# INDEX-OVERLAP SIMILARITY: A VALUE-FREE PROXY FOR MODEL RELATEDNESS

**Anonymous authors**Paper under double-blind review

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

023

025 026

027

029

031

033

034

037

040

041

042

043

044

046

047

048

051

052

## **ABSTRACT**

Measuring client relatedness is central to clustering and personalization in federated learning (FL), but value-based similarities over full weights or gradients are bandwidth-heavy and leak information. We propose Index- $Overlap\ Similarity\ (IOS)$ , a value-free metric that represents each client by the indices of its Top-K salient parameters and scores pairs by the normalized overlap of these supports. We show why IOS preserves alignment: under head-dominance with bounded dispersion, it lower-bounds cosine up to tail error; Top-K is invariant to common layerwise rescalings; and exponential moving averages stabilize supports across rounds. We instantiate IOS for clustered personalized FL, neighbor selection, donor ranking, and oracle distribution alignment. Across FMNIST, CIFAR-10/100, and 20News under Dirichlet and pathological splits, IOS matches or exceeds cosine/Euclidean while sharing only indices. IOS is a simple, scalable primitive for similarity search under communication and privacy constraints.

# 1 Introduction

A fundamental component of machine learning is the comparison of high-dimensional objects, such as model weights, gradients, and data embeddings. It drives important federated learning tasks like personalized aggregation and client clustering (Ghosh et al., 2020; Fallah et al., 2020; Dinh et al., 2020), continual-learning tools like memory retrieval and drift detection (Gama et al., 2014; Lu et al., 2018), and model analysis tools for provenance and similarity search (Indyk & Motwani, 1998; Andoni & Indyk, 2008). Nevertheless, three enduring issues make default full-vector cosine/Euclidean comparisons debilitating at scale and under privacy constraints: (1) Computation/communication scale with model size. Comparing full real-valued vectors requires moving and multiplying arrays whose length equals the number of trainable parameters; even on-device CNNs (1-10M) strain bandwidth at scale, while mid-size backbones like ResNet-18/50 (~11M/~25M) and ViT-B (~80–90M) push per-client payloads into tens-hundreds of MB per round; transformers exacerbate this—BERT-base ( $\sim$ 110M), BERT-large ( $\sim$ 340M), and multi-billion-parameter checkpoints (He et al., 2016; Dosovitskiy et al., 2021; Devlin et al., 2019). (2) Sharing real values enables reconstruction/inference attacks. Logits and partial activations permit model inversion (Fredrikson et al., 2015); repeated round exposures fuel membership and property inference (Shokri et al., 2017; Melis et al., 2019); and gradients/updates can be inverted to recover inputs or labels (Zhu et al., 2019; Geiping et al., 2020). (3) Numeric instability distorts geometry. Layer-wise scale heterogeneity (batch normalization, weight decay, mixed precision) and optimizer-state drift yield poorly calibrated cosine/Euclidean distances on raw weights or gradients across clients and over time (Ioffe & Szegedy, 2015; Loshchilov & Hutter, 2019; Micikevicius et al., 2018).

In over-parameterized networks, salience is predominantly concentrated in a few heads and remains relatively stable: a small subset of model parameters holds most first-order significance, while the long tail is noisy and less indicative of inter-client relatedness (Michel et al., 2019; Voita et al., 2019; Li et al., 2017; Gale et al., 2019). Additionally, models trained on datasets drawn from similar distributions exhibit similar parameter-importance patterns, consistent with representational-similarity findings and client-relatedness in FL (Kornblith et al., 2019; Raghu et al., 2017; Ghosh et al., 2020). Building on this intuition, we choose an alternative approach: we represent each model by the significance of its coordinates rather than their values. Specifically, for each client, we assess the significance of trained model parameters using a diagonal Fisher proxy and thereafter determine the indices of a small set of top-K salient parameters. We calculate similarity as the overlap between

these index sets. This Index-Overlap Similarity (IOS) is deliberately devoid of value; it conveys solely integer identifiers of prominent coordinates. The intersection of coordinates tends to correlate with alignment of learning signals and, importantly, remains resilient to layer-wise rescaling and optimizer peculiarities. The set size K is selected to be a minuscule proportion of the model size, specifically  $K \ll M$  where M denotes the number of trainable parameters. In fact, maintaining only a small fraction of coordinates preserves the majority of the stable "head" structure while reducing computational requirements and data size by more than an order of magnitude, and without revealing real-valued weights or gradients. Only index sets of salient coordinates are shared; no weights, gradients, logits, or activations are exposed. We make no formal privacy claims, but IOS reduces the attack surface relative to value sharing, mitigating risks of gradient/model inversion and membership/property inference (Zhu et al., 2019; Geiping et al., 2020; Shokri et al., 2017).

**Positioning vs. prior work.** The resemblance among trained models is fundamental to numerous machine learning procedures beyond FL. Most current measures function based on values: cosine/Euclidean metrics applied to weights or gradients; prototype/feature similarities from penultimate activations (e.g., FedProto) or representation metrics such as CKA (Tan et al., 2022; Kornblith et al., 2019); and influence/Shapley-style utilities obtained from value-bearing surrogates (Koh & Liang, 2017; Ghorbani & Zou, 2019; Jia et al., 2019). Sparsification and pruning techniques (Lin et al., 2018; Lee et al., 2019; Evci et al., 2020; Han et al., 2015) either learn or enforce sparse *parameters* for enhanced efficiency and occasionally examine *mask overlap* for stability; however, they regard overlap merely as a byproduct of pruning rather than a fundamental similarity primitive. *IOS* is, to our knowledge, the first method that calculates cross-model similarity *without transmitting any parameter values*, facilitating sketch-based scaling and minimizing leakage channels.

Our Contribution. We instantiate IOS for FL, where communication and privacy constraints make dispensing with real values attractive: every stage operates solely on prominent indices, mitigating reconstruction/linkage risks from value sharing. We benchmark IOS against cosine and Euclidean baselines across four applications: (i) Clustered Personalized FL (CPFL): derive an IOS affinity matrix for client clustering to train cluster-specific models, evaluating against cosine-based clustering via FL accuracy. (ii) Neighbor selection for personalized aggregation: build a similarity-weighted neighbor graph, form label-histogram mixtures to match each client, comparing divergence from target labels; (iii) Shapley-style donor ranking: use similarity as a proxy for marginal utility (validation uplift) and assess rank agreement/top-k recall versus KNN-Shapley; (iv) Oracle distribution alignment: test whether the similarity matrix tracks true relatedness from clients' label-distribution divergence using Spearman/Kendall- $\tau$ , noting those baselines share real values while IOS does not.

Our Findings. Across FMNIST, CIFAR-10/100, and 20News under Dirichlet/Patho splits, *IOS* (indices-only) consistently wins in our two target applications. In CPFL it yields *better* accuracy than Cosine/Euclidean (avg +1.5 pp vs. Cosine); for neighbor selection it recovers more oracle neighbors (CIFAR-100, Dir(0.1): R@8=0.67 vs. 0.61/0.53). *IOS* also gives tighter distribution alignment and donor ranking (CIFAR-10, Dir(0.1): JS 0.219 vs. 0.238/0.252; Kendall- $\tau$  0.48, R@5 0.62 vs. 0.36/0.55).

