

# 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 CLOVE: PERSONALIZED FEDERATED LEARNING THROUGH CLUSTERING OF LOSS VECTOR EMBED- DINGS

006  
007     **Anonymous authors**  
008     Paper under double-blind review

## 009     ABSTRACT

011     We propose *CLoVE* (Clustering of Loss Vector Embeddings), a novel algorithm  
012     for Clustered Federated Learning (CFL). In CFL, clients are naturally grouped  
013     into clusters based on their data distribution. However, identifying these clusters is  
014     challenging, as client assignments are unknown. *CLoVE* utilizes client embeddings  
015     derived from model losses on client data, and leverages the insight that clients in  
016     the same cluster share similar loss values, while those in different clusters exhibit  
017     distinct loss patterns. Based on these embeddings, *CLoVE* is able to iteratively  
018     identify and separate clients from different clusters and optimize cluster-specific  
019     models through federated aggregation. Key advantages of *CLoVE* over existing  
020     CFL algorithms are (1) its simplicity, (2) its applicability to both supervised and  
021     unsupervised settings, and (3) the fact that it eliminates the need for near-optimal  
022     model initialization, which makes it more robust and better suited for real-world ap-  
023     plications. We establish theoretical convergence bounds, showing that *CLoVE* can  
024     recover clusters accurately with high probability in a single round and converges  
025     exponentially fast to optimal models in a linear setting. Our comprehensive experi-  
026     ments comparing with a variety of both CFL and generic Personalized Federated  
027     Learning (PFL) algorithms on different types of datasets and an extensive array of  
028     non-IID settings demonstrate that *CLoVE* achieves highly accurate cluster recovery  
029     in just a few rounds of training, along with state-of-the-art model accuracy, across  
030     a variety of both supervised and unsupervised PFL tasks.

## 031     1 INTRODUCTION

032     Federated learning (FL) has emerged as a pivotal framework for training models on decentralized data,  
033     preserving privacy and reducing communication overhead (McMahan et al., 2017; Konečný et al.,  
034     2016; Bonawitz et al., 2019). However, a significant challenge in FL is posed by the heterogeneity  
035     of data distributions across clients (non-IID data) (Zhao et al., 2018), as this can adversely impact  
036     the accuracy and convergence of FL models. To address this issue, Personalized Federated Learning  
037     (PFL) (Tan et al., 2023; Hanzely & Richtárik, 2020) has developed as an active research area, focusing  
038     on tailoring trained models to each client’s local data, thereby improving model performance.

039     Our research focuses on a variant of PFL known as *Clustered Federated Learning (CFL)*. In CFL,  
040     (Ghosh et al., 2020; Werner et al., 2023; Sattler et al., 2021; Mansour et al., 2020a; Chung et al.,  
041     2022; Long et al., 2023; Duan et al., 2022; Vahidian et al., 2023), clients are naturally grouped into  
042      $K$  *true clusters* based on inherent similarities in their local data distributions. However, these cluster  
043     assignments are not known in advance. The objective of CFL is to develop distinct models tailored  
044     to each cluster, utilizing only the local data from clients within that cluster. To achieve this goal,  
045     identifying the cluster assignments for each client is crucial, as it enables the construction of accurate,  
046     cluster-specific models that capture the unique characteristics of each group.

047     Many existing approaches (Awasthi & Sheffet, 2012; Kumar & Kannan, 2010) that provide guarantees  
048     on cluster assignment recovery often rely on the assumption that data of different clusters are well-  
049     separated. However, this assumption is often violated in real-world datasets, as evidenced by the  
050     modest performance of k-means on MNIST clustering even in a centralized setting. Its low Adjusted  
051     Rand Index (ARI)<sup>1</sup> of only approximately 50% (Treder-Tschechlov, 2024) can be attributed to the

052  
053     <sup>1</sup>ARI is a measure of the level of agreement between two clusterings, calculated as the fraction of pairs of  
cluster assignments that agree between the clusterings, adjusted for chance.

054 high variance in the MNIST image data, resulting in many digits being closer to other clusters  
 055 than to their own digit's cluster. Furthermore, in the federated setting, clustering approaches like  
 056 *k-FED* (Dennis et al., 2021) can recover clustering in one shot, but they also rely on the inter-cluster  
 057 separation conditions of Awasthi & Sheffet (2012). However, as we demonstrate (App. D.1.1), even  
 058 *k-FED* achieves low ARI on clustering common mixed linear regression distributions (Yi et al., 2014),  
 059 as the cluster centers of those distributions can also be very close, highlighting the limitations of  
 060 these approaches.

061 An alternative way to recover cluster assignments is to  
 062 leverage the separation of cluster-specific models. For  
 063 instance, the Iterative Federated Clustering Algorithm  
 064 (*IFCA*) (Ghosh et al., 2020) assigns clients to clusters  
 065 based on the model that yields the lowest empirical loss.  
 066 However, *IFCA* has a significant limitation: it requires  
 067 careful initialization of cluster models, ensuring each  
 068 model is sufficiently close to its optimal counterpart.  
 069 Specifically, the distance between an initial model and  
 070 the cluster's optimal model must be strictly less than half  
 071 the minimum separation distance between any two optimal  
 072 models. In practice, achieving this distance requirement  
 073 can be challenging. For example, in our experiments with  
 074 the MNIST dataset, *IFCA* often failed to recover accurate  
 075 cluster assignments. As shown in Fig. 1, this issue  
 076 typically occurs when clients from two or more clusters  
 077 achieve the lowest loss on the same model, and thus get  
 078 assigned to the same cluster. This, in turn, causes the  
 079 algorithm to get stuck and fail to recover accurate clusters.

080 In this paper, we propose the Clustering of Loss Vector  
 081 Embeddings (*CLoVE*) PFL algorithm, a novel approach  
 082 to tackle both the challenges of constructing accurate  
 083 cluster-specific models and recovering the underlying clus-  
 084 ter assignments in a FL setting. Our solution employs  
 085 an iterative process that simultaneously builds multiple  
 086 personalized models and generates client embeddings, rep-  
 087 presented as vectors of losses achieved by these models  
 088 on the clients' local data. Specifically, we leverage the  
 089 clustering of client embeddings to refine the models, while  
 090 also using the models to inform the clustering process.  
 091 The underlying principle of our approach is that clients  
 092 within the same cluster will exhibit similar loss patterns,  
 093 whereas those from different clusters will display distinct loss patterns across models. By analyzing  
 094 and comparing these loss patterns through the clustering of loss vectors, our approach effectively  
 095 separates clients from different clusters while grouping those from the same cluster together. This, in  
 096 turn, leads to improved models for each cluster, as each model predominantly receives updates from  
 097 clients within its corresponding cluster, thereby enhancing accuracy and robustness.

098 *CLoVE* exhibits a distinct advantage over existing CFL methods, including *IFCA*, *CFL-S* (Sattler  
 099 et al., 2021), *k-FED*, *FlexCFL* (Duan et al., 2022), and *FeSEM* (Long et al., 2023). Notably, *CLoVE*  
 100 achieves robust clustering results even from random initializations, a significant advantage over  
 101 *IFCA*, which requires careful model initialization. Unlike *k-FED*, *CLoVE* does not rely on any  
 102 assumptions on the data distribution. Compared to *CFL-S*, which often necessitates hundreds of  
 103 communication rounds to reach convergence, *CLoVE* achieves clustering in fewer than ten rounds  
 104 while reliably handling sparse client participation and stragglers. This accelerated convergence results  
 105 in substantially lower computation and communication overhead compared to both *IFCA* and *CFL-S*.  
 106 Unlike *FlexCFL*, which produces a static clustering in a single round, *CLoVE* dynamically refines  
 107 its cluster assignments, offering the dual benefits of low-cost operation and the capacity to improve  
 clustering quality over time. Moreover, *CLoVE*'s loss vector-based cluster recovery approach avoids  
 the local optima issues inherent in expectation maximization approaches such as *FeSEM*, thereby  
 ensuring optimal cluster recovery with high probability.

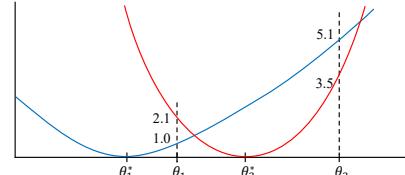


Figure 1: *Model parameters are on the x-axis and model losses on the y-axis. Blue and red curves depict losses for clients of cluster 1 and 2 respectively.  $\theta_1^*$  and  $\theta_2^*$  are optimal models for clients of clusters 1 and 2 respectively. Initial models are  $\theta_1$  and  $\theta_2$ . Loss on model  $\theta_1$  ( $\theta_2$ ) is 1.0 (5.1) and 2.1 (3.5) for clients of clusters 1 and 2, respectively. Loss vectors are [1.0, 5.1] for clients of cluster 1 and [2.1, 3.5] for clients of cluster 2. In *IFCA*, all clients select model  $\theta_1$ , as it minimizes their loss. Consequently they all get assigned to cluster 1. Moreover, because of this assignment, model  $\theta_2$  receives no updates, and therefore its loss stays high; as a result, no clients ever select it in future rounds. Ultimately, this leads to all clients converging to the same cluster, failing to recover the true client clusters. In contrast, clustering based on loss vectors achieves accurate client partitioning, as the loss vectors of clients from different clusters are highly distinct.*

108 We conduct a comprehensive evaluation of our approach, demonstrating its ability to recover clusters  
 109 with high accuracy while constructing per-cluster models in just a few rounds of FL. Furthermore,  
 110 our analytical results show that our approach can achieve single-shot, accurate cluster recovery with  
 111 high probability in the case of linear regression, thereby establishing its efficacy and robustness.  
 112

113 Our technical **contributions** can be summarized as follows:

- 114 • **Introduction of the *CLoVE* PFL algorithm:** We propose a novel approach that simultaneously  
 115 builds multiple personalized models and generates client embeddings based on model  
 116 losses, enabling effective clustering of clients with unknown cluster labels.
- 117 • **Simplicity and wide applicability:** Our method, unlike most others, extends to both  
 118 supervised and unsupervised settings, handling a variety of data distributions, sparse client  
 119 participation, and stragglers through a simple, loss-based embedding approach.
- 120 • **Robustness against initialization:** Unlike existing methods, *CLoVE* does not require careful  
 121 initialization of cluster models, thereby mitigating the risk of inaccurate cluster assignments  
 122 due to initialization sensitivities.
- 123 • **Theoretical guarantees for cluster recovery:** We provide analytical results showing that,  
 124 in the context of linear regression, *CLoVE* can achieve single-shot, accurate cluster recovery  
 125 with high probability.
- 126 • **Efficient cluster recovery and model construction:** We demonstrate through comprehensive  
 127 evaluations that *CLoVE* can recover the true clusters with high accuracy and construct  
 128 per-cluster models within a few rounds of federated learning in a variety of settings.

## 130 2 RELATED WORK

133 A primary challenge in FL is handling heterogeneous data distributions (non-IID data) and varying  
 134 device capabilities. Early optimization methods like *FedAvg* (McMahan et al., 2017) and *FedProx* (Li  
 135 et al., 2019) mitigate some issues by allowing multiple local computations before aggregation, yet  
 136 they typically yield a single global model that may not capture the diversity across user populations  
 137 (Zhao et al., 2018).

138 Personalized Federated Learning (PFL) (Hanzely & Richtárik, 2020; Tan et al., 2023) aims to  
 139 deliver client-specific models that outperform both global models and naïve local baselines. Various  
 140 strategies have been proposed, including fine-tuning global models locally (Fallah et al., 2020)  
 141 or performing adaptive local aggregation (Zhang et al., 2023), meta-learning (Jiang et al., 2023),  
 142 multi-task learning (Smith et al., 2017; Xu et al., 2023), and decomposing models into global and  
 143 personalized components (Deng et al., 2020; Mansour et al., 2020b). However, these methods operate  
 144 in a supervised setting and require labeled data.

145 Clustered Federated Learning (CFL) is an emerging approach in PFL that involves clustering clients  
 146 based on similarity measures to train specialized models for each cluster (e.g. *IFCA* based on model  
 147 loss values, (Werner et al., 2023; Kim et al., 2024) and *FlexCFL* based on gradients). The *CFL-S*  
 148 algorithm utilizes federated multitask learning and gradient embeddings in order to iteratively form  
 149 and refine clusters. Vahidian et al. (2023)'s *PACFL* algorithm derives a set of principal vectors from  
 150 each cluster's data, and *FedProto* (Tan et al., 2022) constructs class prototypes, which allow the server  
 151 to identify distribution similarities.

152 Traditional clustering methods, such as spectral clustering and k-means, have also been adapted for  
 153 federated scenarios. In *FeSEM*, personalization is achieved by deriving the optimal matching of  
 154 users and multiple model centers, and in Chen et al. (2023) through spectral co-distillation. [Licciardi et al. \(2025\) is a hierarchical approach like Sattler et al. \(2021\), and as such the number of rounds for clustering can be very large.](#) Despite their effectiveness, these methods often require complex  
 155 similarity computations, additional communication steps, [or additional amount of state kept at the server.](#)

159 Recent advancements in CFL have focused on enhancing convergence and reducing reliance on  
 160 initial cluster assignments. Improved algorithms (Vardhan et al., 2024) address some limitations of  
 161 earlier approaches that require careful initialization like *IFCA*, but their methods are complex. Our  
*CLoVE* algorithm overcomes these drawbacks using a simple method that utilizes robust loss-based

162 embeddings, converges quickly to the correct clusters, eliminates the need for precise initialization,  
 163 and operates in a variety of both supervised and unsupervised settings.  
 164

165 **3 EMBEDDINGS-DRIVEN FEDERATED CLUSTERING**  
 166

167 The typical FL architecture consists of a server and a number of clients. Let  $M$  be the number  
 168 of clients, which belong to  $K$  clusters. Data for clients of cluster  $j \in [K]$  follow a distinct data  
 169 distribution  $D_j$  (where the notation  $[K]$  denotes the set  $\{1, 2, \dots, K\}$ ). Each client  $i \in [M]$  has  $n_i$   
 170 data points. Let  $f(\theta; z) : \Theta \rightarrow \mathbb{R}$  be the loss function associated with data point  $z$ , where  $\Theta \subseteq \mathbb{R}^d$  is  
 171 the parameter space. The goal in CFL is to simultaneously cluster the clients into  $K$  clusters and train  
 172  $K$  models to minimize the population loss function  $\mathcal{L}_j(\theta) \triangleq \mathbb{E}_{z \sim D_j}[f(\theta; z)]$  for each cluster  $j \in [K]$ .  
 173 The empirical loss for the data  $Z_i$  of client  $i$  is defined as  $\mathcal{L}^i(\theta; Z_i) \triangleq (1/|Z_i|) \sum_{z \in Z_i} f(\theta; z)$ .  
 174

---

175 **Algorithm 1:** Clustering of Loss Vector Embeddings (*CLoVE*) algorithm for PFL

---

176 **Input:** # of clients  $M$ , upper bound on # of models  $\hat{K}$ , # of data points  $n_i$  for client  $i$ ,  $i \in [M]$ ,  
 177 participation rate  $\rho$

178 **Hyperparameters:** learning rate  $\gamma$ , # of rounds  $T$ , # of local epochs  $\tau$  (for model averaging)

179 **Output:** number of clusters  $K^{(T)}$ , final model parameters  $\{\theta_j^{(T)}, j \in [K^{(T)}]\}$

---

180  
 181  
 182 1 Server: initialize model parameters  $\theta_j^{(0)}$  for each of the  $K^{(0)}$  models  $j \in [K^{(0)}]$ , where  $K^{(0)} = \hat{K}$   
 183 2 **for** round  $t = 0, 1, \dots, T - 1$  **do**  
 184   3     $M^{(t)} \leftarrow$  random subset of  $\rho$  fraction of the clients (deals with partial participation/stragglers)  
 185   4    **if** clustering has not stabilized **then**  
 186   5      Server: broadcast model parameters  $\theta_j^{(t)}$  of all  $\hat{K}$  models  $j \in [\hat{K}]$  to all clients in  $M^{(t)}$   
 187   6      **for** each client  $i \in M^{(t)}$  **in parallel** **do**  
 188   7        Run client  $i$ 's data through all  $\hat{K}$  models and get the losses  $\mathcal{L}^{i,(t)} \triangleq (\mathcal{L}^i(\theta_j^{(t)}))$ ,  $j \in [\hat{K}]$   
 189   8        Send loss vector  $\mathcal{L}^{i,(t)}$  to the server  
 190   9      Server (clustering step):  
 191   10     –  $K^{(t+1)} = \text{selectBestK}(\mathcal{L}^{(t)}, \hat{K})$ , where  $\mathcal{L}^{(t)} \triangleq \{\mathcal{L}^{i,(t)}, i \in [M^{(t)}]\}$  is the loss matrix  
 192   11     – Run clustering alg. to group the set of vectors  $\{\mathcal{L}^{i,(t)}, i \in [M^{(t)}]\}$  to form  $K^{(t+1)}$  client clusters  
 193   12     – Map client clusters to models using min-cost matching and assign each client  $i \in M^{(t)}$  to the  
 194   13        corresponding model  $\kappa_i \equiv \kappa_i^{(t)} \in [K^{(t+1)}]$  of its client cluster  
 195   14   **else**  
 196   15      Server:  
 197   16      For new clients: assign them a model using models and centroids saved at stability  
 198   17      Send updated model parameters of its assigned model  $\kappa_i$  to each client  $i \in M^{(t)}$   
 199   18   **for** each client  $i \in M^{(t)}$  **in parallel** **do**  
 200   19      – **Option I** (model averaging): Compute new model parameters  $\tilde{\theta}_i = \text{LocalUpdate}(i, \theta_{\kappa_i}^{(t)}, \gamma, \tau)$   
 201   20        for model  $\kappa_i$  and send them to the server  
 202   21      – **Option II** (gradient averaging): Compute stochastic gradient  $g_i = \hat{\nabla} \mathcal{L}^i(\theta_{\kappa_i}^{(t)})$  using model  $\kappa_i$   
 203   22        and send it to the server  
 204   23   Server (averaging step):  
 205   24    Clients that were assigned to model  $j$  at round  $t$ :  $C_j^{(t)} \triangleq \{i \in M^{(t)} \text{ s.t. } \kappa_i^{(t)} = j\}$   
 206   25    – **Option I** (model averaging):  $\forall j \in [K^{(t+1)}] \text{ s.t. } C_j^{(t)} \neq \emptyset, \theta_j^{(t+1)} = \frac{\sum_{i \in C_j^{(t)}} n_i \tilde{\theta}_i}{\sum_{i \in C_j^{(t)}} n_i}$   
 207   26    – **Option II** (gradient averaging):  $\forall j \in [K^{(t+1)}] \text{ s.t. } C_j^{(t)} \neq \emptyset, \theta_j^{(t+1)} = \theta_j^{(t)} - \gamma \frac{\sum_{i \in C_j^{(t)}} n_i g_i}{\sum_{i \in C_j^{(t)}} n_i}$   
 208  
 209   27   **return**  $K^{(T)}, \{\theta_j^{(T)}, j \in [K^{(T)}]\}$

---

210  
 211  
 212  
 213  
 214 The pseudocode of our *Clustering of Loss Vector Embeddings (CLoVE)* PFL algorithm is shown in  
 215 Alg. 1. The algorithm operates as follows. Initially, the server creates  $\hat{K}$  models, where  $\hat{K}$  is the  
 upper bound on the number of models, given as an input (line 1). Until cluster stability is reached, it

216 sends all  $\hat{K}$  models to a randomly selected fraction  $\rho$  of the participating clients in round  $t$  (lines 3, 217  
 218 5). Each of these clients runs its data through all models and gets the losses, which form one loss  
 219 vector of length  $\hat{K}$  per client. These loss vectors are then sent to the server (lines 6-8). The server, if  
 220 needed, determines the number of clusters  $K^{(t+1)}$  for the next round (e.g. via searching over a range  
 221 of values and choosing the one that yields the highest silhouette score, or using the elbow method) by  
 222 calling Alg. 2. It then groups clients into  $K^{(t+1)}$  clusters (line 11) by clustering their loss vectors  
 223 (e.g. using k-means). It then assigns clients to those models to which their clusters get matched via  
 224 a minimum-cost matching (line 12). For this, a bipartite graph is created, with the nodes being the  
 225 client clusters on the left and the models on the right, and the weight of each edge being the sum of the  
 226 losses of the clients in the cluster on the left on the model on the right. Min-cost bipartite matching is  
 227 found using well known methods (e.g. max-flow-based (Munkres, 1957; Lovász & Plummer, 2009);  
 228 (Jonker & Volgenant, 1987)). Subsequently, each client trains its assigned model  $\kappa_i$  locally on its  
 229 private data (lines 17-19) by running gradient descent denoted by the function `LocalUpdate`( $\cdot$ ) for  
 230  $\tau$  epochs, starting with model  $\kappa_i$ 's parameters, and sends its updates (model parameters in Option  
 231 I and gradients in Option II) to the server. Then, the averaging step follows: the server aggregates  
 232 and applies all the received updates for each model  $j$  from the set  $C_j^{(t)}$  of clients that updated it  
 233 (lines 20-23). In the next iteration, the server sends the computed updated parameters once again  
 234 to the clients (line 5), and the process continues, until the clustering (i.e. assignment of clients to  
 235 clusters/models) stabilizes.  
 236

---

**Algorithm 2: Selection of the number of clusters**


---

**Input:** Loss matrix  $\mathcal{L}$ , initial # of models  $K^{(0)}$

**Output:** Best number of clusters  $K_{best}$

```

239 1 selectBestK( $\mathcal{L}$ ,  $K^{(0)}$ )
240 2 best_score = -1
241 3 for  $K = 2, \dots, K^{(0)}$  do
242 4    $Z_K$  = Clustering of the loss vectors (rows of  $\mathcal{L}$ ) using a clustering algorithm, e.g. Agglomerative
243   Clustering, into  $K$  clusters
244 5    $score = silhouette\_score(\mathcal{L}, Z_K)$ 
245 6   if  $score > best\_score$  then
246 7      $best\_score = score$ 
247 8      $K_{best} = K$ 
248 9 return  $K_{best}$ 

```

---

249 Checking for cluster stability can be done as follows, in a way that the server does not need to keep  
 250 any client-specific state. Each client can keep the history of its assignment to models over the different  
 251 rounds. Once a client determines that its historical assignments have been stable for a set number of  
 252 rounds required for stability, the client declares to the server that it has reached stability. The server  
 253 at each round checks whether all (or a desired percentage) of the participating clients have reached  
 254 stability, and if so, it decides that universal stability has been achieved.

255 After the algorithm reaches a stable state, the server no longer needs to compute loss vectors or  
 256 perform clustering. Instead, it can simply send the already assigned (single) model to each client  
 257 (line 16). For any client that does not yet have an assigned model, the server assigns it the model  
 258 whose cluster centroid is closest to the client's loss vector using the models and centroids saved at the  
 259 point of stability (line 15). At this stage, the algorithm behaves like standard *FedAvg*, so both the  
 260 communication and the computational overhead are minimal.

261 If the number of clusters is unknown a priori, then the choice of  $K^{(t+1)}$  at each round happens as  
 262 follows, as per Alg. 2. A range of candidate  $K$  values is examined, up to a predetermined upper  
 263 bound  $\hat{K}$ . For each candidate  $K$ , the loss vectors are clustered and the silhouette score of the resulting  
 264 clustering is computed. Eventually, the value of  $K$  with the highest silhouette score is returned. Note  
 265 that this procedure is not necessary when the number of clusters  $K$  is known a priori.

