with $\alpha = 100$; (2) 90 clients with $\alpha = 0.05$ and 10 clients with $\alpha = 100$; and (3) 40 clients with $\alpha = 100$, 30 clients 483 with $\alpha = 0.05$, and 30 clients with $\alpha = 0$. Similarly, for the Cifar-100 dataset, we apply a similar 484 splitting approach, obtaining the following partitions: (1) 100 clients with $\alpha = 0$ and 10 clients 485 with $\alpha = 1000$; (2) 90 clients with $\alpha = 0.5$ and 10 clients with $\alpha = 1000$; and (3) 40 clients with

432

437

443

444

445

446

447

Table 3: Comparison of PFL algorithms with and without FedGWC or FeSEM clustering. Empirical results 487 indicate that FedGWC outperforms FeSEM and enhances the performance of pFedME and per-FedAvg when 488 integrated with PFL aggregations.

489										
100	FI mothod		Cifar10		Cifar100			Femnist		
490	r L methou	no clusters	FedGWC	FeSEM	no clusters	FedGWC	FeSEM	no clusters	FedGWC	FeSEM
491	FedAvg	76.2 ± 0.9	75.8 ± 1.1	71.5 ± 1.3	41.6 ± 1.3	53.4 ± 0.4	53.4 ± 1.8	76.0 ± 0.1	$\textbf{76.1} \pm \textbf{0.1}$	$75.4{\scriptstyle\pm}0.5$
492	pFedME	93.4 ± 0.3	93.6 ± 0.1	93.4 ± 0.2	89.0 ± 0.1	93.5 ± 0.1	93.3 ± 0.1	71.2 ± 0.2	$\textbf{72.0} \pm \textbf{0.2}$	$34.5 {\pm}~0.9$
493	per-FedAvg	93.1 ± 0.1	93.4 ± 0.1	93.3 ± 0.1	93.4 ± 0.1	93.5 ± 0.1	93.3 ± 0.1	63.6 ± 0.3	63.9 ± 0.2	$63.6 {\pm} 0.2$

494 Table 4: Clustering with three different splits on CIFAR-10 and CIFAR-100 datasets. FedGWC has su-495 perior clustering quality across different splits. 496 Clustering C Dataset Snlif ۸S ADR

97	Dataset	Split	method	С	AS	ADB
8			IFCA	2	/	/
9		(10, 0, 90)	FeSem FedGWC	3 3	-0.0 ± 0.1 0.1 \pm 0.0	12.0 ± 2.0 0.2 ± 0.0
	CIFAR-10		IFCA	2	/	/
	chritter to	(10, 90, 0)	FeSem FedGWC	3 3	-0.0 ± 0.0 0.2 ± 0.0	12.0 ± 2.0 0.6 \pm 0.0
2			IFCA	2	-0.2 ± 0.0	1.0 ± 0.0
		(40, 30, 30)	FeSem FedGWC	3 3	0.1 ± 0.1 0.6 ± 0.1	20.6 ± 7.1 1.0 \pm 0.4
			IFCA	5	$\textbf{-0.9} \pm 0.0$	1.8 ± 0.0
		(10, 0, 100)	FeSem FedGWC	5 5	-0.8 ± 0.2 0.1 ± 0.1	$\begin{array}{c} 2.6 \pm 0.6 \\ \textbf{0.2} \pm \textbf{0.2} \end{array}$
	CIEAD 100	-	IFCA	5	-0.0 ± 0.0	5.6 ± 1.5
,	CIFAR-100	(10, 90, 0)	FeSem	5	0.2 ± 0.1	12.0 ± 2.0
			FedGWC	5	0.4 ± 0.1	6.4 ± 2.0
5			IFCA	5	-0.2 ± 0.0	1.0 ± 0.0
9		(40, 30, 30)	FeSem FedGWC	5 3	-0.2 ± 0.0 0.4 ± 0.2	33.2 ± 0.0 0.9 ± 0.1

Table 5: Clustering performance of FedGWC on federations with clients from distinct domains, consisting of clean, noisy, and blurred images across CIFAR-10 and CIFAR-100 datasets. Performance is measured usg the Rand Index score (Rand, 1971). It is noted that a and Index value approaching 1 signifies a perfect corspondence between the clustering ground truth and e assigned labels. It is important to observe that a and Index of 1 represents the maximum value this dex can reach, and in four instances, FedGWC sucssfully distinguishes all visual domains.

Dataset	Domain configuration	С	Rand Index
CIFAR-10	(50, 0, 50) (50, 50, 0) (40, 30, 30)	2 2 4	$\begin{array}{c} 1.0 \pm 0.0 \\ 1.0 \pm 0.0 \\ 0.9 \pm 0.0 \end{array}$
CIFAR-100	(50, 0, 50) (50, 50, 0) (40, 30, 30)	2 2 4	$\begin{array}{c} 1.0 \pm 0.0 \\ 1.0 \pm 0.0 \\ 0.6 \pm 0.0 \end{array}$

510

486

511 $\alpha = 1000$, 30 clients with $\alpha = 0.05$, and 30 clients with $\alpha = 0$. Unlike FeSem and IFCA, FedGWC is able to effectively detect varying levels of heterogeneity, as demonstrated by the adjusted Silhuette 512 and Davies-Bouldin reported in Table 4, proving its ability to separate clients according to their data 513 distributions. 514

515 FedGWC detects different visual client domains. Here, we focus on scenarios with nearly 516 uniform class imbalance (high α values) but with different visual domains to investigate how 517 FedGWC forms clusters in such settings. We incorporated various artificial domains (non-perturbed, noisy, and blurred images) into CIFAR-10 and CIFAR-100 datasets under homogeneous conditions 518 $(\alpha = 100.00)$. Our results demonstrate that FedGWC effectively clustered clients according to these 519 distinct domains. Table 5 presents the Rand-Index scores, which assess clustering quality based on 520 known domain labels. The high Rand-Index scores, often approaching the upper bound of 1, indicate 521 that FedGWC successfully separated clients into distinct clusters corresponding to their respective 522 domains. Figure 8 in the appendix visualizes the interaction matrix P, affinity matrix W, and the 523 scatter plot of the clustering with respect to the spectral bi-dimensional embedding on CIFAR-10 for 524 the (40, 30, 30) configuration. This analysis suggests that FedGWC may be applicable for detecting 525 malicious clients in FL, pinpointing a potential direction for future research. 526