# 2 BACKGROUND & PROBLEM SETUP

#### 2.1 BACKGROUND

Federated Learning (FL). We consider n clients with private datasets  $D_i = \{(x_d, y_d)\}_{d=1}^{|D_i|}$  and a shared model parameterization  $w \in \mathbb{R}^M$ . The canonical FL objective is the weighted empirical risk  $F(w) = \sum_{i=1}^n p_i \, F_i(w), \quad F_i(w) = \frac{1}{|D_i|} \sum_{(x,y) \in D_i} \ell(w;x,y), \quad p_i = \frac{|D_i|}{\sum_j |D_j|}.$  In round t, a subset  $S_t$  of clients receives the current global model  $w^{(t)}$ , performs E local SGD steps  $w_i^{(t+1)} \leftarrow w^{(t)} - \eta \sum_{e=1}^E \widehat{\nabla} F_i(w_i^{(t,e-1)})$  and returns  $w_i^{(t+1)}$  to the server. The server aggregates (FedAvg)  $w^{(t+1)} = \sum_{i \in S_t} \bar{p}_i \, w_i^{(t+1)}, \qquad \bar{p}_i = \frac{|D_i|}{\sum_{j \in S_t} |D_j|},$  optionally using update-form aggregation (on  $w_i^{(t+1)} - w^{(t)}$ ) and secure aggregation. We assume a fixed architecture across clients and standard non-IID partitions unless stated otherwise.

**Fisher information–based importance.** For a probabilistic model  $p_{\theta}(y \mid x)$  with loss  $\ell(\theta; x, y) = -\log p_{\theta}(y \mid x)$ , the Fisher information matrix (FIM) is

$$\mathcal{I}(\theta) = \mathbb{E}_x \Big[ \mathbb{E}_{y \sim p_{\theta}(\cdot \mid x)} \big[ \nabla_{\theta} \log p_{\theta}(y \mid x) \nabla_{\theta} \log p_{\theta}(y \mid x)^{\top} \big] \Big].$$

Its diagonal provides a principled, nonnegative per-parameter importance and underpins natural-gradient methods (Amari, 1998). In practice we use the *empirical Fisher* (EF): replace the model expectation with observed labels and estimate  $\operatorname{diag}(\mathcal{I})$  from squared per-sample gradients over minibatches; EF is convenient but can deviate from the true Fisher and may misrepresent second-order geometry (Kunstner et al., 2019). Recent results (Soen & Sun, 2024) give variance bounds and sample-complexity trade-offs for diagonal Fisher estimators, with variance governed by network nonlinearity and parameter grouping. Complementarily, improved EF (iEF) applies diagonal scaling to better approximate natural-gradient behavior while retaining EF's simplicity (Wu et al., 2024). In this work, we adopt diagonal-Fisher / gradient-second-moment surrogates and mitigate estimator noise via mini-batch averaging with an optional EMA.

Client heterogeneity and importance-pattern divergence. Let  $g(w;x,y) = \nabla_w \ell(w;x,y)$  and define a per-parameter importance proxy via the empirical second moment / Fisher diagonal  $s_j = \mathbb{E}[(\partial \ell/\partial w_j)^2]$ . Under a label-mixture model,  $\mathbb{E}_{(x,y)\sim D_i}[g(w;x,y)] = \sum_c \pi_i(c)\,\mu_c(w) + \xi_i$ , so shifts in  $\pi_i$  alter the mean gradient and the induced importance profile. Non-IID data induce gradient dissimilarity, driving weight divergence and slowing FedAvg; SCAFFOLD formalizes bounded dissimilarity and shows drift correction improves convergence (Karimireddy et al., 2020). Empirically, divergence correlates with class-distribution distance (e.g., EMD) (Zhao et al., 2018), and heterogeneous clients exhibit update directions/norms that differ markedly (Wang et al., 2023). Conversely, related distributions exhibit shared task-relevant structure: sparse "winning tickets" transfer across natural-image datasets (Morcos et al., 2019), and fine-tuning on related tasks yields compatible weight-space task vectors that compose (Ilharco et al., 2023). Thus, dissimilar data induce divergent importance/gradient profiles, whereas similar data induce partially overlapping sets of salient parameters—largely independent of the chosen similarity metric.

## 2.2 ASSUMPTIONS AND SCOPE

We adopt an honest-but-curious FL coordinator. All clients share the same parameterization (layer order/tensor shapes), yielding a common index space; cross-architecture similarity is out of scope. The protocol releases only index sets —no real-valued weights, gradients, logits, activations, or example-level metadata. Systems assumptions include authenticated transport and optional secure aggregation; robustness to Byzantine failures/poisoning is orthogonal. Label distributions and oracle similarities are used only for evaluation, never at runtime. Relative to value sharing, *IOS* reduces exposure to reconstruction channels. Formal privacy accounting for index release is beyond scope; the aim is to shrink the attack surface versus sharing values.

# 3 METHOD

IOS is a *value-free*, *index-only* similarity framework: each client computes a local importance signal (e.g., diagonal Fisher/gradient second moment), extracts the Top-K parameter *indices* as its support, and shares only these indices (or compact MinHash signatures). Cross-client similarity is defined by set overlap on supports, enabling exact intersections or scalable LSH-based retrieval—without transmitting any real-valued weights, gradients, or activations. The support size K is chosen locally via importance coverage and stability under communication/privacy budgets.

# 3.1 INDEX-OVERLAP SIMILARITY (IOS)

**Importance accumulation.** For each client i, we form a nonnegative per-parameter importance vector  $s_i \in \mathbb{R}^M_{\geq 0}$  from local data  $D_i$  using the diagonal Fisher (or its empirical second-moment proxy) introduced in § Background:

$$s_{i,j} \approx \frac{1}{T} \sum_{t=1}^{T} \frac{1}{B} \sum_{b=1}^{B} \left( \frac{\partial \ell(w; x_{i,b}^{(t)}, y_{i,b}^{(t)})}{\partial w_{j}} \right)^{2},$$

163

164

166

167

168 169

170

171

172 173

174 175

176 177

178

179

181

182

183

185

186

187

188

189

190

191

192 193 194

195 196

197

200

201

202203204

205

206

207

208209

210211212

213

214

215

optionally stabilized by an EMA  $s_i \leftarrow \beta s_i + (1-\beta)\hat{s}_i$  with  $\beta \in [0,1)$ . This produces a scale-robust, value-nonnegative signal that can be computed entirely on the device.

**Support extraction and Similarity Calculation.** Given a global budget  $K \ll M$ , we define the client's *index-only* representation as the set of its K most salient parameters

$$I_i = \text{TopK}(s_i) \subseteq [M], \qquad |I_i| = K.$$

(Per-layer budgeting is accommodated by choosing  $K_\ell$  with  $\sum_\ell K_\ell = K$  and  $I_i = \bigcup_\ell I_i^{(\ell)}$  where  $I_i^{(\ell)} = \operatorname{TopK}_\ell(s_i^{(\ell)})$ .) Ties are broken deterministically (e.g., by index) for reproducibility. For implementation, we also use a bitmask  $m_i \in \{0,1\}^M$  with  $m_{i,j} = \mathbf{1}\{j \in I_i\}$ . IOS defines similarity purely from set overlap, without transmitting any real-valued weights/gradients. We define similarity S(i,j) (with distance 1- similarity). With  $|I_i| = |I_j| = K$ ,  $S(i,j) = \frac{|I_i \cap I_j|}{K}$ .