## 4 THEORETICAL ANALYSIS OF *CLoVE*'S PERFORMANCE

266 We now provide theoretical evidence supporting the superior performance of *CLoVE*. In the following  
 267 we omit all proofs, which can be found in Appendix A. To analyze the efficacy of *CLoVE*, we consider

270 a mixed linear regression problem (Yi et al., 2014) in a federated learning setting (Ghosh et al., 2022).  
 271 In this context, each client’s data originates from one of  $K$  (where  $K$  is fixed and known a priori)  
 272 distinct clusters, denoted as  $\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_K$ . Within each cluster  $\mathcal{C}_k$ , the data pairs  $z_i = (x_i, y_i)$  are  
 273 drawn from the same underlying distribution  $\mathcal{D}_k$ , which is generated by the linear model:

$$274 \quad 275 \quad y_i = \langle \theta_k^*, x_i \rangle + \epsilon_i.$$

276 Here, the feature  $x_i$  follows a standard normal distribution  $\mathcal{N}(0, I_d)$ , and the additive noise  $\epsilon_i$  follows  
 277 a normal distribution  $\mathcal{N}(0, \sigma^2)$  ( $\sigma \ll 1$ ), both independently drawn. The loss function is the squared  
 278 error:  $f(\theta; x, y) = (\langle \theta, x \rangle - y)^2$ . Under this setting, the parameters  $\theta_k^* \in \mathbb{R}^d$  are the minimizers of  
 279 the population loss for cluster  $\mathcal{C}_k$ .

280 At a high level, our analysis unfolds as follows. We demonstrate that by setting the entries of a client’s  
 281 loss vector (SLV) to the square root of model losses, *CLoVE* can accurately recover client clustering  
 282 with high probability in each round. This, in turn, enables *CLoVE* to accurately match clients to  
 283 models with high probability by utilizing the minimum-cost bipartite matching between clusters  
 284 and models. Consequently, we show that the distance between the models constructed by *CLoVE*  
 285 and their optimal counterparts decreases exponentially with the number of rounds, following the  
 286 application of client gradient updates, thus implying fast convergence of *CLoVE* to optimal models.

287 For our analysis we make some assumptions, including that all  $\theta_i^* \in \mathbb{R}^d$  have unit norms and their  
 288 minimum separation  $\Delta > 0$  is close to 1. We also assume the number of clusters  $K$  is small in  
 289 comparison to the total number of clients  $M$  and the dimension  $d$ . We assume that each client of  
 290 cluster  $k$  uses the same number  $n_k$  of i.i.d. data points independently chosen at each round.

291 For a given model  $\theta$ , the **empirical loss**  $\mathcal{L}_k^i(\theta)$  for a client  $i$  in cluster  $k$  is:

$$293 \quad 294 \quad 295 \quad \frac{1}{n_k} \sum_{j=1}^{n_k} (\langle \theta, x_j \rangle - y_j)^2 = \frac{1}{n_k} \sum_{j=1}^{n_k} (\langle \theta_k^* - \theta, x_j \rangle + \epsilon_j)^2$$

296 For a given collection of models,  $\theta = [\theta_1, \theta_2, \dots]$ , the square root loss vector (SLV) of a client  
 297  $i$  of cluster  $k$  is  $a_k^i(\theta) = [F_k^i(\theta_1), F_k^i(\theta_2), \dots]$ . Here  $F_k^i(\theta) = \sqrt{\mathcal{L}_k^i(\theta)} \sim \frac{\alpha(\theta, \theta_k^*)}{\sqrt{n_k}} \chi(n_k)$ , where  
 298  $\alpha(\theta, \theta_k^*) = \sqrt{\|\theta - \theta_k^*\|^2 + \sigma^2}$ , and  $\chi(n_k)$  denotes a standard chi random variable with  $n_k$  degrees  
 299 of freedom. Let  $c_k(\theta) = [\mathbb{E}[F_k^i(\theta_1)], \mathbb{E}[F_k^i(\theta_2)], \dots]$  denote the distribution mean of the SLVs for a  
 300 cluster  $k$ . In the following  $m_k$  denotes the number of clients in cluster  $k$  and  $n = \min_k n_k$ .

302 We say a collection of  $K$  models  $\theta = [\theta_1, \theta_2, \dots, \theta_K]$  is  $\Delta$ -proximal if each  $\theta_k$  is within a distance at  
 303 most  $\Delta/4$  of its optimal counterpart  $\theta_k^*$ . For our analysis we assume *CLoVE* is initialized with a set  $\theta$   
 304 of  $d$  randomly drawn (hence nearly ortho-normal) models. One of our key results is the following.

305 **Theorem 4.1.** *For both ortho-normal or  $\Delta$ -proximal models collections,  $k$ -means clustering of SLVs,  
 306 with suitably initialized centers, recovers accurate clustering of the clients in one shot. This result  
 307 holds with probability at least  $\xi = 1 - \delta - \frac{1}{\text{polylog}(M)}$ , for any error tolerance  $\delta$ .*

309 We prove this by showing that, for such model collections, the mixture of distributions formed by  
 310 the collection of SLVs from different clients satisfies the proximity condition of Kumar & Kannan  
 311 (2010). Specifically, it is shown in Theorem A.9 that when the number of data points across all clients  
 312 is much larger than  $d, K$ , then for any client  $i$ , with probability at least  $\xi$ :

$$313 \quad 314 \quad 315 \quad \forall j' \neq j, \|a_j^i(\theta) - c_{j'}(\theta)\| - \|a_j^i(\theta) - c_j(\theta)\| \geq \left( \frac{c' K}{m_i} + \frac{c' K}{m_j} \right) \|A(\theta) - C(\theta)\|$$

316 Here  $c'$  is a large constant and client  $i$  is assumed to belong to cluster  $j$ . Theorem 4.1 then follows  
 317 directly from the result of Kumar & Kannan (2010).

318 The proof of Theorem A.9 is based on a sequence of intermediate results. In particular, Lemma A.1  
 319 and Lemma A.2 demonstrate that, for the ortho-normal and  $\Delta$ -proximal model collections, the means  
 320 of the SLV distributions across distinct clusters are sufficiently well-separated:

321 For any  $i \neq j$ ,  $\|c_i(\theta) - c_j(\theta)\| \geq \frac{\Delta}{c}$ . Here,  $c \approx 2$ .

323 Furthermore, Lemma A.4 shows that with high probability, SLVs of each cluster are concentrated  
 324 around their means. We use these results, combined with a result of Dasgupta et al. (2007) to prove

324 Lemma A.7 which establishes that the norm of the matrix  $X(\theta)$  of “centered SLVs”, whose  $i$ -th  
 325 row is  $a_{s(i)}^i(\theta) - c_{s(i)}(\theta)$ , is bounded. That is, for  $\gamma = \sqrt{d \frac{1}{2n} (9 + \sigma^2)}$ , with high probability:  
 326  $\|X(\theta)\| \leq \gamma \sqrt{M} \text{polylog}(M)$ . By combining these results, we are able to derive the proximity  
 327 condition of Theorem A.9.

328 Next, we prove our second key result, that models constructed by *CLoVE* converge exponentially fast  
 329 towards their optimal counterpart. Let  $\theta^i$  denote model collection at round  $i$ .

330 **Theorem 4.2.** *After  $T$  rounds, each model  $\theta_k \in \theta^T$  satisfies  $\|\theta_k - \theta_k^*\| \leq \frac{\Delta}{c_1^T}$ , with high probability,  
 331 for a constant  $c_1$ .*

332 In the first round, the models  $\theta^1$  for *CLoVE* are initialized to be ortho-normal, ensuring accurate cluster  
 333 recovery by Theorem 4.1. Under the assumption that for all clusters  $k$ ,  $n_k m_k \left(\frac{\rho}{\rho+1}\right)^2 \gg \text{polylog}(K)$ ,  
 334 where  $\rho = \Delta^2/\sigma^2$ , Lemma A.16 shows that the min-cost bi-bipartite matching found by *CLoVE*  
 335 yields an optimal client-to-model match with probability close to 1. Furthermore, Lemma A.13  
 336 demonstrates that as a consequence of this optimal match, the updated  $K$  models  $\theta^2$ , obtained by  
 337 applying gradient updates from clients to their assigned models, are  $\Delta$ -proximal with high probability,  
 338 provided that the learning rate  $\eta$  is suitably chosen.

339 In subsequent rounds, the  $\Delta$ -proximality of the models  $\theta^2$  enables the application of a similar  
 340 argument, augmented by Corollary A.14. With high probability, this implies that the updated  $K$   
 341 models  $\theta^3$  not only remain  $\Delta$ -proximal, but also have a distance to their optimal counterparts that is  
 342 a constant fraction  $c_1$  times smaller, resulting from the gradient updates from their assigned clients.  
 343 Thus, by induction, the proof of Theorem 4.2 follows, provided  $\delta T \ll 1$ .

## 344 5 PERFORMANCE EVALUATION

345 We now show how our algorithms compare to the state of the art in both supervised and unsupervised  
 346 settings. In the unsupervised setting in particular, which is understudied in the context of CFL, our  
 347 evaluation is, to the best of our knowledge, the first extensive one. The datasets, models and data  
 348 distributions we use follow the ones used by the state of the art (Vahidian et al., 2023; Duan et al.,  
 349 2022; Xu et al., 2023; Zhang et al., 2023; Ghosh et al., 2020).

### 350 5.1 EXPERIMENTAL SETUP

351 **Datasets and Models** We evaluate our method on three types of tasks – image classification, text  
 352 classification, and image reconstruction – using six widely-used datasets: MNIST, Fashion-MNIST  
 353 (FMNIST), CIFAR-10, FEMNIST, Amazon Reviews, and AG News. For the classification tasks,  
 354 we use three different convolutional neural networks (CNNs): a shared CNN model for MNIST,  
 355 FMNIST and FEMNIST, a deeper CNN for CIFAR-10, a TextCNN for AG News, and an MLP-based  
 356 model (AmazonMLP) for Amazon Reviews. For the unsupervised image reconstruction tasks on  
 357 MNIST and FMNIST, we use a simple autoencoder, while a convolutional autoencoder is employed  
 358 for CIFAR-10. Further details about these datasets and model architectures can be found in App. C.

359 **Dataset Partitioning** As in the state of the art, we extensively cover many high variance and high  
 360 overlap data heterogeneity scenarios through random data partitioning, without replacement, with  
 361 various types of label and feature skews, as well as concept shifts. For label skew, we consider many  
 362 non-overlapping and overlapping label distributions. We present only some of the results here (rest  
 363 are in App. D.1.2 and D.1.3). In case *Label skew 1*, clients within each cluster receive data from  
 364 only two unique classes, with no label overlap between clients in different clusters. Conversely, in  
 365 case *Label skew 2*, each client within each cluster receives samples from  $U$  classes, with at least  
 366  $V$  labels shared between the labels assigned to any two clusters. For the AG News dataset, we use  
 367  $U = 2$  and  $V = 1$ , while for all image datasets we use  $U = 4$  and  $V = 2$ . **As is standard practice**  
 368 **in the FL literature, for *Label skew 1* and *2* (and *3* in App. D.1.2), the data of a class is distributed**  
 369 **amongst the clients that are assigned to that class by sampling from a Dirichlet distribution with**  
 370 **parameter  $\alpha = 0.5$ .** For instance, suppose the clients are divided into three clusters, and a particular  
 371 class contains  $L$  data points. First, we draw three proportions  $p_1, p_2, p_3$  from a Dirichlet distribution  
 372 with concentration parameter  $\alpha = 0.5$  (by construction,  $p_1 + p_2 + p_3 = 1$ ). Then we randomly select  
 373  $Lp_1$  points of the class for the clients in cluster 1,  $Lp_2$  points for the clients in cluster 2, and the  
 374  $Lp_3$  points for the clients in cluster 3. **Label skew 3** is similar, but the data is partitioned into three  
 375 clusters with  $U = 3$  and  $V = 2$ . **Label skew 4** is similar, but the data is partitioned into four clusters with  
 376  $U = 4$  and  $V = 3$ . **Label skew 5** is similar, but the data is partitioned into five clusters with  $U = 5$  and  
 377  $V = 4$ . **Label skew 6** is similar, but the data is partitioned into six clusters with  $U = 6$  and  $V = 5$ . **Label skew 7**  
 378 is similar, but the data is partitioned into seven clusters with  $U = 7$  and  $V = 6$ . **Label skew 8** is similar,  
 379 but the data is partitioned into eight clusters with  $U = 8$  and  $V = 7$ . **Label skew 9** is similar, but the data  
 380 is partitioned into nine clusters with  $U = 9$  and  $V = 8$ . **Label skew 10** is similar, but the data is partitioned  
 381 into ten clusters with  $U = 10$  and  $V = 9$ . **Label skew 11** is similar, but the data is partitioned into eleven  
 382 clusters with  $U = 11$  and  $V = 10$ . **Label skew 12** is similar, but the data is partitioned into twelve clusters with  
 383  $U = 12$  and  $V = 11$ . **Label skew 13** is similar, but the data is partitioned into thirteen clusters with  
 384  $U = 13$  and  $V = 12$ . **Label skew 14** is similar, but the data is partitioned into fourteen clusters with  
 385  $U = 14$  and  $V = 13$ . **Label skew 15** is similar, but the data is partitioned into fifteen clusters with  
 386  $U = 15$  and  $V = 14$ . **Label skew 16** is similar, but the data is partitioned into sixteen clusters with  
 387  $U = 16$  and  $V = 15$ . **Label skew 17** is similar, but the data is partitioned into seventeen clusters with  
 388  $U = 17$  and  $V = 16$ . **Label skew 18** is similar, but the data is partitioned into eighteen clusters with  
 389  $U = 18$  and  $V = 17$ . **Label skew 19** is similar, but the data is partitioned into nineteen clusters with  
 390  $U = 19$  and  $V = 18$ . **Label skew 20** is similar, but the data is partitioned into twenty clusters with  
 391  $U = 20$  and  $V = 19$ . **Label skew 21** is similar, but the data is partitioned into twenty-one clusters with  
 392  $U = 21$  and  $V = 20$ . **Label skew 22** is similar, but the data is partitioned into twenty-two clusters with  
 393  $U = 22$  and  $V = 21$ . **Label skew 23** is similar, but the data is partitioned into twenty-three clusters with  
 394  $U = 23$  and  $V = 22$ . **Label skew 24** is similar, but the data is partitioned into twenty-four clusters with  
 395  $U = 24$  and  $V = 23$ . **Label skew 25** is similar, but the data is partitioned into twenty-five clusters with  
 396  $U = 25$  and  $V = 24$ . **Label skew 26** is similar, but the data is partitioned into twenty-six clusters with  
 397  $U = 26$  and  $V = 25$ . **Label skew 27** is similar, but the data is partitioned into twenty-seven clusters with  
 398  $U = 27$  and  $V = 26$ . **Label skew 28** is similar, but the data is partitioned into twenty-eight clusters with  
 399  $U = 28$  and  $V = 27$ . **Label skew 29** is similar, but the data is partitioned into twenty-nine clusters with  
 400  $U = 29$  and  $V = 28$ . **Label skew 30** is similar, but the data is partitioned into thirty clusters with  
 401  $U = 30$  and  $V = 29$ . **Label skew 31** is similar, but the data is partitioned into thirty-one clusters with  
 402  $U = 31$  and  $V = 30$ . **Label skew 32** is similar, but the data is partitioned into thirty-two clusters with  
 403  $U = 32$  and  $V = 31$ . **Label skew 33** is similar, but the data is partitioned into thirty-three clusters with  
 404  $U = 33$  and  $V = 32$ . **Label skew 34** is similar, but the data is partitioned into thirty-four clusters with  
 405  $U = 34$  and  $V = 33$ . **Label skew 35** is similar, but the data is partitioned into thirty-five clusters with  
 406  $U = 35$  and  $V = 34$ . **Label skew 36** is similar, but the data is partitioned into thirty-six clusters with  
 407  $U = 36$  and  $V = 35$ . **Label skew 37** is similar, but the data is partitioned into thirty-seven clusters with  
 408  $U = 37$  and  $V = 36$ . **Label skew 38** is similar, but the data is partitioned into thirty-eight clusters with  
 409  $U = 38$  and  $V = 37$ . **Label skew 39** is similar, but the data is partitioned into thirty-nine clusters with  
 410  $U = 39$  and  $V = 38$ . **Label skew 40** is similar, but the data is partitioned into forty clusters with  
 411  $U = 40$  and  $V = 39$ . **Label skew 41** is similar, but the data is partitioned into forty-one clusters with  
 412  $U = 41$  and  $V = 40$ . **Label skew 42** is similar, but the data is partitioned into forty-two clusters with  
 413  $U = 42$  and  $V = 41$ . **Label skew 43** is similar, but the data is partitioned into forty-three clusters with  
 414  $U = 43$  and  $V = 42$ . **Label skew 44** is similar, but the data is partitioned into forty-four clusters with  
 415  $U = 44$  and  $V = 43$ . **Label skew 45** is similar, but the data is partitioned into forty-five clusters with  
 416  $U = 45$  and  $V = 44$ . **Label skew 46** is similar, but the data is partitioned into forty-six clusters with  
 417  $U = 46$  and  $V = 45$ . **Label skew 47** is similar, but the data is partitioned into forty-seven clusters with  
 418  $U = 47$  and  $V = 46$ . **Label skew 48** is similar, but the data is partitioned into forty-eight clusters with  
 419  $U = 48$  and  $V = 47$ . **Label skew 49** is similar, but the data is partitioned into forty-nine clusters with  
 420  $U = 49$  and  $V = 48$ . **Label skew 50** is similar, but the data is partitioned into fifty clusters with  
 421  $U = 50$  and  $V = 49$ . **Label skew 51** is similar, but the data is partitioned into fifty-one clusters with  
 422  $U = 51$  and  $V = 50$ . **Label skew 52** is similar, but the data is partitioned into fifty-two clusters with  
 423  $U = 52$  and  $V = 51$ . **Label skew 53** is similar, but the data is partitioned into fifty-three clusters with  
 424  $U = 53$  and  $V = 52$ . **Label skew 54** is similar, but the data is partitioned into fifty-four clusters with  
 425  $U = 54$  and  $V = 53$ . **Label skew 55** is similar, but the data is partitioned into fifty-five clusters with  
 426  $U = 55$  and  $V = 54$ . **Label skew 56** is similar, but the data is partitioned into fifty-six clusters with  
 427  $U = 56$  and  $V = 55$ . **Label skew 57** is similar, but the data is partitioned into fifty-seven clusters with  
 428  $U = 57$  and  $V = 56$ . **Label skew 58** is similar, but the data is partitioned into fifty-eight clusters with  
 429  $U = 58$  and  $V = 57$ . **Label skew 59** is similar, but the data is partitioned into fifty-nine clusters with  
 430  $U = 59$  and  $V = 58$ . **Label skew 60** is similar, but the data is partitioned into六十 clusters with  
 431  $U = 60$  and  $V = 59$ . **Label skew 61** is similar, but the data is partitioned into六十-one clusters with  
 432  $U = 61$  and  $V = 60$ . **Label skew 62** is similar, but the data is partitioned into六十-two clusters with  
 433  $U = 62$  and  $V = 61$ . **Label skew 63** is similar, but the data is partitioned into六十-three clusters with  
 434  $U = 63$  and  $V = 62$ . **Label skew 64** is similar, but the data is partitioned into六十-four clusters with  
 435  $U = 64$  and  $V = 63$ . **Label skew 65** is similar, but the data is partitioned into六十-five clusters with  
 436  $U = 65$  and  $V = 64$ . **Label skew 66** is similar, but the data is partitioned into六十-six clusters with  
 437  $U = 66$  and  $V = 65$ . **Label skew 67** is similar, but the data is partitioned into六十-seven clusters with  
 438  $U = 67$  and  $V = 66$ . **Label skew 68** is similar, but the data is partitioned into六十-eight clusters with  
 439  $U = 68$  and  $V = 67$ . **Label skew 69** is similar, but the data is partitioned into六十-nine clusters with  
 440  $U = 69$  and  $V = 68$ . **Label skew 70** is similar, but the data is partitioned into七十 clusters with  
 441  $U = 70$  and  $V = 69$ . **Label skew 71** is similar, but the data is partitioned into七十-one clusters with  
 442  $U = 71$  and  $V = 70$ . **Label skew 72** is similar, but the data is partitioned into七十-two clusters with  
 443  $U = 72$  and  $V = 71$ . **Label skew 73** is similar, but the data is partitioned into七十-three clusters with  
 444  $U = 73$  and  $V = 72$ . **Label skew 74** is similar, but the data is partitioned into七十-four clusters with  
 445  $U = 74$  and  $V = 73$ . **Label skew 75** is similar, but the data is partitioned into七十-five clusters with  
 446  $U = 75$  and  $V = 74$ . **Label skew 76** is similar, but the data is partitioned into七十-six clusters with  
 447  $U = 76$  and  $V = 75$ . **Label skew 77** is similar, but the data is partitioned into七十-seven clusters with  
 448  $U = 77$  and  $V = 76$ . **Label skew 78** is similar, but the data is partitioned into七十-eight clusters with  
 449  $U = 78$  and  $V = 77$ . **Label skew 79** is similar, but the data is partitioned into七十-nine clusters with  
 450  $U = 79$  and  $V = 78$ . **Label skew 80** is similar, but the data is partitioned into八十 clusters with  
 451  $U = 80$  and  $V = 79$ . **Label skew 81** is similar, but the data is partitioned into八十-one clusters with  
 452  $U = 81$  and  $V = 80$ . **Label skew 82** is similar, but the data is partitioned into八十-two clusters with  
 453  $U = 82$  and  $V = 81$ . **Label skew 83** is similar, but the data is partitioned into八十-three clusters with  
 454  $U = 83$  and  $V = 82$ . **Label skew 84** is similar, but the data is partitioned into八十-four clusters with  
 455  $U = 84$  and  $V = 83$ . **Label skew 85** is similar, but the data is partitioned into八十-five clusters with  
 456  $U = 85$  and  $V = 84$ . **Label skew 86** is similar, but the data is partitioned into八十-six clusters with  
 457  $U = 86$  and  $V = 85$ . **Label skew 87** is similar, but the data is partitioned into八十-seven clusters with  
 458  $U = 87$  and  $V = 86$ . **Label skew 88** is similar, but the data is partitioned into八十-eight clusters with  
 459  $U = 88$  and  $V = 87$ . **Label skew 89** is similar, but the data is partitioned into八十-nine clusters with  
 460  $U = 89$  and  $V = 88$ . **Label skew 90** is similar, but the data is partitioned into九十 clusters with  
 461  $U = 90$  and  $V = 89$ . **Label skew 91** is similar, but the data is partitioned into九十-one clusters with  
 462  $U = 91$  and  $V = 90$ . **Label skew 92** is similar, but the data is partitioned into九十-two clusters with  
 463  $U = 92$  and  $V = 91$ . **Label skew 93** is similar, but the data is partitioned into九十-three clusters with  
 464  $U = 93$  and  $V = 92$ . **Label skew 94** is similar, but the data is partitioned into九十-four clusters with  
 465  $U = 94$  and  $V = 93$ . **Label skew 95** is similar, but the data is partitioned into九十-five clusters with  
 466  $U = 95$  and  $V = 94$ . **Label skew 96** is similar, but the data is partitioned intoninety clusters with  
 467  $U = 96$  and  $V = 95$ . **Label skew 97** is similar, but the data is partitioned intoninety-one clusters with  
 468  $U = 97$  and  $V = 96$ . **Label skew 98** is similar, but the data is partitioned intoninety-two clusters with  
 469  $U = 98$  and  $V = 97$ . **Label skew 99** is similar, but the data is partitioned intoninety-three clusters with  
 470  $U = 99$  and  $V = 98$ . **Label skew 100** is similar, but the data is partitioned intoninety-four clusters with  
 471  $U = 100$  and  $V = 99$ . **Label skew 101** is similar, but the data is partitioned intoninety-five clusters with  
 472  $U = 101$  and  $V = 100$ . **Label skew 102** is similar, but the data is partitioned intoninety-six clusters with  
 473  $U = 102$  and  $V = 101$ . **Label skew 103** is similar, but the data is partitioned intoninety-seven clusters with  
 474  $U = 103$  and  $V = 102$ . **Label skew 104** is similar, but the data is partitioned intoninety-eight clusters with  
 475  $U = 104$  and  $V = 103$ . **Label skew 105** is similar, but the data is partitioned intoninety-nine clusters with  
 476  $U = 105$  and  $V = 104$ . **Label skew 106** is similar, but the data is partitioned intoninety clusters with  
 477  $U = 106$  and  $V = 105$ . **Label skew 107** is similar, but the data is partitioned intoninety-one clusters with  
 478  $U = 107$  and  $V = 106$ . **Label skew 108** is similar, but the data is partitioned intoninety-two clusters with  
 479  $U = 108$  and  $V = 107$ . **Label skew 109** is similar, but the data is partitioned intoninety-three clusters with  
 480  $U = 109$  and  $V = 108$ . **Label skew 110** is similar, but the data is partitioned intoninety-four clusters with  
 481  $U = 110$  and  $V = 109$ . **Label skew 111** is similar, but the data is partitioned intoninety-five clusters with  
 482  $U = 111$  and  $V = 110$ . **Label skew 112** is similar, but the data is partitioned intoninety-six clusters with  
 483  $U = 112$  and  $V = 111$ . **Label skew 113** is similar, but the data is partitioned intoninety-seven clusters with  
 484  $U = 113$  and  $V = 112$ . **Label skew 114** is similar, but the data is partitioned intoninety-eight clusters with  
 485  $U = 114$  and  $V = 113$ . **Label skew 115** is similar, but the data is partitioned intoninety-nine clusters with  
 486  $U = 115$  and  $V = 114$ . **Label skew 116** is similar, but the data is partitioned intoninety clusters with  
 487  $U = 116$  and  $V = 115$ . **Label skew 117** is similar, but the data is partitioned intoninety-one clusters with  
 488  $U = 117$  and  $V = 116$ . **Label skew 118** is similar, but the data is partitioned intoninety-two clusters with  
 489  $U = 118$  and  $V = 117$ . **Label skew 119** is similar, but the data is partitioned intoninety-three clusters with  
 490  $U = 119$  and  $V = 118$ . **Label skew 120** is similar, but the data is partitioned intoninety-four clusters with  
 491  $U = 120$  and  $V = 119$ . **Label skew 121** is similar, but the data is partitioned intoninety-five clusters with  
 492  $U = 121$  and  $V = 120$ . **Label skew 122** is similar, but the data is partitioned intoninety-six clusters with  
 493  $U = 122$  and  $V = 121$ . **Label skew 123** is similar, but the data is partitioned intoninety-seven clusters with  
 494  $U = 123$  and  $V = 122$ . **Label skew 124** is similar, but the data is partitioned intoninety-eight clusters with  
 495  $U = 124$  and  $V = 123$ . **Label skew 125** is similar, but the data is partitioned intoninety-nine clusters with  
 496  $U = 125$  and  $V = 124$ . **Label skew 126** is similar, but the data is partitioned intoninety clusters with  
 497  $U = 126$  and  $V = 125$ . **Label skew 127** is similar, but the data is partitioned intoninety-one clusters with  
 498  $U = 127$  and  $V = 126$ . **Label skew 128** is similar, but the data is partitioned intoninety-two clusters with  
 499  $U = 128$  and  $V = 127$ . **Label skew 129** is similar, but the data is partitioned intoninety-three clusters with  
 500  $U = 129$  and  $V = 128$ . **Label skew 130** is similar, but the data is partitioned intoninety-four clusters with  
 501  $U = 130$  and  $V = 129$ . **Label skew 131** is similar, but the data is partitioned intoninety-five clusters with  
 502  $U = 131$  and  $V = 130$ . **Label skew 132** is similar, but the data is partitioned intoninety-six clusters with  
 503  $U = 132$  and  $V = 131$ . **Label skew 133** is similar, but the data is partitioned intoninety-seven clusters with  
 504  $U = 133$  and  $V = 132$ . **Label skew 134** is similar, but the data is partitioned intoninety-eight clusters with  
 505  $U = 134$  and  $V = 133$ . **Label skew 135** is similar, but the data is partitioned intoninety-nine clusters with  
 506  $U = 135$  and  $V = 134$ . **Label skew 136** is similar, but the data is partitioned intoninety clusters with  
 507  $U = 136$  and  $V = 135$ . **Label skew 137** is similar, but the data is partitioned intoninety-one clusters with  
 508  $U = 137$  and  $V = 136$ . **Label skew 138** is similar, but the data is partitioned intoninety-two clusters with  
 509  $U = 138$  and  $V = 137$ . **Label skew 139** is similar, but the data is partitioned intoninety-three clusters with  
 510  $U = 139$  and  $V = 138$ . **Label skew 140** is similar, but the data is partitioned intoninety-four clusters with  
 511  $U = 140$  and  $V = 139$ . **Label skew 141** is similar, but the data is partitioned intoninety-five clusters with  
 512  $U = 141$  and  $V = 140$ . **Label skew 142** is similar, but the data is partitioned intoninety-six clusters with  
 513  $U = 142$  and  $V = 141$ . **Label skew 143** is similar, but the data is partitioned intoninety-seven clusters with  
 514  $U = 143$  and  $V = 142$ . **Label skew 144** is similar, but the data is partitioned intoninety-eight clusters with  
 515  $U = 144$  and  $V = 143$ . **Label skew 145** is similar, but the data is partitioned intoninety-nine clusters with  
 516  $U = 145$  and  $V = 144$ . **Label skew 146** is similar, but the data is partitioned intoninety clusters with  
 517  $U = 146$  and  $V = 145$ . **Label skew 147** is similar, but the data is partitioned intoninety-one clusters with  
 518  $U = 147$  and  $V = 146$ . **Label skew 148** is similar, but the data is partitioned intoninety-two clusters with  
 519  $U = 148$  and  $V = 147$ . **Label skew 149** is similar, but the data is partitioned intoninety-three clusters with  
 520  $U = 149$  and  $V = 148$ . **Label skew 150** is similar, but the data is partitioned intoninety-four clusters with  
 521  $U = 150$  and  $V = 149$ . **Label skew 151** is similar, but the data is partitioned intoninety-five clusters with  
 522  $U = 151$  and  $V = 150$ . **Label skew 152** is similar, but the data is partitioned intoninety-six clusters with  
 523  $U = 152$  and  $V = 151$ . **Label skew 153** is similar, but the data is partitioned intoninety-seven clusters with  
 524  $U = 153$  and  $V = 152$ . **Label skew 154** is similar, but the data is partitioned intoninety-eight clusters with  
 525  $U = 154$  and  $V = 153$ . **Label skew 155** is similar, but the data is partitioned intoninety-nine clusters with  
 526  $U = 155$  and  $V = 154$ . **Label skew 156** is similar, but the data is partitioned intoninety clusters with  
 527  $U = 156$  and  $V = 155$ . **Label skew 157** is similar, but the data is partitioned intoninety-one clusters with  
 528  $U = 157$  and  $V = 156$ . **Label skew 158** is similar, but the data is partitioned intoninety-two clusters with  
 529  $U = 158$  and  $V = 157$ . **Label skew 159** is similar, but the data is partitioned intoninety-three clusters with  
 530  $U = 159$  and  $V = 158$ . **Label skew 160** is similar, but the data is partitioned intoninety-four clusters with  
 531  $U = 160$  and  $V = 159$ . **Label skew 161** is similar, but the data is partitioned intoninety-five clusters with  
 532  $U = 161$  and  $V = 160$ . **Label skew 162** is similar, but the data is partitioned intoninety-six clusters with  
 533  $U = 162$  and  $V = 161$ . **Label skew 163** is similar, but the data is partitioned intoninety-seven clusters with  
 534  $U = 163$  and  $V = 162$ . **Label skew 164** is similar, but the data is partitioned intoninety-eight clusters with  
 535  $U = 164$  and  $V = 163$ . **Label skew 165** is similar, but the data is partitioned intoninety-nine clusters with  
 536  $U = 165$  and  $V = 164$ . **Label skew 166** is similar, but the data is partitioned intoninety clusters with  
 537  $U = 166$  and  $V = 165$ . **Label skew 167** is similar, but the data is partitioned intoninety-one clusters with  
 538  $U = 167$  and  $V = 166$ . **Label skew 168** is similar, but the data is partitioned intoninety-two clusters with  
 539  $U = 168$  and  $V = 167$ . **Label skew 169** is similar, but the data is partitioned intoninety-three clusters