5 CONCLUSIONS 527

We propose FedGWC, an efficient clustering algorithm for heterogeneous FL settings addressing the 528 challenge of non-IID data and class imbalance. Unlike existing clustered FL methods, FedGWC 529 groups clients by data distributions with flexibility and robustness, simply using the information 530 encoded by the individual empirical loss. FedGWC successfully detects homogeneous clusters, as 531 proved by our proposed novel Wasserstein Adjusted Metric. FedGWC detects splits by removing 532 out-of-distribution clients, thus simplifying the learning task within clusters without increasing com-533 munication overhead or computational cost. Empirical evaluations show that separately training 534 classical FL algorithms on the homogeneous clusters detected by FedGWC consistently improves the performance. Additionally, FedGWC excels over other clustering techniques in grouping clients 536 uniformly with respect to class imbalance and heterogeneity levels, which is crucial to mitigate the 537 effect of non-IIDness FL. Finally, clustering on different class unbalanced and domain unbalanced scenarios, which are correctly detected by FedGWC (see Section 4.2), suggests that FedGWC can 538 also be applied to anomaly client detection and to enhance robustness against malicious attacks in future research.

540	6 Reproducibility Statement
5/12	The reproducibility of the results and the theoretical contributions of this work has been a paramount
543	concern throughout the entire development of this project and while drafting this manuscript. In this
544	section, we provide details and a concise guide to reproduce our results and verify our contributions.
545	• Code Availability: All the code used in our experiments has been included in the supplementary
546	material of the submission. Additionally, we will release online a well-documented and structured
547	final version of the code to allow for easy reproduction of the experiments detailed in the paper.
548	• Datasets and data split: The datasets used in our experiments are publicly available and can
549	be downloaded online. Detailed instructions on accessing and preprocessing the datasets will be
550	provided, along with the final code release. The data splits can be generated directly through
551	scenarios
552	• Theoretical Results • We provide complete proofs for all theorems and propositions presented in
553	the paper in Appendices A and B.
554	We are confident that with these resources, all experimental and theoretical results can be reproduced
555	by the community.
556	
557	References
558	Manoi Ghuhan Ariyazhagan, Vinay Aggarwal, Aaditya Kumar Singh, and Sunay Choudhary. Fed
559	erated learning with personalization lavers. arXiv preprint arXiv:1912.00818. 2019
560	
561	Muhammad Asad, Ahmed Moustafa, and Takayuki Ito. Fedopt: Towards communication efficiency
562	and privacy preservation in federated learning. Applied Sciences, 10(8):2864, 2020.
563	
564	Reith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McWanan, Sarvar Patel Daniel Pamage Agron Sagel and Karn Sath Practical secure aggregation for federated
565	learning on user-held data arXiv preprint arXiv:1611.04482.2016
500	
562	Keith Bonawitz, Fariborz Salehi, Jakub Konečný, Brendan McMahan, and Marco Gruteser. Feder-
560	ated learning with autotuned communication-efficient secure aggregation. In 2019 53rd Asilomar
570	Conference on Signals, Systems, and Computers, pp. 1222–1226. IEEE, 2019.
571	Christopher Briggs Zhong Fon and Datar Andros Federated learning with hierarchical clustering
572	of local undates to improve training on non-iid data. In 2020 international joint conference on
573	neural networks (IJCNN), pp. 1–9. IEEE, 2020.
574	
575	Debora Caldarola, Massimiliano Mancini, Fabio Galasso, Marco Ciccone, Emanuele Rodola, and
576	Barbara Caputo. Cluster-driven graph federated learning over multiple domains. In <i>Proceedings</i>
577	of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops,
578	pp. 2777-2730, June 2021.
579	Debora Caldarola, Barbara Caputo, and Marco Ciccone. Improving generalization in federated
580	learning by seeking flat minima. In European Conference on Computer Vision, pp. 654-672.
581	Springer, 2022.
582	
583	Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H Brendan McMa-
584	nan, virginia Sinui, and Ameet raiwaikai. Leal: A benchmark for rederated settings. arXiv preprint arXiv:1812.01097.2018
585	preprint arxiv.1812.01097; 2018.
586	Xiaoyu Cao, Minghong Fang, Jia Liu, and Neil Zhenqiang Gong. Fltrust: Byzantine-robust feder-
500	ated learning via trust bootstrapping. arXiv preprint arXiv:2012.13995, 2020.
580	
509	Yae Jee Cho, Jianyu Wang, and Gauri Joshi. Towards understanding biased client selection in fed-
591	10375 PMLR 2022
592	10 <i>3 3</i> , 1 1112X , 2022.
593	David L Davies and Donald W Bouldin. A cluster separation measure. <i>IEEE transactions on pattern analysis and machine intelligence</i> , (2):224–227, 1979.

611

612

618

594	Moming Duan, Duo Liu, Xinyuan Ji, Renping Liu, Liang Liang, Xianzhang Chen, and Yujuan
595	Tan. Fedgroup: Efficient federated learning via decomposed similarity-based clustering. In
596	2021 IEEE Intl Conf on Parallel & Distributed Processing with Applications, Big Data &
597	Cloud Computing, Sustainable Computing & Communications, Social Computing & Network-
598	ing (ISPA/BDCloud/SocialCom/SustainCom), pp. 228–237. IEEE, 2021.

- Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning with the-600 oretical guarantees: A model-agnostic meta-learning approach. Advances in neural information 601 processing systems, 33:3557-3568, 2020. 602
- 603 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation 604 of deep networks. In International conference on machine learning, pp. 1126–1135. PMLR, 2017.
- Avishek Ghosh, Jichan Chung, Dong Yin, and Kannan Ramchandran. An efficient framework for 606 clustered federated learning. Advances in Neural Information Processing Systems, 33:19586-607 19597, 2020. 608
- 609 Jenny Hamer, Mehryar Mohri, and Ananda Theertha Suresh. Fedboost: A communication-efficient 610 algorithm for federated learning. In International Conference on Machine Learning, pp. 3973-3983. PMLR, 2020.
- J Harold, G Kushner, and George Yin. Stochastic approximation and recursive algorithm and appli-613 cations. Application of Mathematics, 35(10), 1997. 614
- 615 Kevin Hsieh, Amar Phanishayee, Onur Mutlu, and Phillip Gibbons. The non-iid data quagmire of 616 decentralized machine learning. In International Conference on Machine Learning, pp. 4387-617 4398. PMLR, 2020.
- Guangjing Huang, Xu Chen, Tao Ouyang, Qian Ma, Lin Chen, and Junshan Zhang. Collaboration 619 in participant-centric federated learning: A game-theoretical perspective. IEEE Transactions on 620 Mobile Computing, 22(11):6311-6326, 2022. 621
- 622 Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin 623 Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Ad-624 vances and open problems in federated learning. Foundations and trends® in machine learning, 625 14(1-2):1-210, 2021.
- Leonid V Kantorovich. On the translocation of masses. In Dokl. Akad. Nauk. USSR (NS), volume 37, 627 pp. 199-201, 1942. 628
- 629 Sai Praneeth Karimireddy, Martin Jaggi, Satyen Kale, Mehryar Mohri, Sashank J Reddi, Sebas-630 tian U Stich, and Ananda Theertha Suresh. Mime: Mimicking centralized stochastic algorithms 631 in federated learning. arXiv preprint arXiv:2008.03606, 2020a. 632
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and 633 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In 634 International conference on machine learning, pp. 5132–5143. PMLR, 2020b. 635
- 636 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 637 2009.638
- 639 Yann LeCun. The mnist database of handwritten digits. http://yann. lecun. com/exdb/mnist/, 1998.
- 640 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to 641 document recognition. Proceedings of the IEEE, 86(11):2278-2324, 1998. 642
- 643 Daliang Li and Junpu Wang. Fedmd: Heterogenous federated learning via model distillation. arXiv 644 preprint arXiv:1910.03581, 2019. 645
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. 646 Federated optimization in heterogeneous networks. Proceedings of Machine learning and sys-647 tems, 2:429-450, 2020.

658

661

670

682

683

684 685

686

687

689

690

691

- 648 Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated 649 learning through personalization. In International conference on machine learning, pp. 6357-650 6368. PMLR, 2021.
- Guodong Long, Ming Xie, Tao Shen, Tianyi Zhou, Xianzhi Wang, and Jing Jiang. Multi-center 652 federated learning: clients clustering for better personalization. World Wide Web, 26(1):481–500, 653 2023. 654
- 655 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 656 Communication-efficient learning of deep networks from decentralized data. In Artificial intelli-657 gence and statistics, pp. 1273-1282. PMLR, 2017.
- H Brendan McMahan, FX Yu, P Richtarik, AT Suresh, D Bacon, et al. Federated learning: Strate-659 gies for improving communication efficiency. In Proceedings of the 29th Conference on Neural 660 Information Processing Systems (NIPS), Barcelona, Spain, pp. 5–10, 2016.
- 662 Matias Mendieta, Taojiannan Yang, Pu Wang, Minwoo Lee, Zhengming Ding, and Chen Chen. 663 Local learning matters: Rethinking data heterogeneity in federated learning. In Proceedings of 664 the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8397–8406, 2022.
- 665 Andrew Ng, Michael Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. 666 Advances in neural information processing systems, 14, 2001. 667
- 668 William M Rand. Objective criteria for the evaluation of clustering methods. Journal of the Ameri-669 can Statistical association, 66(336):846–850, 1971.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. The annals of mathemati-671 cal statistics, pp. 400-407, 1951. 672
- 673 Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analy-674 sis. Journal of computational and applied mathematics, 20:53–65, 1987. 675
- Felix Sattler, Klaus-Robert Müller, and Wojciech Samek. Clustered federated learning: Model-676 agnostic distributed multitask optimization under privacy constraints. IEEE transactions on neu-677 ral networks and learning systems, 32(8):3710-3722, 2020. 678
- 679 Benyuan Sun, Hongxing Huo, Yi Yang, and Bo Bai. Partialfed: Cross-domain personalized feder-680 ated learning via partial initialization. Advances in Neural Information Processing Systems, 34: 681 23309-23320, 2021.
 - Canh T Dinh, Nguyen Tran, and Josh Nguyen. Personalized federated learning with moreau envelopes. Advances in neural information processing systems, 33:21394–21405, 2020.
 - Alysa Ziying Tan, Han Yu, Lizhen Cui, and Qiang Yang. Towards personalized federated learning. *IEEE transactions on neural networks and learning systems*, 34(12):9587–9603, 2022.
- Ulrike Von Luxburg. A tutorial on spectral clustering. Statistics and computing, 17:395–416, 2007. 688
 - Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. A novel framework for the analysis and design of heterogeneous federated learning. *IEEE Transactions on Signal* Processing, 69:5234–5249, 2021.
- Rui Ye, Zhenyang Ni, Fangzhao Wu, Siheng Chen, and Yanfeng Wang. Personalized federated 693 learning with inferred collaboration graphs. In International Conference on Machine Learning, 694 pp. 39801-39817. PMLR, 2023.
- 696 Jianqing Zhang, Yang Hua, Hao Wang, Tao Song, Zhengui Xue, Ruhui Ma, and Haibing Guan. 697 Fedala: Adaptive local aggregation for personalized federated learning. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 37, pp. 11237–11244, 2023. 699
- Michael Zhang, Karan Sapra, Sanja Fidler, Serena Yeung, and Jose M Alvarez. Personalized fed-700 erated learning with first order model optimization. In International Conference on Learning Representations, 2021.

702 703	Zaixi Zhang, Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. Fldetector: Defending feder-
704	ated learning against model poisoning attacks via detecting malicious clients. In Proceedings of
705	the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pp. 2343–2353,
706	2022.
707	Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated
707	learning with non-iid data. arXiv preprint arXiv:1806.00582, 2018.
700	
709	
710	
710	
712	
713	
715	
716	
710	
710	
710	
713	
720	
721	
722	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	

756 A THEORETICAL RESULTS FOR FEDGWC

This section provides algorithms, in pseudo-code, to describe FedGWC (see Algorithms 1 and 2). Additionally, here we provide the proofs for the convergence results introduced in Section 3, specifically addressing the convergence (Theorems 3.1 and 3.2) and the formal derivation on the variance bound of the Gaussian weights (Proposition 3.1). In addition, we also present a sufficient condition, under which is guaranteed that the overall sampling rate of the training algorithm does not increase and remain unchanged during the training process (Theorem A.3).