# Algorithm 1 Select K via Importance Coverage (client i)

```
Input: Importance s_i \in \mathbb{R}^M_{\geq 0}; coverage target \tau; cap K_{\max}; stability target \rho_0 and resamples r
Output: K_i^{\star} and support \bar{I_i} = \text{TopK}(s_i, K_i^{\star})
 1: Sort indices j_1,\ldots,j_M by s_{i,j} descending; prefix sums S(t)=\sum_{u=1}^t s_{i,j_u} and S_{\text{tot}}=S(M) 2: lo\leftarrow 1,\ hi\leftarrow K_{\text{max}},\ K^\star\leftarrow K_{\text{max}}
 3: while lo \leq hi do
           mid \leftarrow |(lo + hi)/2|; \ C \leftarrow S(mid)/S_{tot}
           if C \geq \tau then
 5:
                 \overline{K^{\star}} \leftarrow mid; \ hi \leftarrow mid - 1
                                                                                        \triangleright keep smallest K achieving coverage
 6:
 7:
 8:
                 lo \leftarrow mid + 1
 9: if stability target \rho_0 is provided then
           for K \in \{K^*, K^*+1, \min(K^*+2, K_{\max})\} do
10:
                 Estimate \rho_i(K) using r lightweight resamples of s_i by Equation (1)
11:
                 if \rho_i(K) \geq \rho_0 then
12:
                      \hat{K}^{\star} \leftarrow \hat{K}; break
```

# 3.2 SELECTING K FOR IOS

14: **return**  $K_i^{\star} \leftarrow K^{\star}$  and  $I_i = \{j_1, \dots, j_{K^{\star}}\}$ 

IOS represents each client i by the indices of its K most important parameters. Choosing K must balance utility (capturing enough importance mass), stability (insensitivity to estimator noise), and budgets from communication and privacy. Let  $s_i \in \mathbb{R}^M_{\geq 0}$  be the local importance vector (e.g., diagonal Fisher or per-parameter gradient second moment) on client i. Let  $j_1, \ldots, j_M$  be indices sorted by  $s_{i,j}$  in descending order and define the cumulative importance (using  $\ell_1$  mass for Fisher-diagonal):

$$C_i(K) = \frac{\sum_{t=1}^K s_{i,j_t}}{\sum_{t=1}^M s_{i,j_t}} \in [0,1], \qquad C_i(K) \text{ is non-decreasing in } K.$$

To assess robustness of the selected support, we define an *overlap-stability* score. Generate r resamples of the importance estimator (e.g., via bootstrapping mini-batches or adjacent time windows). For each resample b, extract the top-K set  $I_i^{(b)}(K) = \operatorname{TopK}(s_i^{(b)})$ . With  $|I_i^{(b)}(K)| = K$  for all b, define the mean pairwise overlap

$$\rho_i(K) = \frac{2}{r(r-1)} \sum_{b \le b'} \frac{\left| I_i^{(b)}(K) \cap I_i^{(b')}(K) \right|}{K} \in [0, 1]. \tag{1}$$

Thus,  $\rho_i(K)$  quantifies how consistently the same indices appear across resamples; values near 1 indicate stable supports. Given a coverage target  $\tau \in (0,1)$  and an optional stability target  $\rho_0 \in (0,1)$ , select the smallest integer K satisfying

$$C_i(K) \geq \tau$$
 and (if enforced)  $\rho_i(K) \geq \rho_0$ .

Here  $C_i(K)$  guarantees that the selected indices cover a desired fraction of total importance, while  $\rho_i(K)$  ensures reproducibility under estimator noise. The rule is locally computable, value-free externally, and binary-searchable because  $C_i(K)$  is monotone. Stability is only checked near the candidate K to limit overhead. Ties in TopK are broken deterministically by index.

Exact all-pairs overlap costs  $O(n^2K)$ ; we therefore use MinHash-LSH to retrieve candidate neighbors in subquadratic time while preserving *IOS*'s value-free property. Full derivations, and complexity bounds appear in Appx. A.

# 4 THEORETICAL PROPERTIES

We develop the theory of IOS around three pillars: (i) alignment—why overlap of top-K salient indices tracks the cosine similarity of first-order signals; (ii) stability—why the selected supports remain consistent over time under mild noise with EMA smoothing; and (iii) robustness—why IOS is insensitive to common re-scalings and diagonal preconditioning. Additional analyses (coverage–stability trade-off for choosing K, robustness, and complexity) appear in Appx. D.

## 4.1 NOTATION AND ASSUMPTIONS

For client i, let  $s_i \in \mathbb{R}^M_{\geq 0}$  be a nonnegative importance vector, and  $I_i = \operatorname{TopK}(s_i)$  with  $|I_i| = K \ll M$ . Write the head-tail split  $s_i = h_i + t_i$  with  $(h_i)_u = (s_i)_u \, \mathbf{1}\{u \in I_i\}$  and  $t_i = s_i - h_i$ . Define the head-energy fraction  $\alpha_i = \|h_i\|_2^2/\|s_i\|_2^2 \in (0,1]$  and the tail-to-head ratio  $\tau_i = \|t_i\|_2/\|h_i\|_2$ . We say head dominance holds if  $\alpha_i \geq 1 - \varepsilon$  (equivalently,  $\|h_i\|_2 \geq \sqrt{1-\varepsilon} \, \|s_i\|_2$ ) for a small  $\varepsilon \in [0,1)$ . We assume bounded dispersion inside the head:  $\kappa_i := \max_{u \in I_i} (h_i)_u / \min_{u \in I_i} (h_i)_u \leq \kappa$  for a moderate  $\kappa \geq 1$ . For two clients i,j, denote the normalized overlap of their salient supports by  $R_{ij} = \frac{|I_i \cap I_j|}{K} \in [0,1]$ .

#### 4.2 ALIGNMENT: COSINE VS. *IOS*

**Proposition 1** (Cosine lower bound via overlap). *Under head dominance*  $(\alpha_i, \alpha_j \geq 1 - \varepsilon)$  *and bounded dispersion*  $(\kappa_i, \kappa_j \leq \kappa)$ ,

$$\cos(s_i, s_j) = \frac{\langle s_i, s_j \rangle}{\|s_i\|_2 \|s_i\|_2} \ge \frac{(1 - \varepsilon)}{\kappa^2} R_{ij}. \tag{2}$$

**Proof sketch.** With nonnegative entries, the inner product is at least the contribution on the intersecting head block:  $\langle s_i, s_j \rangle \geq \sum_{u \in I_i \cap I_j} (h_i)_u (h_j)_u$ . Let  $a_i = \min_{u \in I_i} (h_i)_u$  and  $b_i = \max_{u \in I_i} (h_i)_u \leq \kappa a_i$ . Then  $\|h_i\|_2^2 \leq K b_i^2 \leq K \kappa^2 a_i^2$ , hence every head entry satisfies  $(h_i)_u \geq \|h_i\|_2/(\kappa \sqrt{K})$ . The intersect term is thus at least  $\frac{|I_i \cap I_j|}{\kappa^2 K} \|h_i\|_2 \|h_j\|_2$ . Head dominance yields  $\|h_i\|_2 \geq \sqrt{1-\varepsilon} \|s_i\|_2$  and likewise for j, which gives equation 2 after normalization. See full proof in Appx. B