378 remaining  $Lp_3$  points for the clients in cluster 3. The selected points are assigned randomly among  
 379 the clients within each respective cluster.  
 380

381 To simulate feature distribution shifts, we apply image rotations ( $0^\circ, 90^\circ, 180^\circ, 270^\circ$ ) to the MNIST,  
 382 CIFAR-10 and FMNIST datasets. Each client within a cluster is assigned data with a specific rotation.  
 383 Concept shift is achieved through label permutation: all labels are distributed across all clusters, but  
 384 clients within each cluster receive data with two unique label swaps. Both test and train data of a  
 385 client have the same distribution. No data mixing is employed for the Amazon Review dataset, i.e.  
 386 we use its original data classes.  
 387

388 **Compared Methods** We compare the following baselines: *Local-only*, where each client trains its  
 389 model locally; *FedAvg*, which learns a single global model; clustered federated learning methods,  
 390 including *CFL-S*, *IFCA*, *FlexCFL*, *FeSEM*, and *PACFL* (Vahidian et al., 2023); and general PFL  
 391 methods *FedProto* (Tan et al., 2022), *Per-FedAvg* (Fallah et al., 2020), *FedALA* (Zhang et al., 2023),  
 392 and *FedPAC* (Xu et al., 2023). The source code we used for these algorithms is linked in Table 7.  
 393

394 **Training Settings** We use an Adam (Kingma & Ba, 2017) optimizer, a batch size of 64 and a  
 395 learning rate  $\gamma = 10^{-3}$ . The number of local training epochs is  $\tau = 1$  and the number of global  
 396 communication rounds is  $T = 100$ . For supervised classification tasks we use cross-entropy loss, and  
 397 for unsupervised tasks we use the MSE reconstruction loss. We vary the number of clients between  
 398 20 and 1000. In most cases, the training data size per client is set to 500. However, for the Amazon  
 399 Reviews and AG News datasets, we use 1000 and 100 samples, respectively. Details on the number  
 400 of clients and training data sizes for all experiments is provided in App. C.  
 401

402 **Result reporting** We report the average performance of the models assigned to the clients on  
 403 their local test data. Specifically, we report test accuracies for supervised classification tasks and  
 404 reconstruction losses for unsupervised tasks. In addition, we report the accuracy of client-to-cluster  
 405 assignments as the Adjusted Rand Index (ARI) between the achieved clustering and the ground-truth  
 406 client groupings. To account for randomness, we run each experiment for 3 values of a randomness  
 407 seed and report the mean and standard deviation. In this section, we only report results for full  
 408 participation, known number of clusters, and no early stopping of clustering. The results for the other  
 409

410 Table 1: Supervised test accuracy results for MNIST, CIFAR-10 and FMNIST

Algorithm	Data mixing	MNIST	CIFAR-10	FMNIST	Data mixing	MNIST	CIFAR-10	FMNIST
FedAvg	Label skew 1	87.4 $\pm$ 0.9	34.0 $\pm$ 3.6	64.7 $\pm$ 1.2	Concept shift	52.4 $\pm$ 0.3	<b>63.1 <math>\pm</math> 1.7</b>	49.1 $\pm$ 0.1
Local-only		99.5 $\pm$ 0.1	85.4 $\pm$ 1.6	98.8 $\pm$ 0.1		94.8 $\pm$ 0.3	45.1 $\pm$ 0.4	78.3 $\pm$ 0.4
Per-FedAvg		97.8 $\pm$ 0.2	75.0 $\pm$ 0.7	97.0 $\pm$ 0.1		80.2 $\pm$ 1.0	34.9 $\pm$ 0.7	70.4 $\pm$ 1.0
FedProto		99.2 $\pm$ 0.0	83.0 $\pm$ 0.3	98.6 $\pm$ 0.0		95.0 $\pm$ 0.4	40.8 $\pm$ 0.3	78.4 $\pm$ 0.3
FedALA		98.9 $\pm$ 0.1	80.0 $\pm$ 0.6	97.8 $\pm$ 0.0		96.7 $\pm$ 0.2	42.0 $\pm$ 1.0	81.1 $\pm$ 0.3
FedPAC		94.4 $\pm$ 0.7	78.3 $\pm$ 0.4	91.7 $\pm$ 1.2		90.7 $\pm$ 1.3	39.0 $\pm$ 1.0	73.8 $\pm$ 0.6
CFL-S		82.5 $\pm$ 21.8	87.4 $\pm$ 0.4	65.8 $\pm$ 22.4		97.5 $\pm$ 0.2	35.7 $\pm$ 3.2	83.7 $\pm$ 0.5
FeSEM		73.2 $\pm$ 0.8	11.2 $\pm$ 1.6	41.8 $\pm$ 0.7		46.2 $\pm$ 0.1	12.7 $\pm$ 0.5	39.7 $\pm$ 0.6
FlexCFL		99.4 $\pm$ 0.0	73.2 $\pm$ 0.3	99.0 $\pm$ 0.1		97.4 $\pm$ 0.1	17.8 $\pm$ 2.1	<b>86.1 <math>\pm</math> 0.2</b>
PACFL		99.1 $\pm$ 0.1	79.1 $\pm$ 1.0	98.7 $\pm$ 0.1		93.3 $\pm$ 0.1	26.8 $\pm$ 1.2	72.6 $\pm$ 1.2
IFCA	Label skew 2	99.1 $\pm$ 0.5	85.7 $\pm$ 6.5	97.8 $\pm$ 1.9	Feature skew	80.0 $\pm$ 0.2	52.7 $\pm$ 4.7	67.0 $\pm$ 6.2
CLoVE		<b>99.8 <math>\pm</math> 0.0</b>	<b>90.5 <math>\pm</math> 0.1</b>	<b>99.1 <math>\pm</math> 0.0</b>		<b>97.6 <math>\pm</math> 0.1</b>	58.4 $\pm$ 0.6	84.3 $\pm$ 0.2
FedAvg		76.4 $\pm$ 0.3	47.7 $\pm$ 4.0	55.2 $\pm$ 1.6		90.0 $\pm$ 0.6	52.4 $\pm$ 2.5	76.5 $\pm$ 0.3
Local-only		94.2 $\pm$ 0.2	55.9 $\pm$ 1.8	86.1 $\pm$ 0.4		93.0 $\pm$ 0.1	38.5 $\pm$ 1.6	75.8 $\pm$ 0.2
Per-FedAvg		94.9 $\pm$ 0.6	48.2 $\pm$ 0.9	78.2 $\pm$ 1.0		87.8 $\pm$ 0.5	22.8 $\pm$ 0.1	66.7 $\pm$ 0.4
FedProto	Label skew 2	94.4 $\pm$ 0.2	56.8 $\pm$ 0.8	82.1 $\pm$ 0.9	Feature skew	92.3 $\pm$ 0.3	34.1 $\pm$ 0.2	74.0 $\pm$ 0.3
FedALA		95.1 $\pm$ 0.3	50.1 $\pm$ 1.6	81.6 $\pm$ 0.7		90.2 $\pm$ 0.6	22.3 $\pm$ 0.4	68.1 $\pm$ 0.5
FedPAC		91.9 $\pm$ 0.5	40.7 $\pm$ 5.2	77.6 $\pm$ 1.2		87.6 $\pm$ 0.4	24.6 $\pm$ 0.2	67.4 $\pm$ 0.2
CFL-S		<b>97.3 <math>\pm</math> 0.3</b>	62.1 $\pm$ 0.4	87.2 $\pm$ 1.2		95.1 $\pm$ 0.2	<b>61.6 <math>\pm</math> 5.1</b>	77.8 $\pm$ 3.5
FeSEM		89.1 $\pm$ 1.6	15.3 $\pm$ 2.5	51.4 $\pm$ 0.2		79.5 $\pm$ 0.4	10.7 $\pm$ 0.5	67.6 $\pm$ 0.5
FlexCFL	Label skew 2	95.6 $\pm$ 0.3	43.6 $\pm$ 0.6	88.7 $\pm$ 0.2	Feature skew	96.7 $\pm$ 0.1	12.2 $\pm$ 0.4	<b>85.2 <math>\pm</math> 0.2</b>
PACFL		92.0 $\pm$ 1.1	40.4 $\pm$ 1.0	81.1 $\pm$ 0.5		89.3 $\pm$ 0.1	17.1 $\pm$ 0.6	70.6 $\pm$ 0.6
IFCA		96.8 $\pm$ 0.2	63.4 $\pm$ 3.4	89.4 $\pm$ 1.3		95.9 $\pm$ 2.0	50.8 $\pm$ 7.0	83.5 $\pm$ 0.3
CLoVE		97.1 $\pm$ 0.3	<b>66.8 <math>\pm</math> 2.3</b>	<b>90.1 <math>\pm</math> 0.3</b>		<b>97.7 <math>\pm</math> 0.1</b>	57.1 $\pm$ 1.8	85.1 $\pm$ 0.2

432

433

Table 2: Supervised ARI results for MNIST, CIFAR-10 and FMNIST

Algorithm	Data mixing	MNIST	CIFAR-10	FMNIST	Data mixing	MNIST	CIFAR-10	FMNIST
CFL-S FeSEM FlexCFL PACFL IFCA CLOVE	Label skew 1	0.71 ± 0.20	<b>1.00 ± 0.00</b>	0.86 ± 0.20	Concept shift	<b>1.00 ± 0.00</b>	0.00 ± 0.00	<b>1.00 ± 0.00</b>
		0.18 ± 0.00	0.00 ± 0.00	0.00 ± 0.00		0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
		<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>		<b>1.00 ± 0.00</b>	0.22 ± 0.24	<b>1.00 ± 0.00</b>
		0.39 ± 0.07	0.60 ± 0.09	<b>1.00 ± 0.00</b>		0.00 ± 0.02	0.00 ± 0.00	0.00 ± 0.00
		0.83 ± 0.12	0.91 ± 0.13	0.92 ± 0.12		0.68 ± 0.00	0.76 ± 0.17	0.55 ± 0.18
		<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>		<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>
CFL-S FeSEM FlexCFL PACFL IFCA CLOVE	Label skew 2	0.00 ± 0.00	0.72 ± 0.00	0.57 ± 0.00	Feature skew	0.48 ± 0.00	0.89 ± 0.15	0.16 ± 0.23
		0.18 ± 0.00	0.00 ± 0.00	0.00 ± 0.00		0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
		<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>		<b>1.00 ± 0.00</b>	0.11 ± 0.08	<b>1.00 ± 0.00</b>
		0.35 ± 0.12	0.32 ± 0.11	0.55 ± 0.03		0.00 ± 0.00	0.00 ± 0.00	0.68 ± 0.11
		0.83 ± 0.12	0.81 ± 0.13	0.92 ± 0.12		0.67 ± 0.28	0.71 ± 0.22	0.55 ± 0.10
		<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>		<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>

445

cases, as well as scaling experiments and ablation studies, are provided in App. D and show that *CLoVE* performs well even in these settings.

448

## 5.2 NUMERICAL RESULTS

449

The test accuracies for image classification tasks are reported in Table 1. The results demonstrate that *CLoVE* performs consistently well across various data distributions, including label and feature skews, as well as concept shifts. While some baseline algorithms exhibit strong performance in specific cases, ***CLoVE* is the only algorithm that consistently ranks among the top performers in almost all scenarios.**

450

Table 2 presents the clustering accuracies of different CFL algorithms for the image datasets under the same experimental setup. *CLoVE* achieves optimal performance in all cases. Although *FlexCFL* also shows near-optimal performance, its accuracy declines when faced with concept shifts, especially on CIFAR-10, as it neglects label information in similarity estimation. In contrast, loss-based approaches like *CLoVE* excel in this metric, as model losses can effectively account for various types of skews. Results for additional types of label skews (e.g. dominant class (Xu et al., 2023)) are included in App. D.1.2.

451

The test accuracies for the textual classification tasks are reported in Table 3 (more metrics in App. D.1.3). While *CLoVE* exhibits good performance here as well, the baseline *FedAvg* performs especially well on the Amazon Review dataset. One reason is that while the dataset is split by product categories, the underlying language and sentiment cues are still fairly similar across all 4 categories and can provide a good signal for classification. On AG News, *CLoVE* outperforms all baselines except *FlexCFL*.