Theorem A.1. Let $\{\alpha_t\}_{t=1}^{\infty}$ be a sequence of positive real values, and $\{\Gamma_k^t\}_{t=1}^{\infty}$ the sequence of Gaussian weights. If $\{\alpha_t\}_{t=1}^{\infty} \in l^2(\mathbb{N})/l^1(\mathbb{N})$, then Γ_k^t converges in L^2 . Furthermore, for $t \to \infty$,

$$\Gamma_k^t \longrightarrow \mu_k \ a.s.$$
 (9)

Proof. At each communication round, we compute the samples $r_i^{t,s}$ from $R_k^{t,s}$ via a Gaussian transformation of the observed loss in Eq. 1. Notice that, due to the linearity of the expectation operator, $\mathbb{E}[\Omega_k^t] = \mu_k$, that is the true, unknown, expected reward. The observed value for the random variable is given by $\omega_k^t = 1/S \sum_{s=1}^{S} r_k^{t,s}$, which is sampled from a distribution centered on μ_k . Each client's weight is updated according to

$$\gamma_k^{t+1} = (1 - \alpha_t)\gamma_t + \alpha_t \omega_k^t \,. \tag{10}$$

Since such an estimator follows a Robbins-Monro algorithm, it is proved to converge in L^2 . In addition, Γ_k^t converges to the expectation $\mathbb{E}[\Omega_k^t] = \mu_k$ with probability 1, provided that α_t satisfies $\sum_{t\geq 1} |\alpha_t| = \infty$, and $\sum_{t\geq 1} |\alpha_t|^2 < \infty$ (Harold et al., 1997).

Theorem A.2. Let $\alpha \in (0,1)$ be a fixed constant, then in the limit $t \to \infty$, the expectation of the weights converges to the individual theoretical reward μ_k , for each client k = 1, ..., K, i.e.,

$$\mathbb{E}[\Gamma_k^t] \longrightarrow \mu_k \ t \to \infty.$$
⁽¹¹⁾

Proof. Recall that $\gamma_k^{t+1} = (1 - \alpha)\gamma_k^t + \alpha \omega_k^t$, where ω_k^t are samples from Ω_k^t . If we substitute backward the value of γ_k^t we can write

$$\gamma_k^{t+1} = (1-\alpha)^2 \gamma_k^{t-1} + \alpha \omega_k^t + \alpha (1-\alpha) \omega_k^{t-1} \,. \tag{12}$$

786 By iterating up to the initialization term γ_k^0 we get the following formulation:

$$\gamma_k^{t+1} = (1-\alpha)^{t+1} \gamma_k^0 + \sum_{\tau=0}^t \alpha (1-\alpha)^\tau \omega_k^{t-\tau} .$$
(13)

Since ω_k^t are independent and identically distributed samples from Ω_k^t , with expected value μ_k , then the expectation of the weight at the *t*-th communication round would be

$$\mathbb{E}[\Gamma_k^t] = \mathbb{E}\left[(1-\alpha)^t \gamma_k^0 + \sum_{\tau=0}^t \alpha (1-\alpha)^\tau \Omega_k^{t-\tau-1} \right] , \qquad (14)$$

that, due to the linearity of expectation, becomes

$$\mathbb{E}[\Gamma_k^t] = (1-\alpha)^t \gamma_k^0 + \sum_{\tau=0}^t \alpha (1-\alpha)^\tau \mu_k \ .$$
(15)

801 If we compute the limit

$$\lim_{t \to \infty} \mathbb{E}[\Gamma_k^t] = \lim_{t \to \infty} (1 - \alpha)^t \gamma_k^0 + \sum_{\tau=0}^\infty \alpha (1 - \alpha)^\tau \mu_k , \qquad (16)$$

and since $\alpha \in (0, 1)$, the first term tends to zero, and also the geometric series converges. Therefore, the expectation of the weights converges to μ_k , namely

$$\lim_{t \to \infty} \mathbb{E}[\Gamma_k^t] = \mu_k \,. \tag{17}$$

L		
L		
L		
-	_	

Proof. From Eq.13, we can show that $\mathbb{V}ar(\Gamma_k^t)$ converges to a value that depends on α and the 814 number of local training iterations S. Indeed

$$\mathbb{V}ar(\Gamma_k^t) = \mathbb{V}ar\left((1-\alpha)^t \gamma_k^0 + \sum_{\tau=0}^t \alpha (1-\alpha)^\tau \Omega_k^{t-\tau-1}\right)$$

$$= \sum_{\tau=0}^t \alpha^2 (1-\alpha)^{2\tau} \mathbb{V}ar(\Omega_k^t) = \frac{1}{S} \sum_{\tau=0}^t \alpha^2 (1-\alpha)^{2\tau} \sigma_k^2$$
(18)

since $\Omega_k^t = 1/S \sum_{s=1}^S R_k^{t,s}$.

If we compute the limit, that exists finite due to the hypothesis $\alpha \in (0, 1)$, we get

$$\lim_{t \to \infty} \mathbb{V}ar(\Gamma_k^t) = \frac{\alpha^2 \sigma_k^2}{S} \sum_{\tau=0}^{\infty} (1-\alpha)^{2\tau} = \frac{\alpha}{2-\alpha} \frac{\sigma_k^2}{S} < \frac{\sigma_k^2}{S} < \sigma_k^2 .$$
(19)

We further demonstrate that the interaction matrix P^t identified by FedGWC is entry-wise bounded from above, as established in the following proposition.

Proposition A.2. The entries of the interaction matrix P^t are bounded from above, namely for any $t \ge 0$ there exists a positive finite constant $C_t > 0$ such that

$$P_{kj}^t \le C_t \ . \tag{20}$$

And furthermore

$$\lim_{t \to \infty} C_t = 1 .$$
⁽²¹⁾

Proof. Without loss of generality we assume that every client of the federation is sampled, and we assume that $\alpha_t = \alpha \in (0, 1)$ for any $t \ge 0$. We recall, from Eq.5, that for any couple of clients $k, j \in \mathcal{P}_t$ the entries of the interaction matrix are updated according to

$$P_{kj}^{t+1} = (1 - \alpha)P_{kj}^{t} + \alpha \omega_k^t \,. \tag{22}$$

If we iterate backward until P_{kj}^0 , we obtain the following update

$$P_{kj}^{t+1} = (1-\alpha)^{t+1} P_{kj}^0 + \sum_{\tau=0}^t \alpha (1-\alpha)^\tau \omega_k^{t-\tau} .$$
⁽²³⁾

We know that, by constructions, the Gaussian rewards $\omega_k^t < 1$ at any time t, therefore the following inequality holds

$$P_{kj}^{t} = (1-\alpha)^{t} P_{kj}^{0} + \sum_{\tau=0}^{t} \alpha (1-\alpha)^{\tau} \omega_{k}^{t-\tau-1} \le (1-\alpha)^{t} P_{kj}^{0} + \sum_{\tau=0}^{t} \alpha (1-\alpha)^{\tau} .$$
(24)