## 4.3 SUPPORT STABILITY OVER TIME

Let  $s_i^{(t)} = \mu_i + \xi_i^{(t)}$ , where  $\mu_i$  is a stationary signal and  $\xi_i^{(t)}$  has independent, mean-zero sub-Gaussian coordinates with proxy variance  $\sigma_i^2$ . *IOS* maintains an exponential moving average (EMA),  $\tilde{s}_i^{(t)} = \beta \, \tilde{s}_i^{(t-1)} + (1-\beta) \, s_i^{(t)}, \quad \beta \in [0,1)$  and selects  $I_i^{(t)} = \mathrm{TopK}(\tilde{s}_i^{(t)})$  (after an arbitrary burn-in). The K-boundary margin is  $\Delta_i := \mu_{i,(K)} - \mu_{i,(K+1)} > 0$  (ties w.r.t.  $\mu_i$  are broken deterministically). The EMA reduces the per-coordinate noise variance to

$$\sigma_i(\beta)^2 = \sigma_i^2 \sum_{s \ge 0} (1 - \beta)^2 \beta^{2s} = \sigma_i^2 \frac{(1 - \beta)^2}{1 - \beta^2} = \sigma_i^2 \frac{1 - \beta}{1 + \beta}.$$
 (3)

**Theorem 1** (Top-K selection stability). Fix client i and let  $\Delta_i > 0$  as above. For any t after burnin,  $\Pr\left(\operatorname{TopK}(\tilde{s}_i^{(t)}) \text{ is unique and equals } \operatorname{TopK}(\mu_i)\right) \geq 1 - 2M \exp\left(-\frac{\Delta_i^2}{8\,\sigma_i(\beta)^2}\right)$ . Equivalently, if  $\Delta_i \geq c\,\sigma_i(\beta)\sqrt{\log M}$  with any  $c > \sqrt{8}$ , then the failure probability decays as  $O(M^{1-c^2/8})$ .

**Proof idea.** Let H be the mean top-K set and  $\bar{H}=[M]\setminus H$ . Consider  $E_1=\{\min_{u\in H}\tilde{g}_{i,u}^{(t)}\geq \mu_{i,(K)}-\Delta_i/2\}$  and  $E_2=\{\max_{v\in \bar{H}}\tilde{g}_{i,v}^{(t)}\leq \mu_{i,(K+1)}+\Delta_i/2\}$ . Sub-Gaussian tails and a union bound over K and (M-K) coordinates give  $\Pr(E_1^c)\leq Ke^{-\Delta_i^2/(8\sigma_i(\beta)^2)}$  and  $\Pr(E_2^c)\leq (M-K)e^{-\Delta_i^2/(8\sigma_i(\beta)^2)}$ ; hence  $1-2Me^{-\Delta_i^2/(8\sigma_i(\beta)^2)}$  overall. Define the temporal self-overlap  $\Gamma_i^{(t)}=\frac{|I_i^{(t)}\cap I_i^{(t-1)}|}{|I_i^{(t)}\cup I_i^{(t-1)}|}\in [0,1]$ . If a fraction  $\rho$  of the K-boundary gaps of  $\mu_i$  exceed  $c\,\sigma_i(\beta)\sqrt{\log M}$  (with  $c>\sqrt{8}$ ), then after burn-in  $\mathbb{E}[\Gamma_i^{(t)}]\geq \rho-O(M^{1-c^2/8})$ . See full proof in Appx. C.

# 5 APPLICATIONS OF *IOS*

We instantiate *IOS* in four FL scenarios. Where prior art uses cosine over value-bearing vectors, we implement the same pipelines and *swap cosine with* IOS to ensure fair comparisons.

- Clustered Personalized FL (CPFL / IFCA-style). We construct an affinity matrix  $A_{ij}$  and run clustering with affinity propagation (AP) (Frey & Dueck, 2007) to obtain client groups  $\{C_1,\ldots,C_K\}$ , followed by cluster-conditioned training. As baselines, we re-implement IFCA (Ghosh et al., 2020) and Clustered FL (Sattler et al., 2020) using their original cosine-based affinities; the IOS variant replaces cosine with value-free overlap of Top-K supports,  $A_{ij}^{(IOS)} = \frac{|I_i \cap I_j|}{K}$ . We report FL accuracy and clustering quality metrics—highlighting regimes where IOS matches accuracy while reducing leakage and bytes.
- Neighbor Selection for Personalized Aggregation (Per-FedAvg + similarity graph). Personalized FL updates are mixed from "nearby" clients via a similarity-weighted graph G with edges  $w_{ij} \propto \widehat{S}(i,j)$ . We instantiate Per-FedAvg (Fallah et al., 2020) with a cosine-based neighbor policy (baseline) and a drop-in IOS policy that computes similarity from the overlap of important indices. Metrics include recall compared to cosine for finding the most relevant neighbors.
- Oracle Distribution Alignment (evaluation-only). To test whether an index-only similarity captures latent relatedness induced by label distributions, we form an oracle  $S^*(i,j) = 1 JS(\pi_i, \pi_j)$  (or 1 Hellinger) from client histograms  $\{\pi_i\}$  and compare it against method-driven similarities computed either by cosine on value-bearing vectors (weights/updates/features) or by *IOS*. We evaluate rank alignment (Spearman/Kendall) and calibration across divergence bins.
- Shapley-Style Donor Ranking (similarity  $\to$  utility proxy). For a target client t, we rank donors  $j \neq t$  via a utility proxy  $\widetilde{u}_t(j) \propto \widehat{S}(t,j)$ . We reproduce kNN-Shapley (Jia et al., 2019) using cosine(Koh & Liang, 2017; Ghorbani & Zou, 2019) in the neighbor stage (baseline), then replace it with IOS:  $\widehat{S}_{IOS}(t,j) = \frac{|I_t \cap I_j|}{K}$ . We report rank correlation with true uplift  $\Delta \mathcal{V}_t(j)$  and top-k recall, noting when IOS maintains fidelity while avoiding value sharing.

# 6 EXPERIMENTS

We evaluate *IOS* on standard FL applications using identical client partitions and architectures—only the similarity differs; extended empirical results appear in Appx. E.

#### 6.1 **SETTING**

**Datasets and models.** Vision benchmarks include FMNIST ( $28 \times 28$  grayscale; 10 classes) with a 2-layer CNN, CIFAR-10 ( $32 \times 32$  RGB; 10 classes) with ResNet18, and *CIFAR-100* ( $32 \times 32$  RGB; 100 classes) with ResNet50. To probe modality-agnostic behavior, we include an optional non-vision task on 20 News with BERT-base. Across clients, the per-client optimal selection size K ( $K_{max} = 20\%$ ) lies in  $K_i^*$  as [10.3, 13.5]% for CNN, [6.1, 8.4]% for ResNet18, [8.2, 10.4]% for ResNet50, and [7.6, 9.3]% for BERT-base; we set K to the client-wise mean in each case, yielding K = 12.2, 7.3, 9.6%, and 8.2%, respectively.