452

Table 5 presents the clustering recovery speed for the image classification datasets under different client data distributions. The first 3 columns show the ARI reached within 10 rounds for each algorithm. *CLoVE* outperforms all baselines here. The next 3 columns show the first round at which an ARI of 0.9 or higher is achieved (dashes mean 0.9 ARI is never achieved). As can be seen, ***CLoVE* consistently reaches such high accuracy within 2-3 rounds, unlike any other baseline.**

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

**Unsupervised Setting** Unlike many of the baselines, *CLoVE* is equally applicable to unsupervised settings, as it relies solely on model losses to identify clients with similar data distributions. Another loss-based approach, *IFCA*, serves as a natural point of comparison in this context. As shown in Table 4, *CLoVE* exhibits consistently good performance compared to *IFCA* and other baselines in unsupervised settings on image reconstruction

Table 3: Supervised test accuracy results for Amazon Review and AG News

Algorithm	Amazon	AG News
FedAvg	<b>99.0 ± 0.1</b>	40.2 ± 0.1
Local-only	81.2 ± 0.4	51.4 ± 1.1
Per-FedAvg	87.7 ± 0.4	48.0 ± 1.1
FedProto	81.3 ± 0.1	18.9 ± 1.5
FedALA	88.3 ± 0.3	49.6 ± 0.1
FedPAC	88.2 ± 0.2	42.4 ± 2.9
CFL-S	87.7 ± 0.8	48.4 ± 2.7
FeSEM	50.9 ± 0.2	19.1 ± 1.5
FlexCFL	84.1 ± 0.8	<b>60.2 ± 5.5</b>
PACFL	82.1 ± 0.3	51.1 ± 1.9
IFCA	85.5 ± 0.5	51.3 ± 3.3
<b>CLoVE</b>	86.8 ± 0.4	55.1 ± 0.2

Table 4: Unsupervised test loss results for MNIST, CIFAR-10 and FMNIST

Algorithm	MNIST	CIFAR-10	FMNIST
FedAvg	0.032 ± 0.001	0.026 ± 0.000	0.030 ± 0.000
Local-only	0.019 ± 0.000	0.022 ± 0.000	0.170 ± 0.000
IFCA	0.024 ± 0.005	<b>0.019 ± 0.000</b>	0.020 ± 0.002
<b>CLoVE</b>	<b>0.014 ± 0.000</b>	0.021 ± 0.000	<b>0.016 ± 0.000</b>

486 tasks using autoencoders, further highlighting its versatility and effectiveness. More experiments in  
 487 an unsupervised setting can be found in App. D.2.  
 488

489 **Robustness to initialization** Our results indicate that one key  
 490 reason *CLoVE* outperforms related baselines such as *IFCA* is its  
 491 robustness to initialization. We evaluate this by varying two factors:  
 492 the initialization of model parameters and the first round client-to-  
 493 model assignment strategy. For model initialization, we consider  
 494 two settings: (1) all models begin with identical parameters (init:s  
 495 – same), and (2) each model is initialized independently using  
 496 PyTorch’s default random initialization (init:d – different). For  
 497 the first round client-to-model assignment, we also explore two  
 498 approaches: (1) assigning clients to models at random (first:r –  
 499 random), and (2) assigning clients by clustering their loss vectors  
 500 (first:e – evaluation-based), followed by bipartite matching as in  
 501 Alg. 1. Fig. 2 compares *CLoVE* and *IFCA* on unsupervised  
 502 MNIST image reconstruction with 10 clusters, each containing 5 clients with 500 samples. The  
 503 results show that *CLoVE* maintains high performance across different initialization methods, unlike  
 504 *IFCA*.  
 505

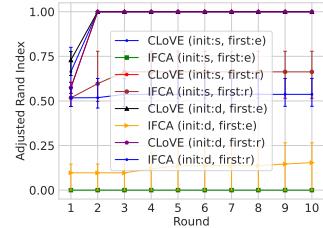


Figure 2: Clustering accuracy of *CLoVE* and *IFCA* over time for the initialization experiments

505 Table 5: Convergence behavior for MNIST, CIFAR-10 and FMNIST

Data mixing	Algorithm	ARI reached in 10 rounds			First round when ARI $\geq 0.9$		
		MNIST	CIFAR-10	FMNIST	MNIST	CIFAR-10	FMNIST
Label skew 1	CFL-S	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	—	41.3 $\pm$ 2.1	—
	FeSEM	0.18 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	—	—	—
	FlexCFL	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.0 <math>\pm</math> 0.0</b>	<b>1.0 <math>\pm</math> 0.0</b>	<b>1.0 <math>\pm</math> 0.0</b>
	PACFL	0.39 $\pm$ 0.07	0.60 $\pm$ 0.09	<b>1.00 <math>\pm</math> 0.00</b>	—	—	<b>1.0 <math>\pm</math> 0.0</b>
	IFCA	0.83 $\pm$ 0.12	0.72 $\pm$ 0.00	0.83 $\pm$ 0.12	—	—	—
	CLoVE	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	2.0 $\pm$ 0.0	2.0 $\pm$ 0.0	2.0 $\pm$ 0.0
Concept shift	CFL-S	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	60.0 $\pm$ 9.4	—	43.0 $\pm$ 5.7
	FeSEM	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	—	—	—
	FlexCFL	<b>1.00 <math>\pm</math> 0.00</b>	0.22 $\pm$ 0.24	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.0 <math>\pm</math> 0.0</b>	—	<b>1.0 <math>\pm</math> 0.0</b>
	PACFL	0.00 $\pm$ 0.02	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	—	—	—
	IFCA	0.68 $\pm$ 0.00	0.76 $\pm$ 0.17	0.55 $\pm$ 0.18	—	—	—
	CLoVE	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	2.0 $\pm$ 0.0	<b>2.3 <math>\pm</math> 0.5</b>	2.0 $\pm$ 0.0

## 521 6 CONCLUSION

522 We introduced *CLoVE*, a simple loss vector embeddings-based framework for personalized clustered  
 523 federated learning with low communication and computation overhead that avoids stringent model  
 524 initialization assumptions and substantially outperforms the state of the art across a range of datasets  
 525 and in a variety of both supervised and unsupervised settings. Further discussion of the overheads  
 526 and the privacy properties of *CLoVE* is provided in App. B. In the future, we plan to further explore  
 527 privacy, as well as fairness and adversarial behavior aspects. Overall, *CLoVE*’s design offers a  
 528 promising and robust approach to scalable model personalization and clustering under heterogeneous  
 529 data distributions.

## 531 7 REPRODUCIBILITY STATEMENT

532 Our code, together with detailed instructions on how to reproduce our experiments, will be open-  
 533 sourced. For the purposes of the reviewing process, after the discussion forums open, we will make  
 534 a comment directed to the reviewers and area chairs and put a link to an anonymous repository  
 535 containing our code, as mentioned in the Author Guide webpage of the conference (<https://iclr.cc/Conferences/2026/AuthorGuide>). Proofs for our theoretical claims are in App.  
 536 A. All datasets we are using are already public. We have provided detailed implementation details in  
 537 Section 5 of the main paper and Appendices C and D.

540 REFERENCES  
541

542 Pranjal Awasthi and Or Sheffet. Improved Spectral-Norm Bounds for Clustering. In *International*  
543 *Workshop on Approximation Algorithms for Combinatorial Optimization*, pp. 37–49. Springer,  
544 2012.

545 K Bonawitz, H Eichner, W Grieskamp, D Huba, A Ingerman, V Ivanov, C Kiddon, J Konecny,  
546 S Mazzocchi, HB McMahan, et al. Towards Federated Learning at Scale: System Design. In  
547 *Proceedings of the 2nd SysML Conference*, pp. 1–15, 2019.

548

549 Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H. Brendan  
550 McMahan, Virginia Smith, and Ameet Talwalkar. LEAF: A Benchmark for Federated Settings,  
551 December 2019. URL <http://arxiv.org/abs/1812.01097>.

552

553 Zihan Chen, Howard H. Yang, Tony Q.S. Quek, and Kai Fong Ernest Chong. Spectral Co-Distillation  
554 for Personalized Federated Learning. In *Proceedings of the 37th International Conference on*  
555 *Neural Information Processing Systems*, NIPS ’23, pp. 8757–8773, Red Hook, NY, USA, December  
556 2023. Curran Associates Inc.

557

558 Jichan Chung, Kangwook Lee, and Kannan Ramchandran. Federated Unsupervised Clustering with  
559 Generative Models. In *AAAI 2022 international workshop on trustable, verifiable and auditable*  
560 *federated learning*, volume 4, 2022.

561

562 Anirban Dasgupta, John Hopcroft, Ravi Kannan, and Pradipta Mitra. Spectral Clustering with  
563 Limited Independence. In *Proceedings of the eighteenth annual ACM-SIAM symposium on*  
564 *Discrete algorithms*, pp. 1036–1045. Citeseer, 2007.

565

566 Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Adaptive Personalized Federated  
567 Learning. *arXiv preprint arXiv:2003.13461*, 2020.

568

569 Don Kurian Dennis, Tian Li, and Virginia Smith. Heterogeneity for the Win: One-Shot Federated  
570 Clustering. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 2611–  
2620. PMLR, July 2021. URL <https://proceedings.mlr.press/v139/dennis21a.html>. ISSN: 2640-3498.

571

572 Moming Duan, Duo Liu, Xinyuan Ji, Yu Wu, Liang Liang, Xianzhang Chen, Yujuan Tan, and  
573 Ao Ren. Flexible Clustered Federated Learning for Client-Level Data Distribution Shift. *IEEE*  
574 *Transactions on Parallel and Distributed Systems*, 33(11):2661–2674, November 2022. ISSN  
575 1558-2183. doi: 10.1109/TPDS.2021.3134263. URL <https://ieeexplore.ieee.org/document/9647969>.

576

577 Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized Federated Learning with  
578 Theoretical Guarantees: A Model-Agnostic Meta-Learning Approach. In *Proceedings of the 34th*  
579 *International Conference on Neural Information Processing Systems*, NIPS ’20, pp. 3557–3568,  
580 Red Hook, NY, USA, December 2020. Curran Associates Inc. ISBN 978-1-7138-2954-6.

581

582 Avishek Ghosh, Jichan Chung, Dong Yin, and Kannan Ramchandran. An Efficient Framework for  
583 Clustered Federated Learning. In *Proceedings of the 34th International Conference on Neural*  
584 *Information Processing Systems*, NIPS ’20, pp. 19586–19597, Red Hook, NY, USA, December  
585 2020. Curran Associates Inc. ISBN 978-1-71382-954-6.

586

587 Avishek Ghosh, Jichan Chung, Dong Yin, and Kannan Ramchandran. An Efficient Framework  
588 for Clustered Federated Learning. *IEEE Transactions on Information Theory*, 68(12):8076–  
589 8091, December 2022. ISSN 1557-9654. doi: 10.1109/TIT.2022.3192506. URL <https://ieeexplore.ieee.org/document/9832954>.

590

591 Antonio Gulli. AG’s Corpus of News Articles. URL [http://groups.di.unipi.it/~gulli/AG\\_corpus\\_of\\_news\\_articles.html](http://groups.di.unipi.it/~gulli/AG_corpus_of_news_articles.html).

592

593 Filip Hanzely and Peter Richtárik. Federated Learning of a Mixture of Global and Local Models.  
arXiv preprint arXiv:2002.05516, 2020.

594 Yihan Jiang, Jakub Konečný, Keith Rush, and Sreeram Kannan. Improving Federated Learning  
 595 Personalization via Model Agnostic Meta Learning, 2023. URL <https://arxiv.org/abs/1909.12488>.  
 596

597 Roy Jonker and Anton Volgenant. A shortest augmenting path algorithm for dense and sparse linear  
 598 assignment problems. *Computing*, 38(4):325–340, 1987.  
 599

600 Heasung Kim, Hyeji Kim, and Gustavo De Veciana. Clustered Federated Learning via Gradient-based  
 601 Partitioning. In *Proceedings of the 41st International Conference on Machine Learning*, ICML’24.  
 602 JMLR.org, 2024.

603 Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization, January 2017.  
 604 URL <http://arxiv.org/abs/1412.6980>. arXiv:1412.6980 [cs].  
 605

606 Jakub Konečný, H. Brendan McMahan, Felix X. Yu, Peter Richtarik, Ananda Theertha Suresh, and  
 607 Dave Bacon. Federated Learning: Strategies for Improving Communication Efficiency. In *NIPS*  
 608 *Workshop on Private Multi-Party Machine Learning*, 2016. URL <https://arxiv.org/abs/1610.05492>.  
 609

610 Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images, 2009. URL <https://www.cs.toronto.edu/~kriz/cifar.html>.  
 611

612 Amit Kumar and Ravindran Kannan. Clustering with Spectral Norm and the k-means Algorithm.  
 613 In *2010 IEEE 51st Annual Symposium on Foundations of Computer Science*, pp. 299–308. IEEE,  
 614 2010.

615 Beatrice Laurent and Pascal Massart. Adaptive Estimation of a Quadratic Functional by Model  
 616 Selection. *Annals of statistics*, pp. 1302–1338, 2000.  
 617

618 Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based Learning Applied to Document  
 619 Recognition. *Proceedings of the IEEE*, 86(11):2278–2324, November 1998. ISSN 1558-2256.  
 620 doi: 10.1109/5.726791. URL <https://ieeexplore.ieee.org/document/726791>.  
 621 Conference Name: Proceedings of the IEEE.

622 Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the Convergence of  
 623 FedAvg on Non-IID Data. *arXiv preprint arXiv:1907.02189*, 2019.  
 624

625 Alessandro Licciardi, Davide Leo, Eros Fanì, Barbara Caputo, and Marco Ciccone. Interaction-  
 626 aware Gaussian weighting for clustered federated learning. In Aarti Singh, Maryam Fazel, Daniel  
 627 Hsu, Simon Lacoste-Julien, Felix Berkenkamp, Tegan Maharaj, Kiri Wagstaff, and Jerry Zhu  
 628 (eds.), *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of  
 629 *Proceedings of Machine Learning Research*, pp. 37642–37666. PMLR, 13–19 Jul 2025. URL  
 630 <https://proceedings.mlr.press/v267/licciardi25a.html>.

631 Guodong Long, Ming Xie, Tao Shen, Tianyi Zhou, Xianzhi Wang, and Jing Jiang. Multi-Center  
 632 Federated Learning: Clients Clustering for Better Personalization. *World Wide Web*, 26(1):  
 633 481–500, January 2023. ISSN 1573-1413. doi: 10.1007/s11280-022-01046-x. URL <https://doi.org/10.1007/s11280-022-01046-x>.  
 634

635 László Lovász and Michael D Plummer. *Matching Theory*, volume 367. American Mathematical  
 636 Society, 2009.

637 Yishay Mansour, Mehryar Mohri, Jae Ro, and Ananda Theertha Suresh. Three Approaches for  
 638 Personalization with Applications to Federated Learning, July 2020a. URL <http://arxiv.org/abs/2002.10619>.  
 639

640 Yishay Mansour, Mehryar Mohri, Jae Ro, and Ananda Theertha Suresh. Three Approaches for  
 641 Personalized Federated Learning. In *Advances in Neural Information Processing Systems*, pp.  
 642 1–13, 2020b.  
 643

644 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.  
 645 Communication-Efficient Learning of Deep Networks from Decentralized Data. In Aarti Singh  
 646 and Jerry Zhu (eds.), *Proceedings of the 20th International Conference on Artificial Intelligence  
 647 and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pp. 1273–1282. PMLR,  
 648 April 2017. URL <https://proceedings.mlr.press/v54/mcmahan17a.html>.

648 James Munkres. Algorithms for the assignment and transportation problems. *Journal of the Society*  
 649 *for Industrial and Applied Mathematics*, 5(1):32–38, 1957.

650

651 Jianmo Ni, Jiacheng Li, and Julian McAuley. Justifying Recommendations using Distantly-Labeled  
 652 Reviews and Fine-Grained Aspects. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan  
 653 (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*  
 654 *and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*,  
 655 pp. 188–197, Hong Kong, China, November 2019. Association for Computational Linguistics. doi:  
 656 10.18653/v1/D19-1018. URL <https://aclanthology.org/D19-1018/>.

657

658 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
 659 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward  
 660 Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang,  
 661 Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning  
 662 Library. In *Proceedings of the 33rd International Conference on Neural Information Processing  
 663 Systems*, number 721, pp. 8026–8037. Curran Associates Inc., Red Hook, NY, USA, December  
 2019.

664

665 Felix Sattler, Klaus-Robert Müller, and Wojciech Samek. Clustered Federated Learning: Model-  
 666 Agnostic Distributed Multitask Optimization Under Privacy Constraints. *IEEE Transactions on  
 667 Neural Networks and Learning Systems*, 32(8):3710–3722, August 2021. ISSN 2162-2388. doi:  
 668 10.1109/TNNLS.2020.3015958. Conference Name: IEEE Transactions on Neural Networks and  
 Learning Systems.

669

670 Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S Talwalkar. Federated Multi-Task  
 671 Learning. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and  
 672 R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Asso-  
 673 ciates, Inc., 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/6211080fa89981f66b1a0c9d55c61d0f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/6211080fa89981f66b1a0c9d55c61d0f-Paper.pdf).

674

675 Alysa Ziying Tan, Han Yu, Lizhen Cui, and Qiang Yang. Towards Personalized Federated Learning.  
 676 *IEEE Transactions on Neural Networks and Learning Systems*, 34(12):9587–9603, 2023. doi:  
 677 10.1109/TNNLS.2022.3160699.

678

679 Yue Tan, Guodong Long, Lu Liu, Tianyi Zhou, Qinghua Lu, Jing Jiang, and Chengqi Zhang.  
 680 FedProto: Federated Prototype Learning across Heterogeneous Clients. *Proceedings of the AAAI  
 681 Conference on Artificial Intelligence*, 36(8):8432–8440, June 2022. ISSN 2374-3468. doi: 10.  
 682 1609/aaai.v36i8.20819. URL <https://ojs.aaai.org/index.php/AAAI/article/view/20819>. Number: 8.

683

684 Dennis Treder-Tschechlov. k-Means Clustering on Image Data using the  
 685 MNIST Dataset, 2024. URL <https://medium.com/@tschechd/k-means-clustering-on-image-data-using-the-mnist-dataset-8101fcc650eb>.  
 686 Accessed: 2025-01-29.

687

688 Saeed Vahidian, Mahdi Morafah, Weijia Wang, Vyacheslav Kungurtsev, Chen Chen, Mubarak Shah,  
 689 and Bill Lin. Efficient Distribution Similarity Identification in Clustered Federated Learning via  
 690 Principal Angles between Client Data Subspaces. *Proceedings of the AAAI Conference on Artificial  
 691 Intelligence*, 37(8):10043–10052, June 2023. ISSN 2374-3468. doi: 10.1609/aaai.v37i8.26197.  
 692 URL <https://ojs.aaai.org/index.php/AAAI/article/view/26197>. Number:  
 693 8.

694

695 Harsh Vardhan, Avishek Ghosh, and Arya Mazumdar. An Improved Federated Clustering Algorithm  
 696 with Model-based Clustering. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856.  
 697 URL <https://openreview.net/forum?id=1ZGA5mSk0B>.

698

699 Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series  
 700 in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019.

701

Mariel Werner, Lie He, Michael Jordan, Martin Jaggi, and Sai Praneeth Karimireddy. Provably  
 Personalized and Robust Federated Learning. *Transactions on Machine Learning Research*, August  
 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=B0uBSSUy0G>.

702 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: A Novel Image Dataset for Bench-  
 703 marking Machine Learning Algorithms, September 2017. URL <http://arxiv.org/abs/1708.07747>. arXiv:1708.07747 [cs].  
 704

705 Jian Xu, Xinyi Tong, and Shao-Lun Huang. Personalized Federated Learning with Feature Alignment  
 706 and Classifier Collaboration. In *The Eleventh International Conference on Learning Representations*,  
 707 2023. URL <https://openreview.net/forum?id=SXZr8aDKia>.  
 708

709 Xinyang Yi, Constantine Caramanis, and Sujay Sanghavi. Alternating Minimization for Mixed Linear  
 710 Regression. In *International Conference on Machine Learning*, pp. 613–621. PMLR, 2014.  
 711

712 Jianqing Zhang, Yang Hua, Hao Wang, Tao Song, Zhengui Xue, Ruhui Ma, and Haibing Guan.  
 713 FedALA: Adaptive Local Aggregation for Personalized Federated Learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(9):11237–11244, June 2023. ISSN 2374-  
 714 3468. doi: 10.1609/aaai.v37i9.26330. URL <https://ojs.aaai.org/index.php/AAAI/article/view/26330>. Number: 9.  
 715

716 Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level Convolutional Networks for Text  
 717 Classification. In *Advances in Neural Information Processing Systems*, volume 28. Curran As-  
 718 sociates, Inc., 2015. URL [https://papers.nips.cc/paper\\_files/paper/2015/hash/250cf8b51c773f3f8dc8b4be867a9a02-Abstract.html](https://papers.nips.cc/paper_files/paper/2015/hash/250cf8b51c773f3f8dc8b4be867a9a02-Abstract.html).  
 719

720 Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, David Civin, and Vikas Chandra. Federated  
 721 Learning with Non-IID Data. In *arXiv preprint arXiv:1806.00582*, pp. 1–10, 2018.  
 722

723

## A THEORETICAL ANALYSIS

### A.1 BACKGROUND

724 We consider a mixed linear regression problem (Yi et al., 2014) in a federated setting (Ghosh et al.,  
 725 2022). Client’s data come from one of  $K$  different clusters, denoted as  $\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_K$ . Cluster  $k$ ’s  
 726 optimal model  $\theta_k^* \in \mathbb{R}^d$  is a vector of dimension  $d$ .  
 727

728 **Assumptions:** We assume all clients of a cluster  $\mathcal{C}_k$  use the same number  $n_k$  of independently drawn  
 729 feature-response pairs  $(x_i, y_i)$  in each federated round. These feature-response pairs are drawn from  
 730 a distribution  $\mathcal{D}_k$ , generated by the model:  
 731

$$735 \quad y_i = \langle \theta_k^*, x_i \rangle + \epsilon_i,$$

736 where  $x_i \sim \mathcal{N}(0, I_d)$  are the features, and  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$  represents additive noise. Both are  
 737 independently drawn. We assume  $\sigma$  is very small. Specifically,  $\sigma \ll 1$ .  
 738

739 We use notation  $[K]$  for the set  $\{1, 2, \dots, K\}$ .  
 740

The loss function  $f(\theta)$  is defined as the square of the error:  
 741

$$742 \quad f(\theta; x, y) = (\langle \theta, x \rangle - y)^2.$$

743 The **population loss** for cluster  $k$  is:  
 744

$$745 \quad \mathcal{L}_k(\theta) = \mathbb{E}_{(x, y) \sim \mathcal{D}_k} \left[ (\langle \theta, x \rangle - y)^2 \right],$$

746 Note that the optimal parameters  $\{\theta_k^*\}_{k=1}^K$  are the minimizers of the population losses  $\{\mathcal{L}_k\}_{k=1}^K$ , i.e.,  
 747

$$748 \quad \theta_k^* = \arg \min_{\theta} \mathcal{L}_k(\theta) \quad \text{for each } k \in [K].$$

749 We assume these optimal models have unit norm. That is  $\|\theta_k^*\| = 1$  for all  $k$ . We also assume the  
 750 optimal models are well-separated. That is, there exists a  $\Delta > 0$  such that for any pair  $i \neq j$ :  
 751

$$752 \quad \|\theta_i^* - \theta_j^*\| \geq \Delta$$

753 For a given  $\theta$ , the **empirical loss** for client  $i$  in cluster  $k$  is:  
 754

$$755 \quad \mathcal{L}_k^i(\theta) = \frac{1}{n_k} \sum_{j=1}^{n_k} (\langle \theta, x_j \rangle - y_j)^2 = \frac{1}{n_k} \sum_{j=1}^{n_k} (\langle \theta_k^* - \theta, x_j \rangle + \epsilon_j)^2$$

756 Let

757 
$$\alpha(\theta, \theta_k^*) = \sqrt{\|\theta - \theta_k^*\|^2 + \sigma^2}.$$
 758

759 Note  $\mathcal{L}_k^i(\theta) \sim \frac{\alpha(\theta, \theta_k^*)^2}{n_k} \chi^2(n_k)$ . That is, it is distributed as a scaled chi-squared random variable with  
760  $n_k$  degrees of freedom, with scaling factor  $\frac{\alpha(\theta, \theta_k^*)^2}{n_k}$ . In the following, we denote the square root of  
761 the empirical loss by  $F_k^i(\theta)$ . That is,  $F_k^i(\theta) = \sqrt{\mathcal{L}_k^i(\theta)}$ . Note:

762 
$$F_k^i(\theta) \sim \frac{\alpha(\theta, \theta_k^*)}{\sqrt{n_k}} \chi(n_k),$$
 763  
764

765 where  $\chi(n_k)$  denotes a chi-distributed random variable with  $n_k$  degrees of freedom.766 Using Sterling approximation,  $\mathbb{E}[\chi(n)] = \sqrt{n-1} [1 - \frac{1}{4n} + O(\frac{1}{n^2})]$ . Also,  $\text{Var}[\chi(n)] = (n-1) - \mathbb{E}[\chi(n)]^2 = \frac{n-1}{2n} [1 + O(\frac{1}{n})]$ . Therefore:

767 
$$\mathbb{E}[F_k^i(\theta)] = \alpha \left[ 1 - \frac{1}{4n_k} + O\left(\frac{1}{n_k^2}\right) \right] \approx \alpha(\theta, \theta_k^*),$$
 768  
769

770 
$$\text{Var}(F_k^i(\theta)) = \frac{\alpha^2}{n_k} \left( \frac{n_k - 1}{2n_k} + O\left(\frac{1}{n_k}\right) \right) \approx \frac{\alpha(\theta, \theta_k^*)^2}{2n_k}.$$
 771  
772