At any round t we can define the costant C_t , as

$$C_t := (1 - \alpha)^t P_{kj}^0 + \alpha \sum_{\tau=0}^t (1 - \alpha)^\tau = (1 - \alpha)^t P_{kj}^0 + 1 - (1 - \alpha)^{t+1} < \infty.$$
(25)

Moreover, since $\alpha \in (0, 1)$, by taking the limit we prove that

$$\lim_{t \to \infty} C_t = \lim_{t \to \infty} (1 - \alpha)^t P_{kj}^0 + 1 - (1 - \alpha)^{t+1} = 1.$$
(26)

	OTTIME I FEdGW_CIUSTEI
1:	Input: $P, n_{max}, \mathcal{K}(\cdot, \cdot)$
2:	Output: cluster labels $y_{n_{cl}}$, and number of clusters n_{cl}
3:	Extract UPVs v_k^j, v_j^k from P for any k, j
4:	$W_{kj} \leftarrow \mathcal{K}(v_k^j, v_j^k)$ for any k, j
5:	for $n = 2, \ldots, n_{max}$ do
6:	$y_n \leftarrow \text{Spectral_Clustering}(W, n)$
7: o.	$DB_n \leftarrow \text{Davies_Bouldin}(W, y_n)$ if min $DB > 1$ then
o. 9:	$\frac{n \min_n DD_n > 1}{n - l} \leftarrow 1$
10:	else
11:	$n_{cl} \leftarrow \arg\min_n DB_n$
12:	end if
13:	end for
Alg	orithm 2 FedGWC
1:	Input: $K, T, S, \alpha_t, \epsilon, \mathcal{P}_t , \mathcal{K}$
2:	Output: $\mathcal{C}^{(1)}, \ldots, \mathcal{C}^{(N_{cl})}$ and $\theta_{(1)}, \ldots, \theta_{(N_{cl})}$
3:	Initialize $N_{cl}^{0} \leftarrow 1$
4:	Initialize $P_{(1)}^0 \leftarrow 0_{K \times K}$
5:	Initialize $MSE_{(1)}^0 \leftarrow 1$
6:	for $t = 0,, T - 1$ do
7:	$\Delta N^t \leftarrow 0$ for each iterations it counts the number of new clusters that are detected
8:	for $n=1,\ldots,N_{cl}^t$ do
9:	Server samples $\mathcal{P}_t^{(n)} \in \mathcal{C}^{(n)}$ and sends the current cluster model $ heta_{(n)}^t$
10:	Each client $k \in \mathcal{P}_t^{(n)}$ locally updates θ_k^t and l_k^t , then sends them to the server
11:	$\omega_k^t \leftarrow \texttt{Gaussian_Rewards}(l_k^t, \mathcal{P}_t^{(n)}), \texttt{Eq. 1}$
12:	$\theta_{(n)}^{t+1} \leftarrow \texttt{FL}_{\texttt{Aggregator}}(\theta_k^t, \mathcal{P}_t^{(n)})$
13:	$P_n^{t+1} \leftarrow \texttt{Update_Matrix}(P_n^t, \omega_k^t, lpha_t, \mathcal{P}_t^{(n)})$, according to Eq. 5
14:	Update $MSE_{(n)}^{t+1}$, according to Eq. 7
15:	if $\mathrm{MSE}_n^{t+1} < \epsilon$ then
16:	Perform FedGW_Cluster $(P_{(n)}^{t+1}, n_{max}, \mathcal{K})$ on $\mathcal{C}^{(n)}$, providing n_{cl} sub-clusters
17:	Update the number of new clusters $\Delta N^t \leftarrow \Delta N^t + n_{cl} - 1$ u
18:	Cluster server splits $P_{(n)}^{t+1}$ filtering rows and columns according to the new clusters
19:	Re-initialize MSE for new clusters to 1
20:	end if
21:	end for Undets the total number of elusters N^{t+1} , $N^{t} \rightarrow N^{t}$
22: 23·	end for

Theorem A.3. (Sufficient Condition for Sample Rate Conservation) Consider K_{min} as the minimum number of clients permitted per cluster, i.e. the cardinality $|\mathcal{C}_n| \geq K_{min}$ for any given cluster $n = 1, \ldots, n_{cl}$, and $\rho \in (0, 1]$ to represent the initial sample rate. There exists a critical threshold $n^* > 0$ such that, if $K_{min} \geq n^*$ is met, the total sample size does not increase.

Proof. Let us denote by ρ_n the participation rate relative to the *n*-th cluster, *i.e.*

$$\rho_n = \max\left\{\rho, \frac{3}{|\mathcal{C}_n|}\right\} \tag{27}$$

because, in order to maintain privacy of the clients' information we need to sample at least three clients, therefore ρ^n is at least 3 over the number of clients belonging to the cluster. The total participation rate at the end of the clustering process is given by

$$\rho^{\text{global}} = \sum_{n=1}^{n_{cl}} \frac{K_n}{K} \tag{28}$$

where K_n denotes the number of clients sampled within the *n*-th cluster. If we focus on the term K_n , recalling Equation 27, we have that

$$K_n = \rho_n |\mathcal{C}_n| = \max\left\{\rho, \frac{3}{|\mathcal{C}_n|}\right\} \times |\mathcal{C}_n| = \max\{\rho |\mathcal{C}_n|, 3\} .$$
⁽²⁹⁾

If we write Equation 29, by the means of the positive part function, denoted by $(x)^+ = \max\{0, x\}$, we obtain that

$$K_n = 3 + \max\{0, \rho |\mathcal{C}_n| - 3\} = 3 + (\rho |\mathcal{C}_n| - 3)^+ .$$
(30)

Observe that we are looking for a threshold value for which $\rho^{\text{global}} = \rho$, *i.e.* the participation rate remains the same during the whole training process.

Let us observe that $K_n = \rho |\mathcal{C}_n| \iff \rho |\mathcal{C}_n| \ge 3 \iff |\mathcal{C}_n| \ge n^* = 3/\rho$. In fact, if we assume that $K_{min} \ge n^*$, then the following chain of equalities holds

$$\rho^{\text{global}} = \sum_{n=1}^{n_{cl}} \frac{K_n}{K} = \frac{1}{K} \sum_{n=1}^{n_{cl}} \rho |\mathcal{C}_n| = \frac{\rho}{K} \sum_{n=1}^{n_{cl}} |\mathcal{C}_n| = \frac{\rho K}{K} = \rho$$

thus proving that $K_{min} \ge n^*$ is a sufficient condition for not increasing the sampling rate during the training process.