**Data partitioning.** We synthesize heterogeneous populations with three regimes. IID partitions split each dataset uniformly at random across N clients. Dirichlet partitions draw client-wise class proportions from  $\mathrm{Dir}(\alpha)$  with  $\alpha \in \{0.1, 0.3\}$  and allocate examples accordingly. Pathological partitions, Patho(n), assign each client a small subset of classes (e.g.,  $n \in \{20\%, 30\%\}$  classes per

324 325

Table 1: CPFL test accuracy (%). Rows are models; columns are distributions with three similarity choices: Cosine, Euclidean, and *IOS*. Bold marks the best.

333334335336

332

341 342 343

356

357

358

350 351

371

372

373

374 375

376

377

Dir(0.3)Dir(0.1)Patho(20%)Patho(30%)Euclidean Euclidean Euclidean Euclidean Cosine Cosine Cosine Cosine SOI SOI SOI Model / Dataset CNN / FMNIST 73.32 71.41 76.23 70.80 70.21 74.79 94.87 91.92 95.18 82.19 81.80 83.15 ResNet18 / CIFAR-10 63.92 61.35 66.13 60.85 58.40 62.1 79.18 75.47 81.44 74.38 73.96 74.81 ResNet50 / CIFAR-100 47.25 41.56 45.40 43.18 48.3 58.12 53.19 53.49 50.06 49.71 58.07 52.98 BERT-base / 20News 49.16 48.04 50.3 44.62 41.51 45.41 56.18 55.03 58.26 53.79 51.80 54.25

Table 2: Cluster quality across heterogeneity regimes. Higher Silhouette/CH, lower DB are better.

Model / Dataset	Method	l	Dir(0	.3)	l	Dir(0	0.1)	P	atho(2	20%)	Patho(30%)		
		Sil↑	DB↓		Sil↑	DB↓	CH↑	Sil↑	DB↓	CH↑	Sil↑	DB↓	CH↑
CNN/FMNIST	Oracle	0.259	1.049	6.490 —	0.400	0.855	12.714 —	0.308	0.921	7.040 —	0.277	1.356	7.986
	IOS	0.212	1.208	5.427 —	0.340	0.993	10.486 —	0.249	1.059	5.701 —	0.234	1.617	6.547
	Cosine	0.207	1.250	5.245 —	0.333	1.016	10.140 —	0.241	1.087	5.569 —	0.226	1.673	6.346
	Euclidean	0.197	1.278	5.074 —	0.321	1.044	9.758 —	0.229	1.120	5.442 —	0.215	1.725	6.114
ResNet18/CIFAR-10	Oracle	0.218	1.185	5.746 —	0.384	0.801	10.361 —	0.373	0.800	9.328 —	0.277	1.356	7.987
	IOS	0.176	1.399	4.636 —	0.322	0.925	8.663 —	0.305	0.944	7.849 —	0.226	1.617	6.479
	Cosine	0.171	1.447	4.539 —	0.311	0.951	8.426 —	0.294	0.966	7.633 —	0.221	1.650	6.308
	Euclidean	0.166	1.489	4.438 —	0.303	0.977	8.020 —	0.279	1.004	7.352 —	0.214	1.693	6.004
ResNet50/CIFAR-100	Oracle	0.017	4.285	1.496 —	0.011	1.072	1.625 —	0.063	2.031	1.818 —	0.028	3.034	1.520
	IOS	0.013	4.995	1.223 —	0.009	1.264	1.308 —	0.053	2.399	1.474 —	0.023	3.491	1.238
	Cosine	0.013	5.186	1.187 —	0.009	1.295	1.256 —	0.051	2.450	1.433 —	0.022	3.599	1.206
	Euclidean	0.013	5.328	1.143 —	0.009	1.325	1.229 —	0.049	2.538	1.363 —	0.021	3.683	1.166
BERT-base/20News	Oracle	0.318	0.962	7.485 —	0.472	0.735	13.102 —	0.345	0.874	8.336 —	0.304	1.184	7.923
	IOS	0.268	1.118	6.214 —	0.416	0.884	10.980 —	0.292	0.944	6.984 —	0.256	1.346	6.582
	Cosine	0.259	1.145	5.980 —	0.404	0.909	10.612 —	0.281	0.970	6.795 —	0.246	1.378	6.331
	Euclidean	0.253	1.168	5.920 —	0.392	0.936	10.400 —	0.277	0.980	6.693 —	0.241	1.406	6.268

client) and distribute examples uniformly within the assigned subset. Unless noted, we use N=20 clients with balanced sample counts, and hold out 10% client-local validation for diagnostics.

Implementation details. All training and importance accumulation are local to clients. Importance vectors aggregate over *all* local batches unless a cap is stated; for stability, we apply an EMA with  $\beta=0.8$ . Models are trained with SGD (momentum 0.9, weight decay  $5\times 10^{-4}$ ) and fixed LR of 0.001; FMNIST runs for 10 local epochs, CIFAR-10/100, and BERT-base for 50; results are averaged over three seeds. Experiments run on a single A100 GPU (40 GB) with a 32-core CPU.

### 6.2 EVALUATING *IOS* APPLICATIONS

Clustered Personalized FL. We follow IFCA/Clustered-FL style clustering and then train cluster-specialized models; all methods share identical data splits, models, and training budgets. Table 1 reports CPFL test accuracy (%) across four heterogeneity regimes on different models. We compare *IOS* to two Cosine-based affinities commonly used in clustered FL pipelines: Cosine and Euclidean indicate Cosine and Euclidean similarities between client gradient/importance vectors aggregated locally. Both transmit full real-valued vectors; *IOS* transmits only top-K index sets. The results indicate that *IOS* consistently outperforms value-based similarities across vision and text, with gains larger under stronger heterogeneity. For example, on CIFAR-100 at Dir(0.1), *IOS* attains 48.3% vs. 45.40% (Cosine) and 43.18% (Euclidean). *IOS* is best in most settings, +1.47 pp over Cosine, with dataset-wise gains of +2.04 (FMNIST), +1.54 (CIFAR-10), +1.20 (CIFAR-100), and +1.12 pp (20News); Euclidean is consistently lower than Cosine and IOS.

Cluster Quality Evaluation. We assess clustering with Silhouette, Davies–Bouldin (DB), and Calinski–Harabasz (CH), where higher Sil/CH and lower DB are better. Using ground-truth label affinities, the *Oracle* is an upper bound. Table 2 shows that *IOS* reliably closes most of the gap to Oracle: per cell it is only  $\sim$ 15–20% lower on Sil/CH and  $\sim$ 14–20% higher on DB. *Cosine* (on

Table 3: **Recall@k vs. Oracle** for neighbor selection. Within each regime, we report Recall@k.