773 For  $N$  models  $\boldsymbol{\theta} = [\theta_1 \ \theta_2 \ \dots \ \theta_N]$ , the square root loss vector (SLV) of a client  $i$  of cluster  $k$   
774 is  $a_k^i(\boldsymbol{\theta}) = [F_k^i(\theta_1), F_k^i(\theta_2), \dots, F_k^i(\theta_N)]$ . Let  $c_k(\boldsymbol{\theta}) = [c_k^1, c_k^2, \dots, c_k^N]$  be the mean of  $k$ -th cluster's  
775 SLVs. Note that  $c_k^j = \mathbb{E}[F_k^i(\theta_j)] = \alpha(\theta_j, \theta_k^*) = \sqrt{\|\theta_j - \theta_k^*\|^2 + \sigma^2}$ .776 All the model collections  $\boldsymbol{\theta}$  considered in this work are assumed to be one of two different types: a)  
777  $N = d$  ortho-normal models, such as  $d$  models drawn randomly from a  $d$ -dimensional unit sphere,  
778 since such vectors are likely to be orthogonal, and b)  $N = K$ ,  $\Delta$ -proximal models in which each  $\theta_k$   
779 is within a distance at most  $\Delta/4$  of its optimal counterpart  $\theta_k^*$ . Furthermore, the norm of any model  
780 in these model collections is required to be bounded. Specifically,  $\|\theta\| \leq 2$  for any model  $\theta \in \boldsymbol{\theta}$ .781 

## A.2 RESULTS

782 We prove that the mean of different cluster SLVs are well-separated for our choice of model  
783 collections. We first show this for a  $N = d$  ortho-normal  $\boldsymbol{\theta}$  model collection.784 **Lemma A.1.** *For ortho-normal models  $\boldsymbol{\theta}$ , for any  $i \neq j$ ,  $\|c_i(\boldsymbol{\theta}) - c_j(\boldsymbol{\theta})\| \geq \frac{\Delta}{c}$ , where  $c \approx 2$ .*785 *Proof.* First, we show the result for the following ortho-normal models:

786 
$$\theta_i = (0, 0, \dots, 1, \dots, 0) \quad \text{with the 1 in the } i\text{-th position.}$$
 787

788 Note,

789 
$$|\|\theta_i^* - \theta_k\|^2 - \|\theta_j^* - \theta_k\|^2| = 2 |\langle (\theta_i^* - \theta_j^*), \theta_k \rangle| = 2 |t_k|,$$

790 where  $t_k$  is the  $k$ -th component of  $\theta_i^* - \theta_j^*$ . Thus

791 
$$|(c_k^i)^2 - (c_k^j)^2| = |\|\theta_i^* - \theta_k\|^2 - \|\theta_j^* - \theta_k\|^2| = 2 |t_k|.$$
 792

793 Also, since  $\forall_i, \|\theta_i^*\| = 1$  and  $\forall_k, \|\theta_k\| = 1$ , it follows for any  $k, i$  that  $c_k^i = \sqrt{\|\theta_k - \theta_i^*\|^2 + \sigma^2} \leq$   
794  $\sqrt{4 + \sigma^2}$ . Thus,

795 
$$|c_k^i - c_k^j| \geq \frac{|t_k|}{\sqrt{4 + \sigma^2}}$$
 796

797 Thus, since  $\sum_k t_k^2 = \|\theta_i^* - \theta_j^*\|^2$ :

798 
$$\|c_i(\boldsymbol{\theta}) - c_j(\boldsymbol{\theta})\|^2 \geq \sum_k \left( \frac{|t_k|}{\sqrt{4 + \sigma^2}} \right)^2 \geq \frac{\|\theta_i^* - \theta_j^*\|^2}{4 + \sigma^2}$$
 799

810 However, since  $\forall i \neq j$ ,  $\|\theta_i^* - \theta_j^*\| \geq \Delta$ , it follows:  
 811

$$812 \quad 813 \quad \|c_i(\boldsymbol{\theta}) - c_j(\boldsymbol{\theta})\|^2 \geq \frac{\Delta^2}{4 + \sigma^2} \\ 814$$

815 since, by assumption,  $\sigma \ll 1$ . Thus, for  $c \approx 2$ ,

$$816 \quad 817 \quad \forall i \neq j, \|c_i(\boldsymbol{\theta}) - c_j(\boldsymbol{\theta})\| \geq \frac{\Delta}{c}. \quad (1) \\ 818$$

819 It is easy to see that the proof holds for any set  $\boldsymbol{\theta}$  of ortho-normal unit vectors. Specifically, since  
 820 vectors drawn randomly from a  $d$ -dimensional unit sphere are nearly orthogonal, for large  $d$ , the  
 821 result holds for them as well.  $\square$

822 We now prove an analogue of the Lemma A.1 for the  $K$  models of  $\boldsymbol{\theta}$  that satisfy the proximity  
 823 condition.

824 **Lemma A.2.** *Let  $\boldsymbol{\theta}$  be collection of  $K$  models that satisfy the  $\Delta$ -proximity condition. Then, for any  
 825  $i \neq j$ ,  $\|c_i(\boldsymbol{\theta}) - c_j(\boldsymbol{\theta})\| > \frac{\Delta}{2}$ .*

826 *Proof.* Recall,  $c_k(\boldsymbol{\theta}) = [c_k^1, c_k^2, \dots, c_k^N]$ , where  $c_k^j = \mathbb{E} [F_k^j(\theta_j)] = \alpha(\theta_j, \theta_k^*) = \sqrt{\|\theta - \theta_k^*\|^2 + \sigma^2}$ .

827 Since  $\|\theta_i - \theta_i^*\| \leq \frac{\Delta}{4}$  and since  $\|\theta_i^* - \theta_j^*\| \geq \Delta$ , by triangle inequality it follows that  $\|\theta_i - \theta_j^*\| \geq \frac{3\Delta}{4}$ .  
 828 Likewise,  $\|\theta_j - \theta_j^*\| \leq \frac{\Delta}{4}$  and  $\|\theta_j - \theta_i^*\| \geq \frac{3\Delta}{4}$ . Thus, since  $\sigma \ll \frac{\Delta}{4}$ :

$$829 \quad 830 \quad |c_j^i - c_i^i| = \left| \sqrt{\|\theta_i - \theta_j^*\|^2 + \sigma^2} - \sqrt{\|\theta_i - \theta_i^*\|^2 + \sigma^2} \right| \\ 831 \quad 832 \quad \geq \left| \|\theta_i - \theta_j^*\|^2 - \sqrt{\|\theta_i - \theta_i^*\|^2 + \left(\frac{\Delta}{4}\right)^2} \right| \geq \frac{3 - \sqrt{2}}{4} \Delta \\ 833 \quad 834$$

835 Likewise,

$$836 \quad |c_i^j - c_j^j| = \left| \sqrt{\|\theta_j - \theta_i^*\|^2 + \sigma^2} - \sqrt{\|\theta_j - \theta_j^*\|^2 + \sigma^2} \right| \geq \frac{3 - \sqrt{2}}{4} \Delta \\ 837 \quad 838$$

839 Since:

$$840 \quad \|c_i(\boldsymbol{\theta}) - c_j(\boldsymbol{\theta})\|^2 \geq |c_j^i - c_i^i|^2 + |c_i^j - c_j^j|^2,$$

841 it follows:

$$842 \quad \|c_i(\boldsymbol{\theta}) - c_j(\boldsymbol{\theta})\|^2 \geq 2 \left( \frac{3 - \sqrt{2}}{4} \right)^2 \Delta^2. \\ 843 \quad 844$$

845 Thus,

$$846 \quad \|c_i(\boldsymbol{\theta}) - c_j(\boldsymbol{\theta})\| > \frac{\Delta}{2} \\ 847 \quad 848$$

$\square$

849 Any ortho-normal model, by definition, satisfies  $\|\theta\| = 1 < 2$ . Now we show the same for any  
 850  $\Delta$ -proximal model.

851 **Lemma A.3.**  $\|\theta\| < 2$ , for any  $\Delta$ -proximal model  $\theta$ .

852 *Proof.* This follows from our assumption  $\|\theta_k^*\| = 1$  for all  $k \in [K]$  and since there exists a  $k \in [K]$   
 853 such that  $\|\theta - \theta_k^*\| \leq \Delta/4$ , where  $\Delta < 2$ . Therefore, by triangle inequality,  $\|\theta\| < 2$ .  $\square$

864 Since, for all models,  $\|\theta\| < 2$ , this implies  $\|\theta - \theta_k^*\| \leq 3$ , for all  $k \in [K]$ . Thus, it follows, for all  
 865  $k \in [K]$ :

$$866 \quad 867 \quad \frac{1}{2n_k} (\|\theta - \theta_k^*\|^2 + \sigma^2) \leq \frac{1}{2n_k} (9 + \sigma^2) \quad (2)$$

868 All our results below hold for any ortho-normal or  $\Delta$ -proximal models collection  $\theta$ , unless noted  
 869 otherwise.  
 870

871 The Lemma below establishes that SLVs of a cluster are concentrated around their means and this  
 872 result holds with high probability.  
 873

874 **Lemma A.4.** *Let client  $i$  be in cluster  $s(i)$ . Let  $\text{Var}_j(a_{s(i)}^i)$  be the variance of the  $j$ -th entry of its  
 875 SLV  $a_{s(i)}^i(\theta)$ . Then, for  $t = \frac{9}{2} \log^2 \frac{NM}{\delta}$  and for any error tolerance  $\delta < 1$ :*

$$876 \quad 877 \quad \mathbb{P} \left[ \forall i, \|a_{s(i)}^i(\theta) - c_{s(i)}(\theta)\|^2 \leq t \max_j \text{Var}_j(a_{s(i)}^i) \right] \geq 1 - \delta.$$

880 *Proof.* In what follows, we use  $k$  to denote the cluster assignment  $s(i)$  for client  $i$ . Let  
 881

$$882 \quad a_k^i(\theta) = [F_k^i(\theta_1), F_k^i(\theta_2), \dots, F_k^i(\theta_N)],$$

883 and therefore

$$884 \quad c_k(\theta) = [\mathbb{E}[F_k^i(\theta_1)], \mathbb{E}[F_k^i(\theta_2)], \dots, \mathbb{E}[F_k^i(\theta_N)]].$$

886 We start by focusing on the  $j$ -th component  $F_k^i(\theta_j)$  of  $a_k^i(\theta)$ . For notational convenience, we will  
 887 omit the index  $j$  and refer to this component as  $F_k^i(\theta)$  in the following proof.  
 888

889 Recall

$$890 \quad 891 \quad \mathcal{L}_k^i(\theta) = (F_k^i(\theta))^2 = \frac{1}{n_k} \sum_{j=1}^{n_k} (\langle \theta_k^* - \theta, x_j \rangle + \epsilon_j)^2.$$

892  $\mathcal{L}_k^i(\theta)$  is scaled chi-squared random variable with  $n_k$  degrees of freedom (i.e.  $\sim \frac{\alpha(\theta, \theta_k^*)^2}{n_k} \chi^2(n_k)$ ).  
 893 Thus, its mean and variance are:

$$894 \quad \mathbb{E} [\mathcal{L}_k^i(\theta)] = \alpha(\theta, \theta_k^*)^2$$

$$895 \quad \text{Var}(\mathcal{L}_k^i(\theta)) = \frac{2}{n_k} \alpha(\theta, \theta_k^*)^4$$

896 Note

$$897 \quad 898 \quad \mathcal{L}_k^i(\theta) \sim \sum_j \frac{1}{n_k} \alpha(\theta, \theta_k^*)^2 y_j^2$$

900 where  $y_j \sim \mathcal{N}(0, 1)$ . Thus,

$$901 \quad 902 \quad \mathcal{L}_k^i(\theta) \sim \sum_j \frac{\mathbb{E} [\mathcal{L}_k^i(\theta)]}{n_k} y_j^2$$

903 Let  $\gamma = \sqrt{\sum_j \left( \frac{\mathbb{E} [\mathcal{L}_k^i(\theta)]}{n_k} \right)^2} = \frac{1}{\sqrt{n_k}} \mathbb{E} [\mathcal{L}_k^i(\theta)]$ . Let  $\alpha = \frac{\mathbb{E} [\mathcal{L}_k^i(\theta)]}{n_k}$ . We apply Lemma 1, Section 4  
 904 of Laurent & Massart (2000) to get:  
 905

$$906 \quad \mathbb{P} [\mathcal{L}_k^i(\theta) - \mathbb{E} [\mathcal{L}_k^i(\theta)] \geq 2\alpha\sqrt{x} + 2\gamma x] \leq e^{-x}$$

907 Here

$$908 \quad 909 \quad 2\alpha\sqrt{x} + 2\gamma x = \mathbb{E} [\mathcal{L}_k^i(\theta)] \left( \frac{2\sqrt{x}}{n_k} + \frac{2x}{\sqrt{n_k}} \right)$$

910 Note that  $\frac{3x}{\sqrt{n_k}} > \frac{2\sqrt{x}}{n_k} + \frac{2x}{\sqrt{n_k}}$  for  $x > 1$  (specifically for any  $x > \frac{4}{n_k}$ ). Thus:  
 911

$$912 \quad 913 \quad \mathbb{P} \left[ \mathcal{L}_k^i(\theta) - \mathbb{E} [\mathcal{L}_k^i(\theta)] \geq \mathbb{E} [\mathcal{L}_k^i(\theta)] \frac{3x}{\sqrt{n_k}} \right] \leq e^{-x}$$

918 With  $3x = 2(1 + \beta)$ :

$$919 \quad \mathbb{P} \left[ \mathcal{L}_k^i(\theta) \geq \mathbb{E} [\mathcal{L}_k^i(\theta)] \left( 1 + 2(1 + \beta) \frac{1}{\sqrt{n_k}} \right) \right] \leq e^{-\frac{2}{3}(1+\beta)}$$

920 Recall  $\mathcal{L}_k^i(\theta) = (F_k^i(\theta))^2$  and

$$921 \quad \mathbb{E} [\mathcal{L}_k^i(\theta)] = \alpha(\theta, \theta_k^*)^2 = \mathbb{E} [F_k^i(\theta)]^2$$

922 Thus:

$$923 \quad \mathbb{P} \left[ F_k^i(\theta)^2 \geq \mathbb{E} [F_k^i(\theta)]^2 \left( 1 + 2(1 + \beta) \frac{1}{\sqrt{n_k}} \right) \right] \leq e^{-\frac{2}{3}(1+\beta)}$$

924 Since

$$925 \quad 1 + 2(1 + \beta) \frac{1}{\sqrt{n_k}} \leq \left( 1 + (1 + \beta) \frac{1}{\sqrt{n_k}} \right)^2,$$

$$926 \quad \mathbb{P} \left[ F_k^i(\theta) \geq \mathbb{E} [F_k^i(\theta)] \left( 1 + (1 + \beta) \frac{1}{\sqrt{n_k}} \right) \right] \leq e^{-\frac{2}{3}(1+\beta)}$$

927 or

$$928 \quad \mathbb{P} \left[ (F_k^i(\theta) - \mathbb{E} [F_k^i(\theta)])^2 \geq (1 + \beta)^2 \frac{\mathbb{E} [F_k^i(\theta)]^2}{n_k} \right] \leq e^{-\frac{2}{3}(1+\beta)}$$

929 Since

$$930 \quad \text{Var}(F_k^i(\theta)) = \frac{\alpha(\theta, \theta_k^*)^2}{2n_k} = \frac{\mathbb{E} [F_k^i(\theta)]^2}{2n_k},$$

$$931 \quad \mathbb{P} \left[ (F_k^i(\theta) - \mathbb{E} [F_k^i(\theta)])^2 \geq 2(1 + \beta)^2 \text{Var}(F_k^i(\theta)) \right] \leq e^{-\frac{2}{3}(1+\beta)}$$

932 We set  $t = 2(1 + \beta)^2$ . Then:

$$933 \quad \mathbb{P} \left[ (F_k^i(\theta) - \mathbb{E} [F_k^i(\theta)])^2 \geq t \text{Var}(F_k^i(\theta)) \right] \leq e^{-\frac{\sqrt{2t}}{3}}$$

934 Note that this result holds for a single component  $F_k^i(\theta_j)$  of  $A_i(\theta)$ . By using union bound over all of its  $N$  components:

$$935 \quad \mathbb{P} \left[ \sum_{j=1}^N (F_k^i(\theta_j) - \mathbb{E} [F_k^i(\theta_j)])^2 \geq t \max_j \text{Var}(F_k^i(\theta_j)) \right] \leq N e^{-\frac{\sqrt{2t}}{3}}$$

936 or

$$937 \quad \mathbb{P} \left[ \|a_k^i(\theta) - c_k(\theta)\|^2 \geq t \max_j \text{Var}(F_k^i(\theta_j)) \right] \leq N e^{-\frac{\sqrt{2t}}{3}}.$$

938 Since  $k = s(i)$  and  $\text{Var}(F_k^i(\theta_j)) = \text{Var}_j(a_{s(i)}^i)$ :

$$939 \quad \mathbb{P} \left[ \|a_{s(i)}^i(\theta) - c_{s(i)}(\theta)\|^2 \geq t \max_j \text{Var}_j(a_{s(i)}^i) \right] \leq N e^{-\frac{\sqrt{2t}}{3}}.$$

940 Applying union bounds for all  $M$  clients and by noticing that for  $t = \frac{9}{2} \log^2 \frac{NM}{\delta}$ :

$$941 \quad MN e^{-\frac{\sqrt{2t}}{3}} = \delta$$

942 we get

$$943 \quad \mathbb{P} \left[ \forall i, \|a_{s(i)}^i(\theta) - c_{s(i)}(\theta)\|^2 \geq t \max_j \text{Var}_j(a_{s(i)}^i) \right] \leq \delta$$

944 or

$$945 \quad \mathbb{P} \left[ \forall i, \|a_{s(i)}^i(\theta) - c_{s(i)}(\theta)\|^2 \leq t \max_j \text{Var}_j(a_{s(i)}^i) \right] \geq 1 - \delta.$$

946  $\square$

972 For models  $\theta$ , Let  $A(\theta)$  be the  $M \times N$  matrix of the SLVs of all  $M$  clients and let  $C(\theta)$  be the  
 973 corresponding  $M \times N$  matrix of SLV means. Thus,  $X(\theta) = A(\theta) - C(\theta)$  is the matrix of centered  
 974 SLVs for all clients. Let  $X_i(\theta) = A_i(\theta) - C_i(\theta)$  be the  $i$ -th row of  $A(\theta) - C(\theta)$ , with  $A_i(\theta)$  being  
 975 the SLV of the  $i$ -th client and  $C_i(\theta)$  the mean of that SLV.

976 We now establish a bound on the directional variance of the rows of  $A(\theta) - C(\theta)$ . First we show the  
 977 following:  
 978

979 **Lemma A.5.** *Let  $X$  be an  $N$ -dimensional vector with components  $x_i$  such that  $\mathbb{E}[x_i] = 0$  for all  $i$ ,  
 980 and let  $v$  be any  $N$ -dimensional vector. Then*

$$982 \mathbb{E} [\langle X, v \rangle^2] \leq \left( \sum_{i=1}^N |v_i| \sqrt{\text{Var}(x_i)} \right)^2. \quad (3)$$

985 *Proof.* Note

$$987 \mathbb{E} [\langle X, v \rangle^2] = \sum_{i=1}^N \sum_{j=1}^N v_i v_j \mathbb{E}[x_i x_j] \leq \sum_{i=1}^N \sum_{j=1}^N |v_i v_j| \mathbb{E}[|x_i x_j|].$$

989 Splitting the sum into diagonal and off-diagonal terms:

$$991 \mathbb{E} [\langle X, v \rangle^2] \leq \sum_{i=1}^N v_i^2 \mathbb{E}[x_i^2] + \sum_{i \neq j} |v_i v_j| \mathbb{E}[|x_i x_j|].$$

994 Applying the Cauchy-Schwarz inequality to the off-diagonal terms:

$$995 |\mathbb{E}[x_i x_j]| \leq \sqrt{\mathbb{E}[x_i^2] \mathbb{E}[x_j^2]} \quad \text{for } i \neq j.$$

997 Thus,

$$999 \mathbb{E} [\langle X, v \rangle^2] \leq \sum_{i=1}^N v_i^2 \mathbb{E}[x_i^2] + \sum_{i \neq j} |v_i v_j| \sqrt{\mathbb{E}[x_i^2] \mathbb{E}[x_j^2]}.$$

1001 Thus,

$$1002 \mathbb{E} [\langle X, v \rangle^2] \leq \left( \sum_{i=1}^N |v_i| \sqrt{\mathbb{E}[x_i^2]} \right)^2.$$

1005 We substitute  $\mathbb{E}[x_i^2] = \text{Var}(x_i)$  (since  $\mathbb{E}[x_i] = 0$ ) to express the bound in terms of variances:

$$1007 \mathbb{E} [\langle X, v \rangle^2] \leq \left( \sum_{i=1}^N |v_i| \sqrt{\text{Var}(x_i)} \right)^2.$$

1010  $\square$

1011 **Lemma A.6.** *Let  $X_i(\theta) = A_i(\theta) - C_i(\theta)$  be the  $i$ -th row of  $A(\theta) - C(\theta)$ , with  $A_i(\theta)$  being the  
 1012 SLV of  $i$ -th client and  $C_i(\theta)$  the mean of that SLV. Let  $i$ -th client be in cluster  $k$ .*

$$1014 \max_{v: \|v\|=1} \mathbb{E} [\langle X_i(\theta), v \rangle^2] \leq N \frac{1}{2n_k} (9 + \sigma^2) \quad (4)$$

1016 *Proof.* Applying Equation 3 of Lemma A.5, it follows:

$$1018 \max_{v: \|v\|=1} \mathbb{E} [\langle X_i(\theta), v \rangle^2] \leq \max_{v: \|v\|=1} \left( \sum_{j=1}^N |v_j| \sqrt{\text{Var}(X_i(\theta_j))} \right)^2,$$

1021 where  $X_i(\theta_j)$  is the  $j$ -th component of  $X_i(\theta)$ , which is the centered square root loss of client  $i$  on  
 1022 model  $\theta_j$ . That is  $X_i(\theta_j) = A_i(\theta_j) - \mathbb{E}[A_i(\theta_j)]$ . Therefore,  $\text{Var}(X_i(\theta_j)) = \text{Var}(A_i(\theta_j))$ . Hence:

$$1024 \max_{v: \|v\|=1} \mathbb{E} [\langle X_i(\theta), v \rangle^2] \leq \max_{v: \|v\|=1} \left( \sum_{j=1}^N |v_j| \sqrt{\text{Var}(A_i(\theta_j))} \right)^2$$

1026 Combining with Equation 2:

1028 
$$\text{Var}(A_i(\theta_j)) = \frac{1}{2n_k} (\|\theta_j - \theta_k^*\|^2 + \sigma^2) \leq \frac{1}{2n_k} (9 + \sigma^2).$$

1029 Thus

1030 
$$\max_{v: \|v\|=1} \mathbb{E} [\langle X_i(\theta), v \rangle^2] \leq N \frac{1}{2n_k} (9 + \sigma^2)$$

1031 since maximum happens when all  $v_j = \frac{1}{\sqrt{N}}$ . □

1032 Using a similar argument, it follows:

1033 
$$\mathbb{E} [\|X_i(\theta)\|^2] \leq N \frac{1}{2n_k} (9 + \sigma^2) \tag{5}$$

1034 We now establish a bound on the matrix operator norm  $\|X(\theta)\| = \|A(\theta) - C(\theta)\|$ . Let  $n = \min_k n_k$   
 1035 be the smallest number of data points any client has. Let  $\gamma = \sqrt{N \frac{1}{2n} (9 + \sigma^2)}$ . Then we can show  
 1036 the following bound:

1037 **Lemma A.7.** *With probability at least  $1 - \frac{1}{\text{polylog}(M)}$ :*

1038 
$$\|X(\theta)\| \leq \gamma \sqrt{M} \text{polylog}(M) \tag{6}$$

1039 *Proof.* With abuse of notation, we will use  $X_i$  to refer to  $i$ -th row of  $\mathbf{X}(\theta)$  and use  $X_{ij}$  to denote the  
 1040  $j$ -th entry of this row. Since  $\text{Var}(\|X_i\|) \leq E[\|X_i\|^2]$ , then by applying Equation 5 and by definition  
 1041 of  $\gamma$ , it follows  $\text{Var}(\|X_i\|) \leq E[\|X_i\|^2] \leq \gamma^2$ . Thus, by Chebyshev's inequality:

1042 
$$\mathbb{P} [\|X_i\| \geq \gamma \sqrt{M} \text{polylog}(M)] \leq \frac{\gamma^2}{M \gamma^2 \text{polylog}(M)}$$

1043 Thus by union bound

1044 
$$\mathbb{P} \left[ \max_{i=1, \dots, M} \|X_i\| \geq \gamma \sqrt{M} \text{polylog}(M) \right] \leq \frac{1}{\text{polylog}(M)}$$

1045 Thus with probability at least  $1 - \frac{1}{\text{polylog}(M)}$ :

1046 
$$\max_{i=1, \dots, M} \|X_i\| < \gamma \sqrt{M} \text{polylog}(M)$$

1047 Now for any unit column vector  $v$

1048 
$$v^T E [\mathbf{X}(\theta)^T \mathbf{X}(\theta)] v = \sum_{i=1}^M E [v^T X_i^T X_i v] = \sum_{i=1}^M E [\|X_i v\|^2] = \sum_{i=1}^M E [\langle X_i, v \rangle^2]$$

1049 Thus from Equation 4 of Lemma A.6 and by definition of  $\gamma$ , it follows:

1050 
$$\max_{v: \|v\|=1} v^T E [\mathbf{X}(\theta)^T \mathbf{X}(\theta)] v \leq \gamma^2 M$$

1051 Since  $E [\mathbf{X}(\theta)^T \mathbf{X}(\theta)]$  is symmetric, its spectral norm equals its largest eigenvalue, which is what is  
 1052 bounded above. Thus

1053 
$$\|E [\mathbf{X}(\theta)^T \mathbf{X}(\theta)]\| \leq \gamma^2 M$$

1054 We now use the following fact (Fact 6.1) from Kumar & Kannan (2010) that follows from the result  
 1055 of Dasgupta et al. (2007).