В THEORETICAL CONSTRUCTION OF THE CLUSTERING METRIC

To address the lack of clustering evaluation metrics suited for FL with distributional heterogeneity and class imbalance, we introduced a theoretically grounded adjustment to standard metrics, derived from the Wasserstein distance, Kantorovich-Rubinstein metric (Kantorovich, 1942). This metric, integrated with popular scores like Silhouette and Davies-Bouldin, enables a modular framework for a posteriori evaluation, effectively comparing clustering outcomes across federated algorithms. In this paragraph, we show how the proposed clustering metric that accounts for class imbalance can be derived from a probabilistic interpretation of clustering.

Definition B.1. Let (M, d) be a metric space, and $p \in [1, \infty]$. The Wasserstein distance between two probability measures \mathbb{P} and \mathbb{Q} over M is defined as

$$W_p(\mathbb{P}, \mathbb{Q}) = \inf_{\gamma \in \Gamma(\mathbb{P}, \mathbb{Q})} \mathbb{E}_{(x, y) \sim \gamma} [d(x, y)^p]^{1/p}$$
(31)

where $\Gamma(\mathbb{P}, \mathbb{Q})$ is the set of all the possible couplings of \mathbb{P} and \mathbb{Q} (see Def. B.2).

Furthermore, we need to introduce the notion of coupling of two probability measures.

Definition B.2. Let (M, d) be a metric space, and \mathbb{P}, \mathbb{Q} two probability measures over M. A coupling γ of \mathbb{P} and \mathbb{Q} is a joint probability measure on $M \times M$ such that, for any measurable subset $A \subset M$,

$$\int_{A} \left(\int_{M} \gamma(dx, dy) \mathbb{Q}(dy) \right) \mathbb{P}(dx) = \mathbb{P}(A),$$

$$\int_{A} \left(\int_{M} \gamma(dx, dy) \mathbb{P}(dx) \right) \mathbb{Q}(dy) = \mathbb{Q}(A).$$
(32)

Let us recall that the empirical measure over M of a sample of observations $\{x_1, \dots, x_N\}$ is defined such that for any measurable set $A \subset M$

$$\mathbb{P}(A) = \frac{1}{N} \sum_{i=1}^{C} \delta_{x_i}(A)$$
(33)

where δ_{x_i} is the Dirac's measure concentrated on the data point x_i .

In particular, we aim to measure the goodness of a cluster by taking into account the distance be-tween the empirical frequencies between two clients' class distributions and use that to properly adjust the clustering metric. For the sake of simplicity, we assume that the distance d over M is the L^2 -norm. We obtain the following theoretical result to justify the rationale behind our proposed metric.

975 Theorem B.1. Let s be an arbitrary clustering score. Then, the class-imbalance adjusted score š
976 is exactly the metric s computed with the Wasserstein distance between the empirical measures over
977 each client's class distribution.

Proof. Let us consider two clients; each one has its own sample of observations $\{x_1, \ldots, x_C\}$ and $\{y_1, \ldots, y_C\}$ where the *i*-th position corresponds to the frequency of training points of class *i* for each client. We aim to compute the *p*-Wasserstein distance between the empirical measures \mathbb{P} and \mathbb{Q} of the two clients, in particular for any dx, dy > 0

$$\mathbb{P}(dx) = \frac{1}{N} \sum_{i=1}^{N} \delta_{x_i}(dx),$$

$$\mathbb{Q}(dy) = \frac{1}{N} \sum_{i=1}^{N} \delta_{y_i}(dy) .$$
(34)

In order to compute $W_p^p(\mathbb{P}, \mathbb{Q})$ we need to carefully investigate the set of all possible coupling measures $\Gamma(\mathbb{P}, \mathbb{Q})$. However, since either \mathbb{P} and \mathbb{Q} are concentrated over countable sets, it is possible to see that the only possible couplings satisfying Eq. 32 are the Dirac's measures over all the possible permutations of x_i and y_i . In particular, by fixing the ordering of x_i , according to the rank statistic $x_{(i)}$, the coupling set can be written as

$$\Gamma(\mathbb{P}, \mathbb{Q}) = \left\{ \frac{1}{C} \delta_{(x_{(i)}, y_{\pi(i)})} : \pi \in \mathcal{S} \right\}$$
(35)

where S is the set of all possible permutations of C elements. Therefore we could write Eq. 31 as follows

$$W_p^p = \min_{\pi \in \mathcal{S}} \int_{M \times M} |x - y|^p \frac{1}{N} \sum_{i=1}^C \delta_{(x_{(i)}, y_{\pi(i)})}(dx, dy)$$
(36)

since S is finite, the infimum is a minimum. By exploiting the definition of Dirac's distribution and the linearity of the Lebesgue integral, for any $\pi \in S$, we get

$$\int_{M \times M} |x - y|^p \frac{1}{C} \sum_{i=1}^C \delta_{(x_{(i)}, y_{\pi(i)})}(dx, dy) = \frac{1}{C} \sum_{i=1}^C \int_{M \times M} |x - y|^p \delta_{(x_{(i)}, y_{\pi(i)})}(dx, dy)$$

$$= \frac{1}{C} \sum_{i=1}^C |x_{(i)} - y_{\pi(i)}|^p .$$
(37)

1012 Therefore, finding the Wasserstein distance between \mathbb{P} and \mathbb{Q} boils down to a combinatorial opti-1013 mization problem, that is, finding the permutation $\pi \in S$ that solves

$$W_{p}^{p}(\mathbb{P},\mathbb{Q}) = \min_{\pi \in \mathcal{S}} \frac{1}{C} \sum_{i=1}^{C} |x_{(i)} - y_{\pi(i)}|^{p} .$$
(38)

The minimum is achieved when $\pi = \pi^*$ that is the permutation providing the ranking statistic, i.e. $\pi^*(y_i) = y_{(i)}$, since the smallest value of the sum is given for the smallest fluctuations. Thus we conclude that the *p*-Wasserstein distance between \mathbb{P} and \mathbb{Q} is given by

$$W_p(\mathbb{P}, \mathbb{Q}) = \left(\frac{1}{C} \sum_{i=1}^{C} |x_{(i)} - y_{(i)}|^p\right)^{1/p}$$
(39)

that is the pairwise distance computed between the class frequency vectors, sorted in order of magnitude, for each client, introduced in Section 3.5, where we chose p = 2.