Dataset / Model	Method	Dir(0.3) (k=4/8/16)	Dir(0.1) (k=4 / 8 / 16)	Patho(20%) (k=4 / 8 / 16)	Patho(30%) (k=4 / 8 / 16)
CNN / FMNIST	Cosine Euclidean IOS	0.56 / 0.74 / 0.90	0.50 / 0.68 / 0.84	0.55 / 0.72 / 0.88 0.48 / 0.65 / 0.81 <b>0.58 / 0.74 / 0.90</b>	0.49 / 0.66 / 0.82
ResNet18 / CIFAR-10	Cosine Euclidean IOS	0.53 / 0.72 / 0.89	0.44 / 0.60 / 0.77	0.50 / 0.66 / 0.81 0.42 / 0.58 / 0.75 <b>0.52 / 0.74 / 0.86</b>	0.43 / 0.59 / 0.76
ResNet50 / CIFAR-100	Cosine Euclidean IOS	0.40 / 0.58 / 0.75	0.36 / 0.53 / 0.70	0.42 / 0.59 / 0.75 0.34 / 0.51 / 0.68 <b>0.44 / 0.63 / 0.79</b>	0.35 / 0.52 / 0.69
BERT-base / 20News	Cosine Euclidean IOS	0.61 / 0.79 / 0.95	0.57 / 0.75 / 0.92	0.62 / 0.80 / 0.96 0.55 / 0.73 / 0.90 <b>0.63 / 0.82 / 0.95</b>	0.56 / 0.74 / 0.91

gradients/parameters) is a further  $\sim$ 1–4% worse than IOS, and *Euclidean* trails Cosine by another  $\sim$ 2–4%. Overall, *IOS* captures client relatedness more consistently than value-based similarities, yielding tighter intra-cluster cohesion and clearer inter-cluster separation.

**Neighbor Selection: Retrieval Quality vs. Oracle.** Personalized FL hinges on selecting the right peers: if the neighborhood assembled for a client does not mirror its underlying data distribution, no aggregation rule can reliably personalize downstream models. Recall@ $k = \frac{1}{n} \sum_{i=1}^{n} \frac{|N_k^S(i) \cap N_k^*(i)|}{k}$  directly measures how many oracle neighbors are recovered, while its trend over k and across heterogeneity regimes reveals robustness. For each client i, let  $N_k^*(i)$  be the oracle neighbor set: the k clients with the smallest Wasserstein distance between the true per-client label histograms (unavailable in practice; used only for evaluation). A method S (Cosine, Euclidean, or IOS) returns  $N_k^S(i)$  via k-NN on its similarity. In this experiment N = 40, and We sweep  $k \in \{4, 8, 16\}$ .

Table 3 shows that *IOS* retrieves oracle neighbors more reliably than value-based baselines, with larger gaps at higher k and stronger non-IID. In the hardest regime (CIFAR-100, Dir(0.1)), *IOS* reaches R@8 = 0.67 vs. 0.61/0.53 (Cosine/Euclidean) and 0.81 at k=16 vs. 0.77/0.70; on CIFAR-10, Dir(0.1) it attains R@8 = 0.73 (+0.05/+0.13) and 0.90 at k=16 vs. 0.83/0.77. Recall rises with k for all methods, yet *IOS* keeps a lead (FMNIST, Dir(0.3), k=16: 0.97 vs. 0.96/0.90). As heterogeneity strengthens (Dir  $0.3 \rightarrow 0.1$ ), all recalls drop but *IOS* degrades less (CIFAR-100 R@8:  $0.70 \rightarrow 0.67$  vs.  $0.66 \rightarrow 0.61$ ). The pattern is modality-agnostic: *IOS* tracks oracle structure in text (20News R@16 up to 0.97), while vision under severe non-IID remains harder yet still favors *IOS*.

Shapley-Style Donor Ranking. In many PFL/DFL schemes, a client aggregates only from a few high-value donors. Inspired by Shapley-value notions of contributor utility in ML (Ghorbani & Zou, 2019; Wang et al., 2020; Lin et al., 2022), we *implement* a Shapley-style donor-ranking evaluation in our codebase: the *oracle* similarity  $S^*(i,j)$  via wasserstein distance over true label histograms induces the ground-truth donor order  $\pi_*(i)$ . For any operational similarity S, we get a method order  $\pi_S(i)$  by sorting row  $S[i,\cdot]$  (excluding i). In our implementation, we replace Cosine with IOS: the default donor ranking uses IOS (index-overlap on Top-K supports) rather than Cosine, and we report agreement with  $\pi_*(i)$  via Kendall's  $\tau$  and Recall@k (averaged over clients and final rounds).

Table 4 shows IOS yields the highest agreement with the oracle, outperforming both Cosine and Euclidean. Gains are consistent across regimes: on CIFAR-10 with Dir(0.1), IOS reaches  $\tau$ =0.48 / R@5 = 0.62 vs. 0.36/0.55 (Cosine) and 0.35/0.53 (Euclidean). On the harder CIFAR-100, Dir(0.1), IOS attains 0.36/0.52 vs. 0.28/0.49 and 0.28/0.48. For FMNIST, Dir(0.3), IOS 0.59/0.68 exceeds 0.48/0.63 and 0.47/0.62; for 20News, Dir(0.3), IOS 0.64/0.72 improves over 0.55/0.68 and 0.54/0.67. Euclidean trails Cosine by  $\approx$ 0.01–0.02 in both  $\tau$  and R@5 throughout, reinforcing that  $indices-only\ IOS$  better preserves oracle donor priority without real-valued sharing.

**Oracle Distribution Alignment.** Personalization quality depends on whether a client's neighbor mixture reproduces its *true* data distribution. For client i with label histogram  $p_i$ , and a method S that selects a k-NN set  $N_k^S(i)$  with weights  $w_{ij} \propto S(i,j)$  (row-normalized), we form the induced mixture

Table 4: Shapley-style donor ranking vs. oracle. Each cell reports Kendall's \u03c4 / Recall@5. Methods: Cosine, Euclidean, and IOS. Oracle ranking is induced by 1-JS on true label distributions.

Dataset / Model	Method	Dir(0.3) (τ / R@5)	Dir(0.1) (τ / R@5)	Patho(20%) (τ / R@5)	Patho(30%) (τ / R@5)
CNN / FMNIST	Cosine	0.48 / 0.63	0.41 / 0.57	0.39 / 0.55	0.40 / 0.56
	Euclidean	0.47 / 0.62	0.43 / 0.56	0.38 / 0.54	0.40 / 0.55
	IOS	<b>0.59</b> / <b>0.68</b>	<b>0.52 / 0.62</b>	<b>0.49 / 0.60</b>	<b>0.51 / 0.61</b>
ResNet18 / CIFAR-10	Cosine	0.44 / 0.60	0.36 / 0.55	0.34 / 0.54	0.35 / 0.54
	Euclidean	0.43 / 0.59	0.35 / 0.53	0.34 / 0.53	0.34 / 0.53
	IOS	<b>0.56 / 0.66</b>	<b>0.48 / 0.62</b>	<b>0.45</b> / <b>0.60</b>	<b>0.46 / 0.61</b>
ResNet50 / CIFAR-100	Cosine	0.32 / 0.52	0.28 / 0.49	0.27 / 0.48	0.28 / 0.49
	Euclidean	0.31 / 0.51	0.28 / 0.48	0.26 / 0.47	0.28 / 0.48
	IOS	<b>0.41 / 0.56</b>	<b>0.36 / 0.52</b>	<b>0.34 / 0.50</b>	<b>0.35 / 0.51</b>
BERT-base / 20News	Cosine	0.55 / 0.68	0.50 / 0.64	0.48 / 0.63	0.49 / 0.64
	Euclidean	0.54 / 0.67	0.49 / 0.62	0.47 / 0.62	0.48 / 0.62
	IOS	<b>0.64 / 0.72</b>	<b>0.59</b> / <b>0.69</b>	<b>0.57</b> / <b>0.67</b>	<b>0.58</b> / <b>0.68</b>

Table 5: Oracle distribution alignment (JS divergence; lower is better) with k=8 neighbors. Bold

indicates	the	best	non-c	orac.	le	method.