1056 **Lemma A.8.** *Let  $Y$  be a  $n \times d$  matrix,  $n \geq d$ , whose rows are chosen independently. Let  $Y_i$   
 1057 denote its  $i$ -th row. Let there exist  $\gamma$  such that  $\max_i \|Y_i\| \leq \gamma \sqrt{n}$  and  $\|E [Y^T Y]\| \leq \gamma^2 n$ . Then  
 1058  $\|Y\| \leq \gamma \sqrt{n} \text{polylog}(n)$ .*

1080 Thus, it follows that with probability at least  $1 - \frac{1}{\text{polylog}(M)}$ :

$$1083 \|\mathbf{X}(\boldsymbol{\theta})\| \leq \gamma \sqrt{M} \text{polylog}(M)$$

□

1086 We now establish one of our main results, that the SLVs satisfy the proximity condition of Kumar &  
1087 Kannan (2010), thereby proving that k-means with suitably initialized centers recovers the accurate  
1088 clustering of the SLVs (i.e. it recovers an accurate clustering of the clients).

1089 Recall that for client  $i$  of cluster  $k$ , its SLV is  $a_k^i(\boldsymbol{\theta}) = [F_k^i(\theta_1), F_k^i(\theta_2), \dots, F_k^i(\theta_N)]$ . Also  $c_k(\boldsymbol{\theta}) =$   
1090  $[c_1^k, c_2^k, \dots, c_N^k]$ , is the mean of  $k$ -th clusters SLVs, where  $c_k^j = \mathbb{E}[F_k^i(\theta_j)]$ . Recall  $X(\boldsymbol{\theta}) =$   
1091  $A(\boldsymbol{\theta}) - C(\boldsymbol{\theta})$  is the matrix of centered SLVs for all clients. Also,  $X_i(\boldsymbol{\theta}) = A_i(\boldsymbol{\theta}) - C_i(\boldsymbol{\theta})$  denotes  
1092 the  $i$ -th row of  $A(\boldsymbol{\theta}) - C(\boldsymbol{\theta})$ . Let  $m_k$  denote the number of clients in cluster  $k$ .

1093 **Theorem A.9.** *When  $nM$ , which is a lower bound on the the number of data points across all clients,  
1094 is much larger in comparison to  $d, K$ , and specifically when*

$$1096 dK^4 \text{polylog}(M) = o(nM), \quad (7)$$

1097 then for all ortho-normal or proximal models  $\boldsymbol{\theta}$ , the following proximity condition of Kumar &  
1098 Kannan (2010) holds for the SLV of any client  $i$ , with probability at least  $1 - \delta - \frac{1}{\text{polylog}(M)}$ , for  
1099 any error tolerance  $\delta$  :

$$1101 \forall j' \neq j, \|a_j^i(\boldsymbol{\theta}) - c_{j'}(\boldsymbol{\theta})\| - \|a_j^i(\boldsymbol{\theta}) - c_j(\boldsymbol{\theta})\| \geq \left( \frac{c'K}{m_i} + \frac{c'K}{m_j} \right) \|A(\boldsymbol{\theta}) - C(\boldsymbol{\theta})\|$$

1104 for some large constant  $c'$ . Here it is assumed that client  $i$  belongs to cluster  $j$ .

1105 *Proof.* Since  $\text{Var}(A_i(\theta_j)) = \frac{1}{2n_k} (\|\theta_i - \theta_k^*\|^2 + \sigma^2)$ , it follows that  $\forall i, j, \text{Var}(A_i(\theta_j)) \leq$   
1106  $\frac{1}{2n} (9 + \sigma^2) = \gamma^2/N$ , where  $n = \min_k n_k$ . Combining this with Equation A.4 established in  
1107 Lemma A.4, it follows with high probability (at least  $1 - \delta$ ):

$$1109 \forall i, \|A_i(\boldsymbol{\theta}) - C_i(\boldsymbol{\theta})\|^2 \leq \frac{9}{2} \log^2 \frac{NM}{\delta} \max_j \text{Var}(A_i(\theta_j)) \leq \frac{9}{2} \log^2 \frac{NM}{\delta} \frac{\gamma^2}{N},$$

1112 where  $A_i(\boldsymbol{\theta})$  and  $C_i(\boldsymbol{\theta})$  are the  $i$ -th rows of matrices  $A(\boldsymbol{\theta})$  and  $C(\boldsymbol{\theta})$  respectively.

1113 Without loss of generality, let the SLV of client  $i$  belonging to cluster  $j$  be in the  $i$ -th row of the  
1114 matrix  $A(\boldsymbol{\theta})$ . Thus, it follows:

$$1116 \|a_j^i(\boldsymbol{\theta}) - c_j(\boldsymbol{\theta})\| = \|A_i(\boldsymbol{\theta}) - C_i(\boldsymbol{\theta})\| \leq \frac{3}{\sqrt{2}} \frac{1}{\sqrt{N}} \log \frac{NM}{\delta} \gamma.$$

1118 Let

$$1119 \zeta = \frac{3}{\sqrt{2}} \frac{1}{\sqrt{N}} \log \frac{NM}{\delta} \gamma$$

1121 Since models in  $\boldsymbol{\theta}$  collection are ortho-normal or  $\Delta$ -proximal, from Lemma A.1 and Lemma A.2 it  
1122 follows that for any  $j' \neq j$ ,  $\|c_{j'}(\boldsymbol{\theta}) - c_j(\boldsymbol{\theta})\| \geq \Delta/c$ , where  $c \approx 2$ .

1124 Thus, by triangle inequality, for any  $j' \neq j$ :

$$1125 \|a_j^i(\boldsymbol{\theta}) - c_{j'}(\boldsymbol{\theta})\| - \|a_j^i(\boldsymbol{\theta}) - c_j(\boldsymbol{\theta})\| \geq \frac{\Delta}{c} - 2\zeta \quad (8)$$

1127 We now show that:

$$1129 \frac{\Delta}{c} - 2\zeta \geq \left( \frac{c'K}{m_i} + \frac{c'K}{m_j} \right) \gamma \sqrt{M} \text{polylog}(M)$$

1131 By assumption,  $N \leq d \ll M$ . Thus, in

$$1133 2\zeta + \left( \frac{c'K}{m_i} + \frac{c'K}{m_j} \right) \gamma \sqrt{M} \text{polylog}(M) = \frac{3}{\sqrt{2}} \frac{1}{\sqrt{N}} \log \frac{NM}{\delta} \gamma + \left( \frac{c'K}{m_i} + \frac{c'K}{m_j} \right) \gamma \sqrt{M} \text{polylog}(M),$$

the second term dominates. Also, we make the assumption that number of clients in each cluster is not too far from the average number of clients per cluster  $M/K$ . That is, there exists a small constant  $c''$  such that the number of clients  $m_k$  in any cluster  $k$  is at least  $\frac{1}{c''} \frac{M}{K}$ . Hence:

$$2\zeta + \left( \frac{c'K}{m_i} + \frac{c'K}{m_j} \right) \gamma \sqrt{M} \text{polylog}(M) = O \left( \gamma \frac{K^2}{\sqrt{M}} \text{polylog}(M) \right)$$

Since  $\sigma \ll 1$  is very small and  $K \leq N \leq d$ , it follows that  $\gamma \leq \sqrt{\frac{2d}{n}}$ . Thus:

$$2\zeta + \left( \frac{c'K}{m_i} + \frac{c'K}{m_j} \right) \gamma \sqrt{M} \text{polylog}(M) = O \left( \frac{\sqrt{d}}{\sqrt{n}} \frac{K^2}{\sqrt{M}} \text{polylog}(M) \right)$$

Since by assumption  $\Delta$  is close to 1 and

$$dK^4 \text{polylog}(M) = o(nM),$$

it follows that  $\frac{\Delta}{c}$  (where  $c \approx 2$ ) is much larger than  $O \left( \frac{\sqrt{d}}{\sqrt{n}} \frac{K^2}{\sqrt{M}} \text{polylog}(M) \right)$ .

This establishes, with probability at least  $1 - \delta$ :

$$\frac{\Delta}{c} \geq 2\zeta + \left( \frac{c'K}{m_i} + \frac{c'K}{m_j} \right) \gamma \sqrt{M} \text{polylog}(M)$$

Applying Equation 6 of Lemma A.7, it follows with probability at least  $1 - \delta - \frac{1}{\text{polylog}(M)}$  that:

$$\frac{\Delta}{c} - 2\zeta \geq \left( \frac{c'K}{m_i} + \frac{c'K}{m_j} \right) \|A(\boldsymbol{\theta}) - C(\boldsymbol{\theta})\|$$

Combining with Equation 8, the proximity result follows. That is, with probability at least  $1 - \delta - \frac{1}{\text{polylog}(M)}$ :

$$\forall j' \neq j, \|a_j^i(\boldsymbol{\theta}) - c_{j'}(\boldsymbol{\theta})\| - \|a_j^i(\boldsymbol{\theta}) - c_j(\boldsymbol{\theta})\| \geq \left( \frac{c'K}{m_i} + \frac{c'K}{m_j} \right) \|A(\boldsymbol{\theta}) - C(\boldsymbol{\theta})\|$$

□

This yields one of our main results:

**Theorem A.10.** *For ortho-normal or  $\Delta$ -proximal models collection  $\boldsymbol{\theta}$ , with probability at least  $1 - \delta - \frac{1}{\text{polylog}(M)}$ , k-means with suitably initialized centers recovers the accurate clustering of the SLVs (and hence an accurate clustering of the clients) in one shot.*

*Proof.* Follows from Lemma A.1 and Theorem A.9 and by applying the result of Kumar & Kannan (2010). □

*Remark A.11.* There have been further improvements to the proximity bounds of Kumar & Kannan (2010) (e.g. Awasthi & Sheffet (2012)) that can be used to improve our bounds (Equation 7). However, since our goal in this work is mainly to establish feasibility, we leave any improvements of those bounds for future work.

*Remark A.12.* In the proof above we used loss vectors of square roots of losses (SLV). This is because it allows us to recover the accurate clustering of SLVs in one round of k-means. For loss vectors (LV) formed by losses instead of square roots of losses, we are only able to prove that the distance bound as established in Lemma A.4 holds for a large fraction of the rows  $A(\boldsymbol{\theta})$ , which means that k-means can recover accurate clustering of not all but still a large number of clients. This is still fine for our approach, since we do not need to apply updates from all clients to build the models for the next round. However, it does mean that the analysis gets more complex and we therefore leave that for future work.

1188 In the following, let the set of clients in cluster  $k$  be denoted by  $Q_k$ . Let  $i \in Q_k$  be a client in cluster  $k$ .  
 1189 Let  $H_k^i$  denote the set of  $n_k$  data points and additive noise values  $(x_j, y_j, \epsilon_j)$  of client  $i$ . Its empirical  
 1190 loss is

$$1191 \mathcal{L}_k^i(\theta; H_k^i) = \frac{1}{n_k} \sum_{(x_j, y_j, \epsilon_j) \in H_k^i} (\langle \theta, x_j \rangle - y_j)^2 = \frac{1}{n_k} \sum_{(x_j, y_j, \epsilon_j) \in H_k^i} (\langle \theta_k^* - \theta, x_j \rangle + \epsilon_j)^2$$

1194 Let  $H_k$  be the  $m_k n_k \times d$  matrix with rows being the features of clients in cluster  $k$ :

$$1196 \mathcal{H}_k = \left[ x_j : (x_j, y_j, \epsilon_j) \in \bigcup_{j \in Q_k} H_k^j \right]$$

1200 Note that all rows of  $H_k$  are drawn independently from  $\mathcal{N}(0, I_d)$ . Let  $\xi_k$  be the  $m_k n_k \times 1$  column  
 1201 vector of additive noise parameters of clients in cluster  $k$ :

$$1202 \xi_k = \left[ \epsilon_j : (x_j, y_j, \epsilon_j) \in \bigcup_{j \in Q_k} H_k^j \right]$$

1206 Note that all entries of  $\xi_k$  are drawn independently from  $\mathcal{N}(0, \sigma^2)$ .

1207

1208 **Lemma A.13.** *After applying gradient updates of clients belonging to a cluster  $k$  to a model*  
 1209  $\theta$ ,  $\|\theta\| \leq 2$ , *the new model  $\theta_k$  satisfies  $\|\theta_k^* - \theta_k\| \leq \Delta/c$  with probability at least  $1 - c_6 e^{-c_7 d}$ , for*  
 1210 *some constants  $c_6, c_7$ . Furthermore, with probability at least  $1 - c_6 K e^{-c_7 d}$ ,  $\|\theta_k^* - \theta_k\| \leq \frac{\Delta}{c}$ , for all*  
 1211 *clusters  $k$ . Here learning parameter  $\eta$  is assumed to be at most  $\frac{1}{2} (1 - \frac{\Delta}{3c})$ .*

1212

1213 *Proof.* Note that

$$1215 \nabla \mathcal{L}_k^i(\theta; H_k^i) = \frac{1}{n_k} \sum_{(x_j, y_j) \in H_k^i} (2x_j x_j^T (\theta - \theta_k^*) - 2x_j \epsilon_j)$$

1217 Gradient updates from clients  $Q_k$  of cluster  $k$  on model  $\theta$  for learning rate  $\eta$  yield model  $\theta_k$ :

$$1219 \theta_k = \theta - \frac{\eta}{m_k} \sum_{i \in Q_k} \nabla \mathcal{L}_k^i(\theta; H_k^i)$$

1222 Thus, as rows of matrix  $H_k$  contain features of all clients in cluster  $k$ , and since entries of vector  $\xi_k$   
 1223 contain corresponding additive noise parameters of those clients:

$$1224 \theta_k = \theta - \frac{\eta}{m_k n_k} (2H_k^T H_k (\theta - \theta_k^*) - 2H_k^T \xi_k)$$

1225 Thus,

$$1226 \theta_k - \theta_k^* = \theta - \theta_k^* - \frac{2\eta}{m_k n_k} H_k^T H_k (\theta - \theta_k^*) + \frac{2\eta}{m_k n_k} H_k^T \xi_k$$

1227 Or

$$1228 \theta_k - \theta_k^* = \theta - \theta_k^* - \frac{2\eta}{m_k n_k} \mathbb{E}[H_k^T H_k] (\theta - \theta_k^*) + \frac{2\eta}{m_k n_k} (\mathbb{E}[H_k^T H_k] - H_k^T H_k) (\theta - \theta_k^*) + \frac{2\eta}{m_k n_k} H_k^T \xi_k$$

1229 Note  $\mathbb{E}[H_k^T H_k] = m_k n_k I$ . Therefore:

$$1230 \theta_k - \theta_k^* = (1 - 2\eta) (\theta - \theta_k^*) + \frac{2\eta}{m_k n_k} (\mathbb{E}[H_k^T H_k] - H_k^T H_k) (\theta - \theta_k^*) + \frac{2\eta}{m_k n_k} H_k^T \xi_k$$

1231 Taking norms:

$$1232 \|\theta_k - \theta_k^*\| \leq |1 - 2\eta| \|\theta - \theta_k^*\| + \frac{2\eta}{m_k n_k} \|\mathbb{E}[H_k^T H_k] - H_k^T H_k\| \|\theta - \theta_k^*\| + \frac{2\eta}{m_k n_k} \|H_k^T \xi_k\|$$

1233

1234 We now bound the various terms on the RHS.

1242 Since  $\|\theta\| \leq 2$  and  $\|\theta_k^*\| \leq 1$ :

$$1243 \quad \|\theta - \theta_k^*\| \leq 3.$$

1244 As shown in Ghosh et al. (2022), following the result of Wainwright (2019), with probability at least  
1245  $1 - 2e^{-\frac{d}{2}}$ , the operator norm of the matrix  $\mathbb{E}[H_k^T H_k] - H_k^T H_k$  is bounded from above:

$$1246 \quad \|\mathbb{E}[H_k^T H_k] - H_k^T H_k\| \leq 6\sqrt{dm_k n_k}.$$

1247 By assumption  $d \ll m_k n_k$ . Therefore,

$$1248 \quad \|\mathbb{E}[H_k^T H_k] - H_k^T H_k\| = o(m_k n_k).$$

1249 and therefore  $\frac{2\eta}{m_k n_k} \|\mathbb{E}[H_k^T H_k] - H_k^T H_k\| \|\theta - \theta_k^*\| \leq c_2 \Delta$ , for some  $c_2 \ll 1$ .

1250 As shown in Ghosh et al. (2022), with probability at least  $1 - c_4 e^{-c_5 d}$ , for some constants  $c_4, c_5$ .

$$1251 \quad \|H^T \xi_k\| \leq c_1 \sigma \sqrt{dm_k n_k} = o(\sigma m_k n_k),$$

1252 for some constant  $c_1$ . Since  $\sigma \ll 1$ ,  $\frac{2\eta}{m_k n_k} \|H^T \xi_k\| \leq c_3 \Delta$ , for some  $c_3 \ll 1$ .

1253 Combining all the bounds, with probability at least  $1 - c_6 e^{-c_7 d}$ , for some constants  $c_6, c_7$ :

$$1254 \quad \|\theta_k - \theta_k^*\| \leq 3|1 - 2\eta| + (c_2 + c_3)\Delta.$$

1255 For  $\eta \leq \frac{1}{2}(1 - \frac{\Delta}{3c})$ ,  $3|1 - 2\eta| \leq \frac{\Delta}{c}$ . Thus, with appropriate choice of  $\eta$ , with probability at least  
1256  $1 - c_6 e^{-c_7 d}$ :

$$1257 \quad \|\theta_k - \theta_k^*\| \leq \frac{\Delta}{c} + (c_2 + c_3)\Delta < \frac{\Delta}{c}.$$

1258 Thus, by applying union bound, it follows that with probability at least  $1 - c_6 K e^{-c_7 d}$ ,  $\|\theta_k^* - \theta_k\| \leq \frac{\Delta}{c}$ ,  
1259 for all clusters  $k$ .  $\square$

1260 We set  $c = 4$  and  $\eta = \frac{1}{2}(1 - \frac{\Delta}{12})$ . Thus,  $\|\theta_k^* - \theta_k\| \leq \frac{\Delta}{4}$  for all clusters  $k$  with probability at least  
1261  $1 - c_6 K e^{-c_7 d}$ .