1026 С PRIVACY OF FEDGWC 1027

1028 In the framework of FedGWC, clients are required to send only the empirical loss vectors $l_k^{r,s}$ to the server (Cho et al., 2022). While concerns might arise regarding the potential leakage of sensitive 1029 information from sharing this data, it is important to clarify that the server only needs to access 1030 aggregated statistics, working on aggregated data. This ensures that client-specific information re-1031 mains private. Privacy can be effectively preserved by implementing the Secure Aggregation proto-1032 col (Bonawitz et al., 2016), which guarantees that only the aggregated results are shared, preventing 1033 the exposure of any raw client data. 1034

1035

D COMMUNICATION AND COMPUTATIONAL OVERHEAD OF FEDGWC 1036

1037 FedGWC minimizes communication and computational overhead, aligning with the requirements of 1038 scalable FL systems (McMahan et al., 2016). On the client side, the computational cost remains 1039 unchanged compared to the chosen FL aggregation, e.g. FedA, as clients are only required to com-1040 municate their local models and a vector of empirical losses after each round. The size of this loss 1041 vector, denoted by S, corresponds to the number of local iterations (*i.e.* the product of local epochs and the number of batches) and is negligible w.r.t. the size of the model parameter space, $|\Theta|$. In our 1042 experimental setup, S = 8, ensuring that the additional communication overhead from transmitting 1043 loss values is negligible in comparison to the transmission of model weights. 1044

All clustering computations, including those based on interaction matrices and Gaussian weighting, 1045 are performed exclusively on the server. This design ensures that client devices are not burdened 1046 with additional computational complexity or memory demands. The interaction matrix P used in 1047 FedGWC is updated incrementally and involves sparse matrix operations, which significantly reduce 1048 both memory usage and computational costs. 1049

These characteristics make FedGWC particularly well-suited for cross-device scenarios involving 1050 large federations and numerous communication rounds. Moreover, by operating on scalar loss val-1051 ues rather than high-dimensional model parameters, the clustering process in FedGWC achieves 1052 computational efficiency while maintaining effective grouping of clients. The server-side process-1053 ing ensures that the method remains scalable, even as the number of clients and communication 1054 rounds increases. Consequently, FedGWC meets the fundamental objectives of FL by minimizing 1055 costs while preserving privacy and maintaining high performance.

1056

1064

1066

10

1057 Ε METRICS USED FOR EVALUATION 1058

E.1 SILHOUETTE SCORE 1059

Silhouette Score is a clustering metric that measures the consistency of points within clusters by 1061 comparing intra-cluster and nearest-cluster distances (Rousseeuw, 1987). Let us consider a metric space (M, d). For a set of points $\{x_1, \ldots, x_N\} \subset M$ and clustering labels $\mathcal{C}_1, \ldots, \mathcal{C}_{n_{cl}}$. The 1062 Silhouette score of a data point x_i belonging to a cluster C_i is defined as 1063

$$s_i = \frac{b_i - a_i}{\max\{a_i, b_i\}} \tag{40}$$

1067 where the values b_i and a_i represent the average intra-cluster distance and the minimal average 1068 outer-cluster distance, *i.e.* 1069

$$a_{i} = \frac{1}{|C_{i}| - 1} \sum_{x_{j} \in C_{i} \setminus \{x_{i}\}} d(x_{i}, x_{j})$$

$$b_{i} = \min_{j \neq i} \frac{1}{|C_{j}|} \sum_{x_{j} \in C_{j}} d(x_{i}, x_{j})$$

$$(41)$$

$$b_{i} = \min_{j \neq i} \frac{1}{|C_{j}|} \sum_{x_{j} \in C_{j}} d(x_{i}, x_{j})$$

The value of the Silhouette score ranges between -1 and +1, *i.e.* $s_i \in [-1, 1]$. In particular, a Silhouette score close to 1 indicates well-clustered data points, 0 denotes points near cluster boundaries, 1077 and -1 suggests misclassified points. In order to evaluate the overall performance of the clustering, 1078 a common choice, that is the one adopted in this paper, is to average the score value for each data 1079 point.

1080 E.2 DAVIES-BOULDIN SCORE

The Davies-Bouldin Score is a clustering metric that evaluates the quality of clustering by measuring the ratio of intra-cluster dispersion to inter-cluster separation (Davies & Bouldin, 1979). Let us consider a metric space (M, d), a set of points $\{x_1, \ldots, x_N\} \subset M$, and clustering labels $C_1, \ldots, C_{n_{cl}}$. The Davies-Bouldin score is defined as the average similarity measure R_{ij} between each cluster C_i and its most similar cluster C_j :

$$DB = \frac{1}{n_{cl}} \sum_{i=1}^{n_{cl}} \max_{j \neq i} R_{ij}$$
(42)

where R_{ij} is given by the ratio of intra-cluster distance S_i to inter-cluster distance D_{ij} , *i.e.*

$$R_{ij} = \frac{S_i + S_j}{D_{ij}} \tag{43}$$

1093 with intra-cluster distance S_i defined as

$$S_i = \frac{1}{|\mathcal{C}_i|} \sum_{x_k \in \mathcal{C}_i} d(x_k, c_i)$$
(44)

where c_i denotes the centroid of cluster C_i , and $D_{ij} = d(c_i, c_j)$ is the distance between centroids of clusters C_i and C_j . A lower Davies-Bouldin Index indicates better clustering, as it reflects wellseparated and compact clusters. Conversely, a higher DBI suggests that clusters are less distinct and more dispersed.

1102 E.2.1 RAND INDEX

1087

1088

1090 1091

1094 1095

1107

1108 1109

1103Rand Index is a clustering score that measures the outcome of a clustering algorithm with respect to
a ground truth clustering label (Rand, 1971). Let us denote by a the number of pairs that have been
grouped in the same clusters, while by b the number of pairs that have been grouped in different
clusters, then the Rand-Index is defined as

$$RI = \frac{a+b}{\binom{N}{2}} \tag{45}$$

where N denotes the number of data points. In our experiments we opted for the Rand Index score to
evaluate how the algorithm was able to separate clients in groups of the same level of heterogeneity
(which was known a priori and used as ground truth). A Rand Index ranges in [0, 1], and a value of
1 signifies a perfect agreement between the identified clusters and the ground truth.