1262 **Corollary A.14.** *After applying the gradient updates as in Lemma A.13, with probability at least  
1263  $1 - c_6 K e^{-c_7 d}$ ,  $\|\theta_k - \theta_k^*\|$  is a constant fraction smaller than  $\|\theta - \theta_k^*\|$ . The result applies to all  
1264 models  $\theta$  whose norm is less than 2.*

1265 *Proof.* Note that Equation A.2 can be restated as:

$$1266 \quad \|\theta_k - \theta_k^*\| \leq |1 - 2\eta| \|\theta - \theta_k^*\| + (c_2 + c_3)\Delta.$$

1267 For our choice of  $\eta$ ,

$$1268 \quad \|\theta_k - \theta_k^*\| \leq \frac{\Delta}{3c} \|\theta - \theta_k^*\| + (c_2 + c_3)\Delta.$$

1269 Since  $\Delta$  is close to 1 and since  $c \geq 4$ , and since  $(c_2 + c_3)\Delta$  is negligibly small, it follows that  
1270  $\|\theta_k - \theta_k^*\|$  is a constant fraction smaller than  $\|\theta - \theta_k^*\|$ , as long as  $\|\theta - \theta_k^*\|$  is much larger than  
1271  $(c_2 + c_3)\Delta$ .  $\square$

1272 In the following,  $s(i) \in [1, K]$  denotes the ground-truth cluster assignment for client  $i$ . Let  $Q$  be a  
1273 clustering of the clients. That is,  $Q = [Q_1, Q_2, \dots, Q_K]$ , where  $Q_k$  is the set of clients in the  $k$ -th  
1274 cluster of  $Q$ . Let  $G(Q; \theta) = (U \cup V, E)$  be a fully-connected bi-partite graph. Here,  $U$  and  $V$  are  
1275 sets of nodes of size  $K$  each, where  $k$ -th node of  $U$  (or  $V$ ) corresponds to the  $k$ -th cluster of  $Q$ .  $E$  is  
1276 the set of edges, with weight  $w(k, j)$  of edge  $(k, j)$  set to the total loss of clients in  $Q_k$  on model  $\theta_j$   
1277 of  $\theta$ . That is:

$$1278 \quad w(k, j) = \sum_{i \in Q_k} \mathcal{L}_{s(i)}^i(\theta; H_{s(i)}^i).$$

1296 **Assumption A.15.**  $n_k m_k (\frac{\rho}{\rho+1})^2 \gg \text{polylog}(K)$ , for  $\rho = \frac{\Delta^2}{\sigma^2}$ . In the following,  
 1297  $\theta = [\theta_1 \ \theta_2 \ \dots \ \theta_K]$  refers to a collection of  $K$  models that satisfy the  $\Delta$ -proximity  
 1298 condition  $\forall k \in [K], \|\theta_k - \theta_k^*\| \leq \frac{\Delta}{c}$ , for  $c \geq 4$ .  
 1299

1300 **Lemma A.16.** Let  $\theta$  be a collection of  $K$  models that satisfy the  $\Delta$ -proximity condition. Let  $Q$  be an  
 1301 accurate clustering of clients to clusters. Then, the probability that min-cost perfect matching  $M$  of  
 1302  $G(Q; \theta)$  is  $M = \{(k, k) : Q_k \in Q\}$  is at least  $1 - c_1 e^{-c_2 \frac{1}{\log K^2} n_k m_k (\frac{\rho}{\rho+1})^2}$ .  
 1303

1304 *Proof.* Let clients of a cluster  $Q_k$  achieve lower loss on a model  $\theta_j$  than on model  $\theta_k$ , for  $j \neq k$ .  
 1305 Note that the probability of this is:  
 1306

$$1307 \mathbb{P} \left[ \sum_{i \in Q_k} \mathcal{L}_{s(i)}^i (\theta_j; H_{s(i)}^i) < \sum_{i \in Q_k} \mathcal{L}_{s(i)}^i (\theta_k; H_{s(i)}^i) \right]$$

1311 Since  $\sum_{i \in Q_k} \mathcal{L}_{s(i)}^i (\theta; H_{s(i)}^i)$  is scaled chi-squared random variable with  $n_k m_k$  degrees of freedom  
 1312 (i.e.  $\sim \frac{m_k}{n_k} (\|\theta - \theta_k^*\|^2 + \sigma^2) \chi^2(n_k m_k)$ ), this probability is:  
 1313

$$1314 \mathbb{P} [(\|\theta_j - \theta_k^*\|^2 + \sigma^2) \chi^2(n_k m_k) < (\|\theta_k - \theta_k^*\|^2 + \sigma^2) \chi^2(n_k m_k)]$$

1316 Since  $\|\theta_k - \theta_k^*\| \leq \frac{\Delta}{4} = \frac{\Delta}{2} - \frac{\Delta}{4}$  and since  $\|\theta_k^* - \theta_j^*\| \geq \Delta$ , by triangle inequality it follows that  
 1317  $\|\theta_k - \theta_j^*\| \geq \frac{3\Delta}{4} = \frac{\Delta}{2} - \frac{\Delta}{4}$ . Thus, from the result of Ghosh et al. (2022), which is obtained by  
 1318 applying the concentration properties of chi-squared random variables (Wainwright, 2019), it follows  
 1319 that for some constants  $c_1$  and  $c_2$ :

$$1320 \mathbb{P} [(\|\theta_j - \theta_k^*\|^2 + \sigma^2) \chi^2(n_k m_k) < (\|\theta_k - \theta_k^*\|^2 + \sigma^2) \chi^2(n_k m_k)] \leq c_1 e^{-c_2 n_k m_k (\frac{\rho}{\rho+1})^2}$$

1322 Hence:  
 1323

$$1324 \mathbb{P} \left[ \sum_{i \in Q_k} \mathcal{L}_{s(i)}^i (\theta_j; H_{s(i)}^i) \leq \sum_{i \in Q_k} \mathcal{L}_{s(i)}^i (\theta_k; H_{s(i)}^i) \right] \leq c_1 e^{-c_2 n_k m_k (\frac{\rho}{\rho+1})^2}$$

1327 Thus, by union bound it follows:  
 1328

$$1329 \mathbb{P} \left[ \exists k \neq j \text{ s.t. } \sum_{i \in Q_k} \mathcal{L}_{s(i)}^i (\theta_j; H_{s(i)}^i) \leq \sum_{i \in Q_k} \mathcal{L}_{s(i)}^i (\theta_k; H_{s(i)}^i) \right] \leq K^2 c_1 e^{-c_2 n_k m_k (\frac{\rho}{\rho+1})^2}$$

1332 Thus:  
 1333

$$1334 \mathbb{P} \left[ \forall k \neq j, \sum_{i \in Q_k} \mathcal{L}_{s(i)}^i (\theta_k; H_{s(i)}^i) < \sum_{i \in Q_k} \mathcal{L}_{s(i)}^i (\theta_j; H_{s(i)}^i) \right]$$

$$1335 \geq 1 - K^2 c_1 e^{-c_2 n_k m_k (\frac{\rho}{\rho+1})^2}$$

$$1336 \geq 1 - c_1 e^{-c_2 \frac{1}{\log K^2} n_k m_k (\frac{\rho}{\rho+1})^2}$$

1340 Thus, with probability close to 1, clients in cluster  $Q_k$ , achieve lowest loss on model  $\theta_k$ , for all  $k$ .  
 1341

1342  $\square$   
 1343

### 1344 A.3 COMPARISON WITH *IFCA*

1345 We demonstrate that the cluster recovery condition of *IFCA* implies loss vector separation for *CLoVE*.  
 1346 We assume a gap of at least  $\delta > 0$  between the loss a client achieves on its best model and its second  
 1347 best model. The following Lemma shows that vectors of clients belonging to different clusters are at  
 1348 least  $\sqrt{2}\delta$  apart.  
 1349

1350  
**Lemma A.17.** Let clients  $i$  and  $j$  achieve lowest loss on models  $k_1$  and  $k_2$  respectively ( $k_1 \neq k_2$ ).  
 1351 Let there be at least a  $\delta$  gap between the lowest loss and the second lowest loss achieved by a client  
 1352 on any model. Then  $\|L_i - L_j\| \geq \sqrt{2}\delta$ , where  $L_i$  and  $L_j$  are the loss vectors for clients  $i$  and  $j$ .  
 1353

1354  
 1355 *Proof.* Let  $l(x, y)$  denote loss of client  $x$  on model  $y$ . Then:

1356 
$$\|L_i - L_j\|^2 \geq (l(i, k_1) - l(j, k_1))^2 + (l(i, k_2) - l(j, k_2))^2$$

1358 However, because of  $\delta$  loss gap:

1360 
$$l(i, k_2) \geq l(i, k_1) + \delta$$

1361 and

1363 
$$l(j, k_2) \leq l(j, k_1) - \delta$$

1364 Thus:

1365 
$$\|L_i - L_j\|^2 \geq (l(i, k_1) - l(j, k_1))^2 + (l(i, k_1) - l(j, k_1) + 2\delta)^2$$

1367 The right side is minimized when  $(l(i, k_1) - l(j, k_1))^2 = (l(i, k_1) - l(j, k_1) + 2\delta)^2$ , or when  
 1368  $l(i, k_1) - l(j, k_1) = -\delta$ . Thus:

1369 
$$\|L_i - L_j\|^2 \geq 2\delta^2$$

1371 Hence, the result follows. □

1373 Note that in *IFCA* the best model for each cluster is initialized to be close to its optimal counterpart.  
 1374 Consequently,  $\delta$  can be much larger than 0. Lemma A.17 therefore implies significant distance  
 1375 between the loss vectors of clients of different clusters.  
 1376

## 1378 B DISCUSSION

### 1380 B.1 COMMUNICATION AND COMPUTATION COSTS

1382 Compared with other PFL methods like *IFCA* that exchange only model weights, *CLoVE* additionally  
 1383 sends the loss vectors. The size of the loss vectors is  $\Theta(K)$  per client. This is small compared to  
 1384 the model weights themselves, which are  $\Theta(K * R)$  per client, where  $R$  is the number of model  
 1385 parameters (or gradients). Also, in practice,  $K$  is usually small. So the overhead of communicating  
 1386 the loss vectors is not significant.

1387 Furthermore, thanks to *CLoVE*'s early stopping feature, all models need only to be sent to all clients  
 1388 for the rounds until stability is reached. As shown in our experiments (see App. D.1.5), stability is  
 1389 achieved in a few rounds with high probability. After that, only a single model needs to be sent to  
 1390 each client. From then on, each cluster performs regular *FedAvg*, so the cost is similar to *FedAvg*,  
 1391 which is the best possible (other baselines also have a computation and communication cost at least  
 1392 as high as that of *FedAvg*).

1393 Another technique for further reducing the size of the communicated models in the first few rounds is  
 1394 to leverage weight-sharing techniques from multi-task learning, as suggested also in *IFCA*: the  $K$   
 1395 models can share a few initial layers that are meant to learn the common properties of the data, and  
 1396 the rest of the layers are meant to learn the distinct features of each cluster. This way, when the server  
 1397 sends the models to each client, it needs to only send the common part once alongside the distinct  
 1398 parts, thus achieving savings in communication costs.

1399 Finally, one might argue that *CLoVE* might result in high peak memory usage, due to clients evaluating  
 1400 all the models before cluster stability is reached. This can be alleviated as follows: the  $K$  models do  
 1401 not need to be evaluated by the client all at the same time. Instead, if memory is a constraint, the  
 1402 client can be requesting and evaluating each of the  $K$  models sequentially one at a time, storing in  
 1403 memory only the loss produced by each model. This way, it can form the loss vector while only  
 holding one model in memory at a time.

1404  
1405 B.2 PRIVACY

1406 Although privacy is not the focus of this paper, it is an important concern in FL and we now briefly  
 1407 discuss it. Compared to standard FL algorithms such as *FedAvg* that share model weights/gradients,  
 1408 clients in *CLoVE* additionally share the average model losses. Importantly, these loss values are  
 1409 not the raw model outputs; they are computed based on the output and the (unknown to the server)  
 1410 input or labels, and also averaged for all datapoints of the client before being sent to the server.  
 1411 While this aggregate measure is generally not considered privacy-sensitive, its privacy implications  
 1412 remain underexplored. Further research is needed to assess whether, and under what conditions,  
 1413 such aggregate statistics could lead to unintended information leakage, not with regard to *CLoVE* in  
 1414 particular, but more generally. Moreover, several baselines (e.g. Tan et al. (2022); Xu et al. (2023);  
 1415 Vahidian et al. (2023)) similarly share 1-D vectors or summaries of client data (we only send losses,  
 1416 not summaries of data). This process is deemed irreversible by those baselines, and they even mention  
 1417 that further privacy-preserving techniques can be employed.

1418  
1419 B.3 COMPARISON WITH FEDGWC PAPER

1420 Licciardi et al. (2025) is a hierarchical approach like Sattler et al. (2021), and as such the number  
 1421 of rounds for clustering can be very large, as is also evident from the analytical guarantees that are  
 1422 only in the asymptotic regime as the number of rounds grows without bound; in contrast, *CLoVE*  
 1423 achieves correct clustering in 2-3 rounds, as shown both analytically and empirically. Moreover, in  
 1424 their approach, amount of state kept at server is quadratic in number of clients. With our approach,  
 1425 the amount of state at the server is not a function of the number of clients, but rather it is linear on the  
 1426 number of clusters, which is typically much smaller than number of clients. Additionally, FedGWC’s  
 1427 clustering happens on the basis of losses computed on a single model. Specifically, the vectors used  
 1428 in Eq. 4 for computing entries of matrix  $W$  are based on the losses (and rewards) of clients on the  
 1429 single model that is sent to the clients in that round. On the other hand, *CLoVE* performs clustering  
 1430 by comparing clients’ losses across multiple models. This can be more effective, as losses of multiple  
 1431 models amplify the difference in data distributions among clients belonging to different clusters. This  
 1432 is why *CLoVE* is able to recover clusters so much faster.

1433  
1434 C DETAILS ON DATASETS, MODELS AND COMPUTE USED  
1435

1436 **Datasets** We use the MNIST dataset (LeCun et al., 1998), containing handwritten digit images  
 1437 from 0 to 9 (10 classes; 60K points training, 10K testing), the CIFAR10 dataset (Krizhevsky, 2009),  
 1438 containing color images from 10 categories (airplane, bird, etc.) (50K points training, 10K testing),  
 1439 and the Fashion-MNIST (FMNIST) image dataset (Xiao et al., 2017), containing images of clothing  
 1440 items from 10 categories (T-shirt, dress, etc.) (60K points training, 10K testing). Beyond image  
 1441 datasets, we also use two text datasets: Amazon Review Data (Ni et al., 2019) (denoted by Amazon  
 1442 Review), containing Amazon reviews from 1996-2018 for products from 4 categories (2834 for  
 1443 books, 1199 for DVDs, 1883 for electronics, and 1755 for kitchen and houseware) and the consumer  
 1444 sentiment for each product based on the reviews (binary – positive/negative), and AG News (Gulli;  
 1445 Zhang et al., 2015), containing news articles (30000 training and 1900 testing) from 4 classes (world,  
 1446 sports, business, and science). Finally, we use the Federated EMNIST (FEMNIST) dataset (Caldas  
 1447 et al., 2019), an FL dataset with data points being handwritten letters or digits from a particular  
 1448 human writer. All these datasets are also used by the baselines.

1449 **Models** For MNIST and FMNIST unsupervised image reconstruction, we used a fully connected  
 1450 autoencoder with Sigmoid output activation. For CIFAR-10 unsupervised image reconstruction,  
 1451 we used a Convolutional Autoencoder with a fully connected latent bottleneck and a Tanh output  
 1452 activation. For supervised MNIST, FMNIST and FEMNIST image classification we used a simple  
 1453 convolutional neural network (CNN) classifier designed for grayscale images. Each convolutional  
 1454 layer uses: 5×5 kernels, Stride 1 (default) followed by ReLU and 2×2 max pooling. Its output is  
 1455 a logit vector. For CIFAR-10 image classification, we used a deep convolutional neural network  
 1456 classifier, designed for 32×32 RGB images. It uses batch normalization for stable training and dropout  
 1457 in fully connected layers to reduce overfitting. For AG News we used a TextCNN with an embedding  
 1458 layer, and for Amazon News we used a simple multi-layer perceptron (MLP) model.

Note that we are the first to propose an algorithm that can cluster correctly *using autoencoders in a fully unsupervised setting*. The datasets we use can all be clustered successfully using simple models, thanks to *CLoVE*'s superior ability to separate the clients, so there is no need to employ heavier models for these tasks. (If a task requires a certain type of model, *CLoVE* should again work well, as it is only based on the loss that the model outputs and the model itself is used as a black box).

**Experiment setup** We provide the details on the number of clients and datapoints used for each of our experiments in Table 6.

Table 6: Details on number of clients and datapoints (S: Supervised, U: Unsupervised)

Experiment	Number of clients			Number of datapoints per client		
	MNIST	CIFAR-10	FMNIST	MNIST	CIFAR-10	FMNIST
S: Label skew 1	25	15	25	500	500	500
S: Label skew 2	25	15	25	500	500	500
S: Label skew 3	25	25	25	500	500	500
S: Label skew 4	50	30	50	500	500	500
S: Feature skew	40	20	40	500	500	500
S: Concept shift	20	12	20	1000	1000	1000
S: Amazon Review		20			1000	
S: AG News		30			100	
S: Ablation (CIFAR-10)		25			500	
U: Initialization (MNIST)		50			500	
U: Linear		25			1000	
U: Everything else		50			1000	

**GitHub links** The GitHub links for the code for each of the baseline algorithms used in the paper are shown in Table 7.

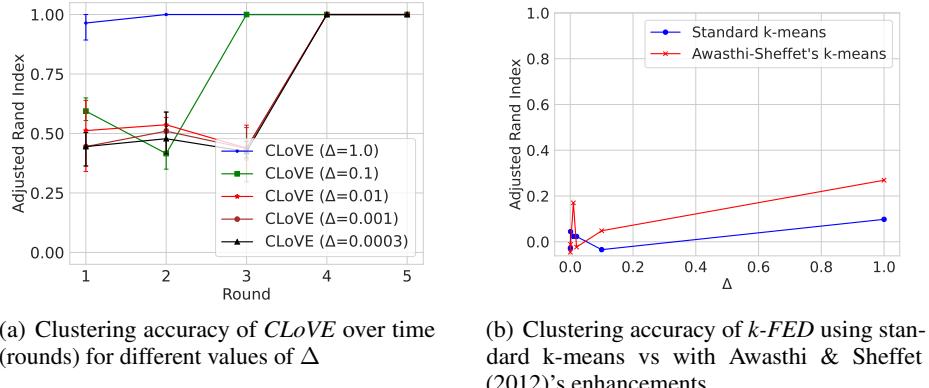
Table 7: GitHub links for comparison algorithms

Algorithm	Source code
Per-FedAvg	<a href="https://github.com/TsingZ0/PFLlib">https://github.com/TsingZ0/PFLlib</a>
FedProto	<a href="https://github.com/TsingZ0/PFLlib">https://github.com/TsingZ0/PFLlib</a>
FedALA	<a href="https://github.com/TsingZ0/PFLlib">https://github.com/TsingZ0/PFLlib</a>
FedPAC	<a href="https://github.com/TsingZ0/PFLlib">https://github.com/TsingZ0/PFLlib</a>
CFL-S	<a href="https://github.com/felisat/clustered-federated-learning">https://github.com/felisat/clustered-federated-learning</a>
FeSEM	<a href="https://github.com/morningD/FlexCFL">https://github.com/morningD/FlexCFL</a>
FlexCFL	<a href="https://github.com/morningD/FlexCFL">https://github.com/morningD/FlexCFL</a>
PACFL	<a href="https://github.com/MMorafah/PACFL">https://github.com/MMorafah/PACFL</a>
FedAvg	self-implemented
Local-only	self-implemented
IFCA	self-implemented
CLoVE	self-implemented

**Compute used** We train the neural networks and run all our evaluations on a server with 98 cores with Intel(R) Xeon(R) Gold 6248R CPU @ 3.00GHz, and 256GB RAM. For our experiments, we used Python 3.12 and PyTorch (Paszke et al., 2019), version 2.7.0.

1512 **D ADDITIONAL EXPERIMENTAL RESULTS**  
15131514 **D.1 ADDITIONAL SUPERVISED RESULTS**  
15151516 **D.1.1 MIXED LINEAR REGRESSION EXPERIMENTS: COMPARISON WITH k-FED**  
1517

1518 We conduct a study to evaluate the impact of closeness of the different clusters optimal model  
1519 parameters on *CLoVE*’s clustering performance. We select a supervised setting, where each cluster’s  
1520 data is generated by a linear model, as described in Sec. 4. We randomly select optimal model  
1521 parameters  $\theta_k^* \in \mathbb{R}^d$  from a unit  $d$ -dimensional sphere, ensuring that the pairwise distance between  
1522 them falls within the range  $[\Delta, 5\Delta]$ . Our experimental setup consists of  $K = 5$  clusters and  $M = 25$   
1523 clients; each client has  $n = 1000$  data points, and we use  $\tau = 25$  local epochs. Figure 3(a) illustrates  
1524 the clustering accuracy of *CLoVE* over multiple rounds, with  $\Delta$  varying from 0.0003 to 1.0.  
1525

1535 **Figure 3: Linear regression experiment results**  
1536

1537 As shown, *CLoVE* consistently recovers the correct cluster assignment. Moreover, as expected, the  
1538 convergence rate of *CLoVE* slows down when the optimal model parameters are closer together  
1539 (i.e. when  $\Delta$  is smaller). In contrast, the performance of *k-FED* (Dennis et al., 2021) on the same  
1540 data, even with Awasthi-Sheffet’s k-means enhancements (Awasthi & Sheffet, 2012), is significantly  
1541 worse (Fig. 3(b)). This discrepancy can be attributed to the fact that the cluster centers (i.e. the mean  
1542 of their data) are very close to the origin, highlighting the limitations of *k-FED* in such scenarios.  
1543 Our results demonstrate that *CLoVE* outperforms *k-FED*, especially when the cluster centers are not  
1544 well-separated.  
1545

1546 **D.1.2 ADDITIONAL LABEL SKEWS**  
1547

1548 Table 8 presents image classification test accuracy results for MNIST, FMNIST, CIFAR-10 datasets  
1549 on additional types of label skews. The corresponding ARI accuracy results are presented in Table 9  
1550 and the clustering convergence results in Table 10. These results include two cases *Label skew 3*  
1551 and *Label skew 4*. The case *Label skew 3* is a data distribution in which there is significant overlap  
1552 between labels assigned to clients of different clusters. Specifically, there are 5 clusters with 5 clients  
1553 each. Among the first 4 clusters, all clients of a cluster get samples from all 10 classes except one.  
1554 The missing class is uniquely selected for each of the 4 clusters. The clients of the 5th cluster get  
1555 samples from all 10 classes. The case *Label skew 4* is a data distribution in which  $s\%$  of data (we use  
1556 one third) for clients belonging to a cluster is uniformly sampled from all classes, and the remaining  
1557  $(100 - s)\%$  comes from a dominant class that is unique for each cluster (Xu et al., 2023). Each  
1558 cluster has 5 clients each of which is assigned 500 samples. For *Label skew 3*, the data of a class is  
1559 distributed amongst the clients that are assigned that class by sampling from a Dirichlet distribution  
1560 with parameter  $\alpha = 0.5$ , according to the procedure described in Section 5.  
1561

1562 We observe that in some cases for *Label skew 3* and *Label skew 4*, *FedAvg* is achieving the highest  
1563 test accuracy. This can be attributed to the fact that under these cases, each client holds data from  
1564 almost every data class, which makes the client distributions very similar. Thus, naturally, the single  
1565

model that *FedAvg* trains is also a good model for the individual clients. Moreover, the model of *FedAvg* is trained on the data of all clients, while the other algorithms only use the data of a single cluster to train each of their models. This gives *FedAvg* an advantage, as it trained on much more data for the same number of rounds.

Table 8: Supervised test accuracy results for MNIST, CIFAR-10 and FMNIST

Data mixing	Algorithm	MNIST	CIFAR-10	FMNIST
Label skew 3	FedAvg	<b>99.0 ± 0.0</b>	<b>84.1 ± 0.5</b>	<b>88.8 ± 0.3</b>
	Local-only	85.9 ± 0.6	41.1 ± 1.5	71.1 ± 0.4
	Per-FedAvg	95.3 ± 0.3	41.6 ± 0.7	74.9 ± 0.7
	FedProto	88.8 ± 0.2	39.9 ± 0.4	74.1 ± 0.2
	FedALA	94.7 ± 0.3	40.0 ± 0.7	75.1 ± 0.5
	FedPAC	95.6 ± 0.2	41.4 ± 0.6	75.5 ± 0.8
	CFL-S	97.9 ± 0.1	61.2 ± 0.6	84.3 ± 1.2
	FeSEM	94.0 ± 0.2	13.3 ± 0.2	72.2 ± 1.9
	FlexCFL	92.2 ± 0.3	20.8 ± 0.5	79.7 ± 0.5
	PACFL	82.8 ± 0.9	27.1 ± 3.8	68.1 ± 2.1
Label skew 4	IFCA	94.2 ± 0.2	57.3 ± 1.9	80.6 ± 1.0
	CLoVE	94.2 ± 0.2	56.9 ± 0.5	79.5 ± 0.5
	FedAvg	<b>98.5 ± 0.1</b>	60.7 ± 0.8	85.9 ± 0.2
	Local-only	94.8 ± 0.2	67.1 ± 0.5	85.5 ± 0.1
	Per-FedAvg	95.6 ± 0.1	63.7 ± 0.8	83.7 ± 0.1
	FedProto	93.3 ± 0.2	45.3 ± 0.7	76.5 ± 0.3
	FedALA	96.8 ± 0.1	69.8 ± 0.3	87.3 ± 0.1
	FedPAC	96.1 ± 0.2	66.3 ± 0.5	84.8 ± 0.3
	CFL-S	97.8 ± 0.1	71.3 ± 1.1	88.2 ± 0.4
	FeSEM	94.9 ± 0.3	11.9 ± 1.5	80.6 ± 1.1
Label skew 5	FlexCFL	97.2 ± 0.1	66.8 ± 0.0	89.5 ± 0.1
	PACFL	93.4 ± 0.2	67.4 ± 0.3	84.1 ± 0.2
	IFCA	98.0 ± 0.1	64.3 ± 0.9	89.4 ± 0.2
	CLoVE	97.8 ± 0.2	<b>72.6 ± 0.3</b>	<b>89.9 ± 0.1</b>

Table 9: Supervised ARI results for MNIST, CIFAR-10 and FMNIST

Data mixing	Algorithm	MNIST	CIFAR-10	FMNIST
Label skew 3	CFL-S	0.00 ± 0.00	0.77 ± 0.18	0.00 ± 0.00
	FeSEM	0.18 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
	FlexCFL	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>
	PACFL	0.20 ± 0.08	0.00 ± 0.00	0.15 ± 0.13
	IFCA	0.92 ± 0.12	0.69 ± 0.08	0.65 ± 0.15
	CLoVE	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>
	CFL-S	0.00 ± 0.00	0.61 ± 0.03	0.29 ± 0.02
	FeSEM	0.12 ± 0.00	0.00 ± 0.00	0.11 ± 0.02
	FlexCFL	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>
	PACFL	0.53 ± 0.06	0.66 ± 0.07	0.92 ± 0.02
Label skew 4	IFCA	0.55 ± 0.03	0.36 ± 0.22	0.50 ± 0.05
	CLoVE	0.60 ± 0.10	<b>1.00 ± 0.00</b>	0.95 ± 0.07

1620

1621 Table 10: Convergence behavior for MNIST, CIFAR-10 and FMNIST

1622 Data mixing	1623 Algorithm	1624 ARI reached in 10 rounds			1625 First round when ARI $\geq 0.9$		
		1626 MNIST	1627 CIFAR-10	1628 FMNIST	1629 MNIST	1630 CIFAR-10	1631 FMNIST
1629 Label skew 2	CFL-S	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	—	—	—
	FeSEM	0.18 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	—	—	—
	FlexCFL	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.0 <math>\pm</math> 0.0</b>	<b>1.0 <math>\pm</math> 0.0</b>	<b>1.0 <math>\pm</math> 0.0</b>
	PACFL	0.35 $\pm$ 0.12	0.32 $\pm$ 0.11	0.55 $\pm$ 0.03	—	—	—
	IFCA	0.83 $\pm$ 0.12	0.72 $\pm$ 0.00	0.92 $\pm$ 0.12	—	—	—
	CLoVE	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	2.0 $\pm$ 0.0	2.0 $\pm$ 0.0	2.0 $\pm$ 0.0
1630 Label skew 3	CFL-S	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	—	—	—
	FeSEM	0.18 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	—	—	—
	FlexCFL	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.0 <math>\pm</math> 0.0</b>	<b>1.0 <math>\pm</math> 0.0</b>	<b>1.0 <math>\pm</math> 0.0</b>
	PACFL	0.20 $\pm$ 0.08	0.00 $\pm$ 0.00	0.15 $\pm$ 0.13	—	—	—
	IFCA	0.92 $\pm$ 0.12	0.69 $\pm$ 0.08	0.65 $\pm$ 0.15	—	—	—
	CLoVE	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	2.0 $\pm$ 0.0	2.0 $\pm$ 0.0	2.0 $\pm$ 0.0
1635 Label skew 4	CFL-S	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	—	—	—
	FeSEM	0.12 $\pm$ 0.00	0.00 $\pm$ 0.00	0.11 $\pm$ 0.02	—	—	—
	FlexCFL	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.0 <math>\pm</math> 0.0</b>	<b>1.0 <math>\pm</math> 0.0</b>	<b>1.0 <math>\pm</math> 0.0</b>
	PACFL	0.53 $\pm$ 0.06	0.66 $\pm$ 0.07	0.92 $\pm$ 0.02	—	—	—
	IFCA	0.54 $\pm$ 0.03	0.70 $\pm$ 0.12	0.50 $\pm$ 0.05	—	—	—
	CLoVE	0.90 $\pm$ 0.07	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	2.0 $\pm$ 0.0	2.0 $\pm$ 0.0	2.0 $\pm$ 0.0
1641 Feature skew	CFL-S	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	—	—	—
	FeSEM	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	—	—	—
	FlexCFL	<b>1.00 <math>\pm</math> 0.00</b>	0.11 $\pm$ 0.08	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.0 <math>\pm</math> 0.0</b>	—	<b>1.0 <math>\pm</math> 0.0</b>
	PACFL	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.68 $\pm$ 0.11	—	—	—
	IFCA	0.67 $\pm$ 0.28	0.71 $\pm$ 0.22	0.55 $\pm$ 0.10	—	—	—
	CLoVE	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	<b>1.00 <math>\pm</math> 0.00</b>	2.0 $\pm$ 0.0	<b>4.3 <math>\pm</math> 0.9</b>	2.0 $\pm$ 0.0

1646

1647

1648 Table 11: Supervised ARI results for Amazon Review and AG News

1649 Algorithm	1650 Amazon Review	1651 AG News
CFL-S	0.07 $\pm$ 0.08	0.39 $\pm$ 0.07
FeSEM	0.00 $\pm$ 0.01	0.14 $\pm$ 0.00
FlexCFL	0.00 $\pm$ 0.02	0.02 $\pm$ 0.00
PACFL	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00
IFCA	0.29 $\pm$ 0.20	0.64 $\pm$ 0.12
CLoVE	<b>0.73 <math>\pm</math> 0.24</b>	<b>0.94 <math>\pm</math> 0.09</b>

1656

1657

1658

## 1659 D.1.3 ADDITIONAL DETAILS FOR TEXT CLASSIFICATION

1660

1661 The ARI accuracy and speed of convergence results for the Amazon Review and AG News datasets  
1662 are presented in Table 11 and Table 12, respectively.

1663

1664 D.1.4 THE CASE OF UNKNOWN NUMBER OF CLUSTERS  $K$ 

1665

1666

1667

1668

1669

1670

1671

1672

1673

As mentioned in Section 3, *CLoVE* does not necessarily require the true number of clusters to be known in advance and given as an input, but it can converge to an appropriate number of clusters during its execution. To demonstrate this feature of *CLoVE*, we perform the following experiment on the FEMNIST dataset. We use  $M = 100$  clients and assign to each client all of the data from one of the human writers in FEMNIST. We use a CNN and random first-round weight initialization. In this case, the number of true clusters (which human writers have a similar writing style) is unknown. *CLoVE* searches over different possibilities for the number of clusters/models, achieving an accuracy of  $68.2 \pm 1.3$ . This demonstrates the versatility of our algorithm in general settings where the number of true clusters is not given as input.

1674

1675

1676

1677

1678

1679

1680

1681

1682

1683

1684

1685

1686

Table 12: Convergence behavior for the Amazon Review and AG News datasets

Algorithm	ARI reached in 10 rounds		First round when ARI $\geq 0.9$	
	Amazon	AG News	Amazon	AG News
CFL-S	$0.00 \pm 0.00$	$0.00 \pm 0.00$	—	—
FeSEM	$0.00 \pm 0.01$	$0.14 \pm 0.00$	—	—
FlexCFL	$0.00 \pm 0.02$	$0.02 \pm 0.00$	—	—
PACFL	$0.00 \pm 0.00$	$0.00 \pm 0.00$	—	—
IFCA	$0.29 \pm 0.20$	$0.64 \pm 0.12$	—	—
CLoVE	<b><math>0.72 \pm 0.04</math></b>	<b><math>0.94 \pm 0.09</math></b>	—	<b><math>2.0 \pm 0.0</math></b>

## D.1.5 EARLY STOPPING OF CLUSTERING

As mentioned in Section 3, *CLoVE* has the ability to stop clustering (and the associated sending of all models by the server to each client and the inference by each client on all the models to generate the loss vectors) early once it stabilizes, so as to save in communication and computation costs. In this section, we demonstrate the fact that *CLoVE* indeed can reach a stable clustering within a few rounds. We run the same supervised classification experiments on CIFAR-10 for different modes of data mixing as before, but now allowing early stopping when the clustering stabilizes (we define stability as 3 rounds of no change). In all experiments, *CLoVE* achieves an ARI of 1.0. The results for the test accuracy and the round at which cluster stability is reached are shown in Table 13. We observe that *CLoVE* is able to reach a stable client-to-cluster assignment within 4-7 rounds in all scenarios.

Table 13: Early stopping experiments (CIFAR-10)

Data mixing	Test accuracy	Round of stability
Label skew 1	$90.7 \pm 0.2$	$4.0 \pm 0.0$
Label skew 2	$67.2 \pm 1.5$	$4.0 \pm 0.0$
Feature skew	$58.2 \pm 1.2$	$6.3 \pm 1.2$
Concept shift	$59.4 \pm 0.4$	$4.0 \pm 0.0$

## D.1.6 PARTIAL PARTICIPATION

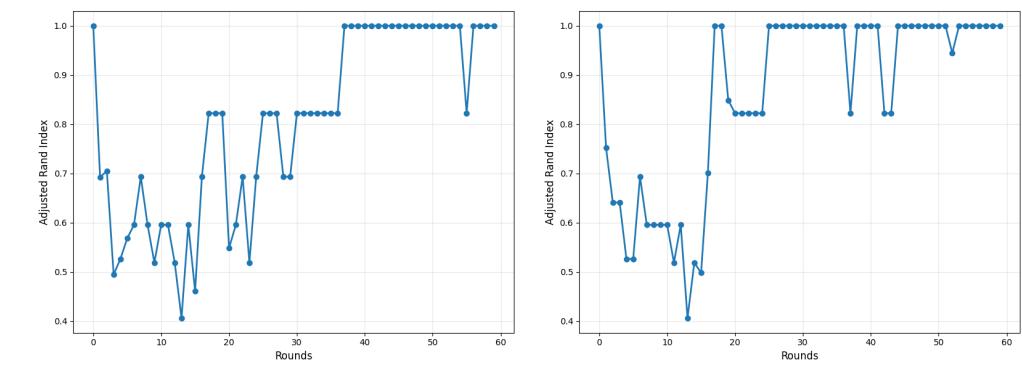
We now experiment with different levels of client participation (parameter  $\rho$ ). We use MNIST in a supervised setting with a CNN,  $K = 5$  clusters,  $M = 125$  clients (25 per cluster),  $n = 100$  datapoints/client, and random initialization of the model weights. The results are shown in Table 14. Note that for partial participation the ARI at each round is measured using only the clients that participated in that round and comparing their clustering with the ground truth clustering projected on those clients. We see that *CLoVE* very quickly and consistently achieves perfect ARI and almost perfect accuracy even for very small (e.g. 10%) client participation rates.

Table 14: *CLoVE*’s performance under partial participation

Participation rate $\rho$	Test accuracy	Final ARI	ARI reached in 10 rounds	First round when ARI $\geq 0.9$
1.0	$99.5 \pm 0.1$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$2.0 \pm 0.0$
0.9	$99.6 \pm 0.1$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$2.0 \pm 0.0$
0.75	$99.4 \pm 0.2$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$2.0 \pm 0.0$
0.5	$99.6 \pm 0.1$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$2.0 \pm 0.0$
0.25	$99.4 \pm 0.1$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$2.0 \pm 0.0$
0.1	$99.6 \pm 0.2$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$3.0 \pm 0.8$

We also perform an experiment for the extreme case where participation is so low that the number of participating clients is lower than the true number of clusters. We examine *CLoVE*’s behavior with early stopping after cluster stability (Fig. 4(a)) and without early stopping (Fig. 4(b)). We use MNIST with  $K = 10$  clusters,  $m = 2$  clients per cluster and  $\rho = 40\%$  participation rate (so only 8 clients

1728  
 1729 participate at each round, which is lower than the true number of clusters, 10). The ARI is evaluated  
 1730 only for clients that have participated at least once in the training process until that round. In the first  
 1731 few rounds, the ARI is higher, then it drops, and then again rises to the optimal. This is because in  
 1732 the beginning very few clients have participated until that point, so there are not many wrong cluster  
 1733 assignments. As more clients join, initially the number of incorrect assignments increases, but then,  
 1734 as *CLoVE* enhances the assignments over time, the ARI increases once again. As can be seen, after a  
 1735 few tens of rounds, *CLoVE* manages to almost reach optimal clustering, even though this setting is  
 1736 very challenging.



1737  
 1738 Figure 4: ARI for very few participants – fewer than the number of clusters – with early stopping  
 1739 (left) and without early stopping (right).  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747

## 1752 D.2 ADDITIONAL UNSUPERVISED RESULTS

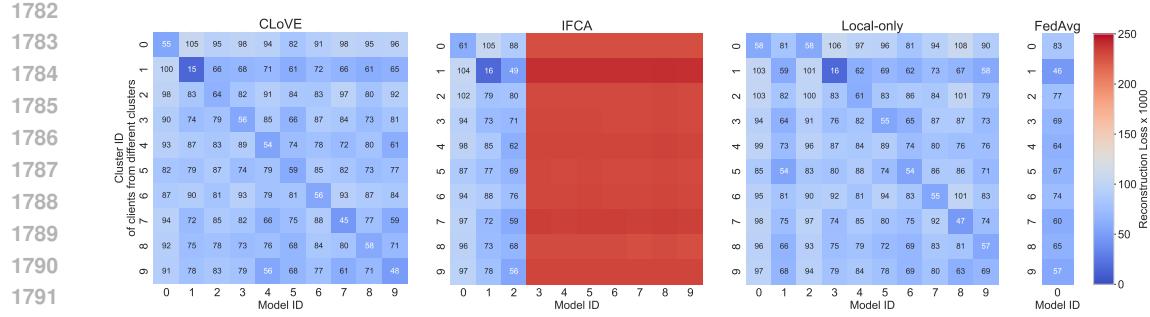
1753  
 1754 For this section’s experiments, we have used a batch size of 64.  
 1755

### 1756 D.2.1 BENEFIT OF PERSONALIZATION

1757 We now examine the reconstruction loss performance of the trained models and show the benefit  
 1758 of our personalized models over traditional FL models. For MNIST and  $K = 10$ ,  $M = 50$  and  
 1759  $n = 250$ , we run 4 different algorithms: *CLoVE*, *IFCA*, *FedAvg*, *Local-only*. In *FedAvg*, i.e. the  
 1760 original FL algorithm, a single model is trained and aggregated by all clients, regardless of their true  
 1761 cluster (equivalently, *FedAvg* is *CLoVE* with  $K = 1$ ). *Local-only* means federated models but no  
 1762 aggregation: each client trains its own model only on its local data, and does not share it with the rest  
 1763 of the system.  
 1764

1765 For this experiment, we consider that each of the 50 clients has a testset of the same data class of the  
 1766 client’s trainset, and we feed its testset through all models (10 models for *CLoVE* and *IFCA*, 1 for  
 1767 *FedAvg* and 50 for *Local-only*), creating the loss matrices shown as heatmaps in Fig. 5 (run for a single  
 1768 seed value). The vertical axis represents clients, one selected from each true cluster, and the horizontal  
 1769 represents the different models. Clients and models are labeled so that the diagonal represents the  
 1770 true client to model/cluster assignment. Numbers on the heatmap express the reconstruction loss  
 1771 (multiplied by a constant for the sake of presentation) for each client testset and model combination.  
 1772 Very high values are hidden and the corresponding cells are colored red. For the local case, we only  
 1773 show 10 of the 50 models, the ones that correspond to the selected clients.  
 1774

1775 We observe that in all cases, the reconstruction loss of each client’s assigned model by *CLoVE* (the  
 1776 diagonal) is the smallest in the client’s row (i.e. the smallest among all models), which confirms that  
 1777 the models *CLoVE* assigned to each client are actually trained to reconstruct the client’s data (MNIST  
 1778 digit) well, and not reconstruct well the other digits. On the contrary, *IFCA* gets stuck and assigns  
 1779 all clients to 3 clusters (models) – the first 3 columns. The rest of the models remain practically  
 1780 untouched (untrained), hence their high (red) loss values on all clients’ data. The 3 chosen models  
 1781 are not good for the reconstruction of MNIST digits, because they don’t achieve the lowest loss on  
 their respective digits.

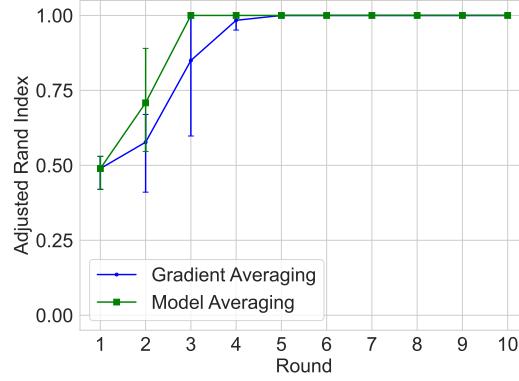


1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794 Figure 5: Heatmaps of reconstruction loss of MNIST digits by models clustered and trained by each  
1795 algorithm

1796  
1797  
1798 Comparing *CLoVE* and *FedAvg*, we see that for each row (client from each cluster), the loss achieved  
1799 by *CLoVE* is lower than that achieved by *FedAvg*. This demonstrates the **benefit of our personalized**  
1800 **FL approach over traditional FL**: *FedAvg* uses a single model to train on local data of each client,  
1801 and aggregates updated parameters only for this single model. *CLoVE*, on the other hand, clusters  
1802 the similar clients together and trains a separate model on the data of clients belonging to each  
1803 cluster, thus effectively learning – and tailoring the model to – the cluster’s particular data distribution.  
1804 Finally, comparing *CLoVE* with *Local-only*, we see that the diagonal losses achieved by *CLoVE* on  
1805 the selected models for each client are lower than the corresponding diagonal losses by *Local-only*.  
1806 This shows the benefit of federated learning, as *CLoVE* models are trained on data from multiple  
1807 clients from the same cluster, unlike *Local-only* where clients only train models based on their limited  
1808 amount of private data.

#### 1810 D.2.2 MODEL AVERAGING VERSUS GRADIENT AVERAGING

1811  
1812 In this experiment, we compare the two variants of our *CLoVE* algorithm, i.e. *model averaging* and  
1813 *gradient averaging*. We run an MNIST experiment with  $K = 10$ ,  $M = 50$  and  $n = 100$ , using both  
1814 averaging methods. The results are shown in Fig. 6. We observe that *CLoVE* is able to arrive at the  
1815 correct clustering with both methods, with gradient averaging needing slightly more rounds to do so  
1816 than model averaging.



1817  
1818  
1819 Figure 6: Clustering accuracy (ARI) of *CLoVE* over time (rounds) for MNIST for model averaging  
1820 vs gradient averaging.

1836 D.3 SCALING WITH THE NUMBER OF CLIENTS  
1837

1838 In this experiment, we examine the effect of increasing the number of clients. We use 10 clusters,  
1839 and 100-200 datapoints per client. The results for MNIST, CIFAR-10 and FMNIST in the supervised  
1840 setting (using CNNs) are shown in Table 15, and for MNIST in the unsupervised setting (using  
1841 autoencoders) are shown in Table 16. In the supervised setting, we see that *CLoVE*'s final ARI is  
1842 always perfect, and both the ARI and the test accuracy of *CLoVE* consistently exceed those of *IFCA*,  
1843 no matter how big the number of clients is. Similarly, in the unsupervised setting we observe that  
1844 *IFCA* consistently achieves a bad ARI, while *CLoVE* always manages to reach perfect clustering  
1845 accuracy and scales successfully when the number of clients is large.

1846 Table 15: Scaling with the number of clients – Supervised setting

Algorithm	# of clients $M$	Final ARI			Test accuracy		
		MNIST	CIFAR-10	FMNIST	MNIST	CIFAR-10	FMNIST
IFCA	50	0.92 ± 0.11	0.85 ± 0.11	0.74 ± 0.22	99.2 ± 0.8	83.7 ± 5.3	96.8 ± 2.3
	100	0.85 ± 0.11	0.65 ± 0.25	0.93 ± 0.11	99.1 ± 0.5	72.7 ± 12.7	97.4 ± 2.6
	250	0.87 ± 0.18	0.85 ± 0.10	0.85 ± 0.10	99.4 ± 0.7	80.9 ± 7.7	98.6 ± 0.5
	500	0.93 ± 0.10	0.85 ± 0.10	0.85 ± 0.10	97.1 ± 2.9	77.3 ± 5.8	89.2 ± 7.7
	1000	0.85 ± 0.10	0.93 ± 0.10	0.75 ± 0.21	98.3 ± 1.0	78.0 ± 7.2	92.0 ± 3.6
CLoVE	50	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	99.8 ± 0.0	90.8 ± 0.2	99.1 ± 0.0
	100	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	99.8 ± 0.0	90.3 ± 0.6	99.3 ± 0.0
	250	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	99.9 ± 0.0	91.2 ± 1.5	99.3 ± 0.0
	500	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	99.6 ± 0.1	82.3 ± 1.5	95.3 ± 4.4
	1000	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	99.6 ± 0.1	84.6 ± 1.2	96.5 ± 3.0

1859  
1860 Table 16: Scaling with the number of clients – Unsupervised setting

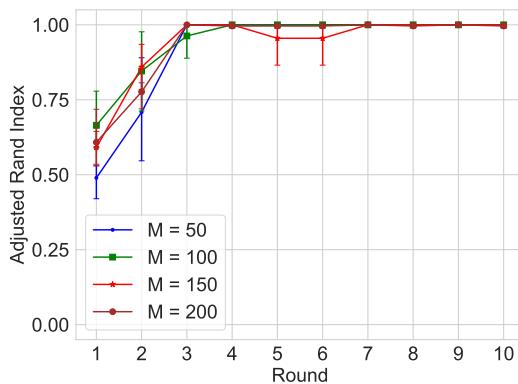
Algorithm	# of clients $M$	Test loss	Final ARI
IFCA	100	0.055 ± 0.001	0.091 ± 0.064
	200	0.055 ± 0.001	0.079 ± 0.060
	500	0.059 ± 0.002	0.210 ± 0.149
	1000	0.067 ± 0.002	0.186 ± 0.137
CLoVE	100	0.065 ± 0.001	1.000 ± 0.000
	200	0.063 ± 0.000	1.000 ± 0.000
	500	0.065 ± 0.000	1.000 ± 0.000
	1000	0.073 ± 0.000	1.000 ± 0.000

1890  
1891

## D.4 CONVERGENCE DYNAMICS

1892  
1893  
1894  
1895

We now show the convergence dynamics of *CLoVE* for MNIST in the unsupervised setting. We use  $K = 10$  clusters, and  $n = 100$  datapoints per client. We vary the number of clients from 50 to 200. The behavior of the ARI can be seen in Figure 7. We observe that *CLoVE* reaches a perfect ARI in just a few rounds, for all values of the number of clients  $M$ .

1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909

1910

Figure 7: Clustering accuracy (ARI) of *CLoVE* over time (rounds) for MNIST for different numbers of clients

1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924

## D.5 ABLATION STUDIES

We analyze the impact of changing some of key components of *CLoVE* one at a time. This includes:

- **No Matching:** Order the clusters by their smallest client ID and assign models to them in that order (i.e. cluster  $k$  gets assigned model  $k$ ). This is in lieu of using the bi-partite matching approach of Alg. 1.
- **Agglomerative clustering:** Use a different method than k-means to cluster the loss vectors.
- **Square root loss:** Cluster vectors of square roots of model losses instead of the losses themselves.

1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934

Our experiments are for image classification for the CIFAR-10 dataset using 25 clients that are partitioned uniformly among 5 clusters using a label skew with very large overlap. The clients of the first four cluster are assigned samples from every class except one class which is chosen to be different for each cluster. The clients of the fifth cluster are assigned samples from all classes. The partitioning of the data of each class among the clients that are assigned to that class is by sampling from a Dirichlet( $\alpha$ ) distribution with  $\alpha = 0.5$ . The test accuracy achieved is shown in Table 17. We find that “No Matching” has the most impact, as it reduces the test accuracy by almost  $4.3\% \pm 1\%$ . Switching to Agglomerative clustering results in a slight improvement in the mean test accuracy, but the result is not statistically significant. Likewise, square root loss results in no significant change in performance. This shows that *CLoVE* is robust in its design.

1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

Table 17: Ablation studies

	Test accuracy
Original CLoVE	$58.3 \pm 1.4$
No matching	$55.8 \pm 0.6$
Agglomerative clustering	$59.3 \pm 1.5$
Square root loss	$58.1 \pm 3.3$