¹¹¹⁴ **F** ESTIMATING THE DIRICHLET PARAMETER α

In this appendix, we provide a detailed explanation of the procedure used to estimate the Dirichlet parameter β for each split identified by our algorithm. This analysis is conducted a posteriori to statistically evaluate the capability of FedGWC in grouping clients with similar levels of data heterogeneity.

1119 Consider a fixed cluster of clients, denoted by C_i . For each client $k \in C_i$, let $z_k \in \mathbb{N}^C$ represent 1120 the vector of sample counts per class (with C as the total number of classes). Given z_k , we can 1121 compute the likelihood function, assuming z_k is drawn from a Dirichlet distribution parameterized 1122 by the skew parameter α_k , such that $\mathcal{L}(\alpha_k; z_k)$ represents the likelihood function for z_k . 1123 To estimate α_k we apply the maximum likelihood estimator by solving:

To estimate α_k , we employ the maximum likelihood estimator by solving:

$$\hat{\alpha}_k \in rg\max_{\alpha > 0} \mathcal{L}(\alpha; z_k)$$

1127 This estimation is achieved through stochastic optimization (e.g., ADAM or SGD). For each clus-1128 ter, we obtain an estimate $\hat{\alpha}_k$, which allows us to compare average values, standard errors, and 1129 confidence intervals of α across clusters detected by FedGWC, providing a quantitative measure of 1130 heterogeneity among clusters.

1131

1125 1126

1132

1134 G DATASETS AND IMPLEMENTATION DETAILS

To simulate a realistic FL environment with heterogeneous data distribution, we conduct experi-1136 ments on CIFAR-100 (Krizhevsky et al., 2009). As a comparison, we also run experiments on 1137 the simpler CIFAR-10 dataset Krizhevsky et al. (2009). CIFAR-10 and CIFAR-100 are distributed 1138 among K clients using a Dirichlet distribution (by default, we use $\alpha = 0.05$ for CIFAR-10 and 1139 $\alpha = 0.5$ for CIFAR-100) to create highly imbalanced and heterogeneous settings. By default, we 1140 use K = 100 clients with 500 training and 100 test images. The classification model is a CNN 1141 with two convolutional blocks and three dense layers. Additionally, we perform experiments on 1142 the Femnist dataset LeCun (1998), partitioned among 400 clients using a Dirichlet distribution with 1143 $\alpha = 0.01$. In these experiments, we employ LeNet5 as the classification model LeCun et al. (1998). 1144 Local training on each client uses SGD with a learning rate of 0.01, weight decay of $4 \cdot 10^{-4}$, and batch size 64. The number of local epochs is 1, resulting in 7 batch iterations for CIFAR-10 and 1145 CIFAR-100 and 8 batch iterations for Femnist. The number of communication rounds is set to 3,000 1146 for Femnist, 10,000 for CIFAR-10 and 20,000 for CIFAR-100, with a 10% client participation rate 1147 per cluster. For FedGWC we employ constant value $\alpha_t = \alpha$ equal to the participation rate, *i.e.* 10%. 1148 As the performance metric, we use the balanced accuracy, *i.e.* the average accuracy of each cluster 1149 model evaluated on the test sets associated to each client in the same cluster. 1150

1151 H SENSITIVE ANALYSIS BETA VALUE RBF KERNEL

1153 This section provides a sensitivity analysis for the β hyper-parameter of the RBF kernel adopted for 1154 FedGWC. The results of this tuning are shown in Table 6.

1155

Table 6: A sensitivity analysis on the RBF kernel hyper-parameter β is conducted. We present the balanced accuracy for FedGWC on the Cifar10, Cifar100, and Femnist datasets for $\beta \in \{0.1, 0.5, 1.0, 2.0, 4.0\}$. It is noteworthy that FedGWC demonstrates robustness to variations in this hyperparameter.

β	Cifar10	Cifar100	Femnist
0.1	74.2	49.9	76.0
0.5	74.9	53.4	76.0
1.0	75.1	49.5	76.0
2.0	75.6	50.9	75.6
4.0	72.6	52.6	76.1

1164 1165 1166

I EVALUATION OF IFCA AND CFL ALGORITHMS WITH DIFFERENT NUMBER OF CLUSTERS OF CLUSTERS

This section shows the tuning of the number of clusters for the IFCA and CFL algorithms, which cannot automatically detect this value. The results of this tuning are shown in Table 7.

1172

J FURTHER EXPERIMENTS

Figure 5 illustrates the clustering results corresponding to varying degrees of heterogeneity, as described in Section 4.2. As per FedGWC, the detection of clusters based on different levels of heterogeneity in the Cifar10 dataset is achieved. Specifically, an examination of the interaction matrix reveals a clear distinction between the two groups. In Figure 6, we show that in class-balanced scenarios with small heterogeneity, like Cifar10 with $\alpha = 100$, FedGWC successfully detects one single cluster. Indeed, in homogeneous scenarios such as this one, the model benefits from accessing more data from all the clients.

Figure 4 shows how the MSE converges to a small value as the rounds increase for a Cifar10 experiment.

As Figure 7 illustrates, FedGWC partitions the CIFAR-100 dataset into clients based on class distri-

butions. Each cluster's distribution is distinct and non-overlapping, demonstrating the algorithm's

efficacy in partitioning data with varying degrees of heterogeneity. In Figure 8, we report the domain detection on Cifer 100, where 40 clients have clean images, 30 have poicy images, and 30 have

main detection on Cifar100, where 40 clients have clean images, 30 have noisy images, and 30 have blurred images. Table 5 shows that FedGWC performs a good clustering, effectively separating the

1187 different domains.



Table 7: Performance of for baseline algorithms for clustering in FL FeSEM, and IFCA, w.r.t. the number of

Figure 3: Cluster evolution with respect to the recursive splits in FedGWC on Cifar100, projected on the 1238 spectral embedded bi-dimensional space. From left to right, top to bottom, we can see that FedGWC splits the 1239 client into cluster, until a certain level of intra-cluster homogeneity is reached 1240

1241





Figure 8: FedGWC in the presence of domain imbalance. Three domains on Cifar100: clean clients (unlabeled), noisy clients (+), and blurred clients (x). *Left*: is the interaction matrix P at convergence from which it is possible to see client relations. *Center*: The affinity matrix W computed with respect to the UPVs extracted from P, and on which FedGW_Clustering is performed. We can see that FedGWC clusters the clients according to the domain, as proved by results in Table 5.