		Dir	(0.3)			Dir	(0.1)			Patho	(20%)		Patho(30%)				
Dataset / Model	Oracle	Cosine	Euclidean	SOI	Oracle	Cosine	Euclidean	SOI	Oracle	Cosine	Euclidean	SOI	Oracle	Cosine	Euclidean	SOI	
CNN / FMNIST			0.233				0.265		0.217		0.279			0.242			
ResNet18 / CIFAR-10			0.226					0.219			0.275			0.236			
ResNet50 / CIFAR-100		0.278		0.254			0.329				0.352			0.305			
20News / BERT-base	0.123	0.141	0.156	0.134	0.136	0.153	0.100	0.145	0.145	0.163	0.174	0.155	0.132	0.156	0.168	0.148	

  $\hat{p}_i^S(k) = \sum_{j \in N_b^S(i)} w_{ij} \, p_j$ . We evaluate alignment via Jensen–Shannon divergence  $\mathrm{JS}ig(p_i, \hat{p}_i^S(k)ig)$ (lower is better), averaged over clients and the last 10 rounds. As an *oracle* upper bound, we compute  $N_k^*(i)$  using the k neighbors that minimize the Wasserstein distance to  $p_i$  (unavailable in practice; used only for evaluation). The result in Table 5 for fix k=8 indicates that (i) IOS consistently yields tighter alignment to client distributions. Across all settings, IOS is the best non-oracle: e.g., on CIFAR-10 and Dir(0.1), IOS reduces JS to 0.219 vs. 0.238 (Cosine,  $\downarrow 0.019$ ) and 0.252 (Euclidean,  $\downarrow 0.033$ ), approaching the oracle's 0.195 within 0.024. On the harder CIFAR-100, Dir(0.1), IOS attains 0.285, improving over Cosine by 0.027 and over Euclidean by 0.044, and within 0.020 of oracle. (ii) Non-IID severity increases mismatch for all, but IOS degrades least. Moving from Dir(0.3) to Dir(0.1) increases JS by  $\sim 0.02$ –0.04; IOS's increments are systematically smaller than Cosine/Euclidean (e.g., FMNIST: +0.029 for IOS vs. +0.033 / +0.032). (iii) Modality trend holds. Text classification (BERT-base/20News) exhibits lower JS overall, yet IOS preserves a clear gap (Dir(0.3): 0.134 vs. 0.141 / 0.156), showing that indices-only geometry transfers beyond vision.

#### DISCUSSION & CONCLUSION.

IOS is a value-free, index-only similarity —a drop-in replacement for value-based (parameter/gradient) affinities when sharing real numbers is undesirable. While our evaluation focuses on FL, the abstraction extends beyond federation: IOS applies to any domain that needs modelto-model similarity (e.g., model hubs, ensemble selection, continual/transfer learning, checkpoint curation) where value sharing is costly or risky. By measuring support overlap of salient parameters, it captures the inductive bias most predictive of relatedness, yielding consistent gains across applications: clustered PFL (faster convergence, higher accuracy), neighbor selection (higher Recall@k), donor ranking (higher Kendall- $\tau$ ), and distribution alignment (lower JS). Exchanging indices only shrinks the attack surface (gradient leakage, model inversion, membership/property inference). Overall, IOS recovers near-oracle structure with lightweight communication and a strong privacy posture—advancing scalable personalization in decentralized learning.

## REFERENCES

- Shun-ichi Amari. Natural gradient works efficiently in learning. *Neural Computation*, 10(2):251–276, 1998.
- Alexandr Andoni and Piotr Indyk. Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions. *Communications of the ACM*, 51(1):117–122, 2008.
  - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.
  - Canh T. Dinh, Nguyen H. Tran, and Tuan D. Nguyen. Personalized federated learning with moreau envelopes. In *NeurIPS*, 2020.
    - Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021.
  - Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the lottery: Making all tickets winners. In *NeurIPS*, 2020.
    - Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning: A model-agnostic meta-learning approach. In *NeurIPS*, 2020.
  - Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. In *CCS*, 2015.
  - Brendan J. Frey and Delbert Dueck. Clustering by passing messages between data points. *Science*, 315(5814):972–976, February 2007. doi: 10.1126/science.1136800.
  - Trevor Gale, Erich Elsen, and Sara Hooker. The state of sparsity in deep neural networks. *arXiv* preprint arXiv:1902.09574, 2019.
  - João Gama, Indrè Žliobaitė, Albert Bifet, Mykola Pechenizkiy, Indre Žliobaite, et al. A survey on concept drift adaptation. *ACM Computing Surveys*, 46(4):44:1–44:37, 2014.
  - Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. Inverting gradients how easy is it to break privacy in federated learning? In *NeurIPS*, 2020.
  - Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine learning. In *ICML*, 2019.
  - Avishek Ghosh, Jichan Chung, Dong Yin, and Kannan Ramchandran. An efficient framework for clustered federated learning. In *NeurIPS*, 2020.
  - Song Han, Jeff Pool, John Tran, and William J. Dally. Learning both weights and connections for efficient neural networks. In *NeurIPS*, 2015.
  - Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
  - Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. In *International Conference on Learning Representations (ICLR)*, 2023. URL https://openreview.net/forum?id=6t0Kwf8-jrj.
  - Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In *STOC*, 1998.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
  - Ruoxi Jia, Daphne Dao, Boxin Wang, Frank Hubis, et al. Efficient task-specific data valuation for nearest neighbor algorithms. *PVLDB*, 12(11):1610–1623, 2019.

- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J. Reddi, Sebastian U. Stich, and Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for federated learning. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119 of *Proceedings of Machine Learning Research*, pp. 5132–5143. PMLR, 2020.
- Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *ICML*, 2017.
  - Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *ICML*, 2019.
  - Frederik Kunstner, Lukas Balles, and Philipp Hennig. Limitations of the empirical fisher approximation for natural gradient descent. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
  - Namhoon Lee, Thalaiyasingam Ajanthan, and Philip H. S. Torr. Snip: Single-shot network pruning based on connection sensitivity. In *ICLR*, 2019.
  - Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. In *International Conference on Learning Representations (ICLR)*, 2017.
  - X Lin et al. Measuring and learning data quality for fl via shapley values. In AAAI, 2022.
  - Yujun Lin, Song Han, Huizi Mao, Yu Wang, and William J. Dally. Deep gradient compression: Reducing the communication bandwidth for distributed training. In *ICLR*, 2018.
  - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In ICLR, 2019.
  - Jie Lu, Anjin Liu, Fan Dong, Feng Gu, João Gama, and Guangquan Zhang. Learning under concept drift: A review. *IEEE Transactions on Knowledge and Data Engineering*, 31(12):2346–2363, 2018.
  - Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov. Exploiting unintended feature leakage in collaborative learning. In *IEEE Symposium on Security and Privacy*, 2019.
  - Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
  - Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, et al. Mixed precision training. In *ICLR*, 2018.
  - Ari S. Morcos, Haonan Yu, Michela Paganini, and Yuandong Tian. One ticket to win them all: Generalizing lottery ticket initializations across datasets and optimizers. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
  - Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability. *arXiv preprint arXiv:1706.05806*, 2017.
  - Felix Sattler, Klaus-Robert Müller, Wojciech Samek, and Thomas Wiegand. Clustered federated learning: Model-agnostic distributed multi-task optimization under privacy constraints. In *IEEE Transactions on Neural Networks and Learning Systems*, 2020. doi: 10.1109/TNNLS.2020. 3015958.
  - Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *IEEE Symposium on Security and Privacy*, 2017.
  - Alexander Soen and Ke Sun. Trade-offs of diagonal fisher information matrix estimators. In *Advances in Neural Information Processing Systems*, volume 37, 2024.
  - Yue Tan, Guodong Long, Lu Liu, Tianyi Zhou, Qinghua Lu, Jing Jiang, and Chengqi Zhang. Fed-proto: Federated prototype learning across heterogeneous clients. In *AAAI*, 2022.

J Wang et al. Measuring the influence of clients in federated learning. In NeurIPS Workshop on FL, 2020. Lin Wang, Yongxin Guo, Tao Lin, and Xiaoying Tang. DELTA: Diverse client sampling for fast federated learning. In Advances in Neural Information Processing Systems, volume 36, 2023. Xiaodong Wu, Wenyi Yu, Chao Zhang, and Philip Woodland. An improved empirical fisher ap-proximation for natural gradient descent. In Advances in Neural Information Processing Systems, volume 37, 2024. Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated learning with non-iid data. arXiv preprint arXiv:1806.00582, 2018. Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients. In *NeurIPS*, 2019. 

Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head

the 57th Annual Meeting of the Association for Computational Linguistics (ACL), 2019.

self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *Proceedings of* 

# A HASH-LSH FOR SCALABLE CANDIDATE RETRIEVAL

Why Hash-LSH. Computing all-pairs overlap between n clients' supports  $\{I_i\}_{i=1}^n \operatorname{costs} O(n^2K)$  even with sorted lists/bitsets. We use MinHash with locality-sensitive hashing (LSH) to retrieve candidate neighbors in subquadratic time and bytes: each client publishes a short signature (or only band buckets), the server (or a peer) looks up candidates that collide in at least one band, and exact similarity is computed only on this small candidate set. This preserves IOS's value-free property (only indices/signatures, no real-valued parameters) and scales to large n. For a support set  $I \subseteq [M]$ , define h independent minhashes  $\phi_k(I) = \min\{H_k(j) : j \in I\}$  using 2-universal (approx. minwise) hash functions  $H_k : [M] \to \{0, \dots, 2^{64} - 1\}$ . The signature is  $\Phi(I) = (\phi_1(I), \dots, \phi_h(I)) \in \mathbb{N}^h$ . For any two sets  $I_i, I_j$  with S(i, j), an an unbiased estimator with  $\operatorname{Var}[\hat{s}] = s(1-s)/h$  is  $\mathbb{P}\big[\phi_k(I_i) = \phi_k(I_j)\big] = s \Rightarrow \hat{s} = \frac{1}{h} \sum_{k=1}^h \mathbf{1}\{\phi_k(I_i) = \phi_k(I_j)\}$ . Signature cost is O(hK) time and O(h) words per client.

# Algorithm 2 IOS–MinHash–LSH (build & query)

**Input:** Supports  $\{I_i\}_{i=1}^n$ , hash family  $\{H_k\}_{k=1}^h$ , bands b, rows per band r, exact index store for  $I_i$  1: **Build:** For each client i:

- 2: Compute signature  $\Phi(I_i)$  where  $\phi_k(I_i) = \min_{j \in I_i} H_k(j)$
- 3: For each band u=1..b, form key  $B_u(I_i)=(\phi_{(u-1)r+1},\ldots,\phi_{ur})$  and insert i into table  $\mathcal{T}_u[B_u(I_i)]$
- 4: **Query**(q): Given a query support  $I_q$
- 5: Compute  $\Phi(I_q)$  and band keys  $\hat{B}_u(I_q)$
- 6: Candidates  $C \leftarrow \bigcup_{u=1}^{b} \mathcal{T}_{u}[B_{u}(I_{q})] \setminus \{q\}$
- 7: For each  $j \in C$ : fetch  $I_j$  (if not local), compute  $S(q,j) = \frac{|I_q \cap I_j|}{K}$
- 8: Return top-k by S(q, j) (or all with  $S(q, j) \ge s_{\min}$ )

## B ALIGNMENT PROOFS (OVERLAP-ONLY)

# B.1 AUXILIARY BOUND

**Lemma 1** (Per-coordinate lower bound from dispersion). If  $\kappa_i \leq \kappa$ , then for every  $u \in I_i$ ,

$$(h_i)_u \geq \frac{\|h_i\|_2}{\kappa\sqrt{K}}.$$

*Proof.* Let  $a = \min_{u \in I_i} (h_i)_u$  and  $b = \max_{u \in I_i} (h_i)_u \le \kappa a$ . Then  $||h_i||_2^2 = \sum_{u \in I_i} (h_i)_u^2 \le Kb^2 \le K\kappa^2 a^2$ , so  $a \ge ||h_i||_2 / (\kappa \sqrt{K})$ .

### B.2 Proof of Proposition 1

Because  $s_i, s_j$  are nonnegative,

$$\langle s_i, s_j \rangle \geq \sum_{u \in I_i \cap I_j} (h_i)_u (h_j)_u \geq s_{ij} \cdot \frac{\|h_i\|_2}{\kappa \sqrt{K}} \cdot \frac{\|h_j\|_2}{\kappa \sqrt{K}} = \frac{s_{ij}}{\kappa^2 K} \|h_i\|_2 \|h_j\|_2,$$

using Lemma 1. Head dominance gives  $||h_i||_2 \ge \sqrt{1-\varepsilon} ||s_i||_2$  and similarly for j. Divide by  $||s_i||_2 ||s_j||_2$  to obtain equation 2.

# C STABILITY PROOFS (EMA + UNION BOUND)

#### C.1 EMA VARIANCE

Let  $\eta^{(t)} = \sum_{s\geq 0} (1-\beta)\beta^s \xi^{(t-s)}$  be the EMA noise at any coordinate (client index suppressed). Since the  $\xi^{(t)}$  are independent, mean-zero, sub-Gaussian with proxy variance  $\sigma^2$ ,

$$Var(\eta^{(t)}) = \sigma^2 \sum_{s \ge 0} (1 - \beta)^2 \beta^{2s} = \sigma^2 \frac{(1 - \beta)^2}{1 - \beta^2} = \sigma^2 \frac{1 - \beta}{1 + \beta}.$$

## C.2 CONCENTRATION AT THE BOUNDARY

Let H be the set of the K largest means of  $\mu$  and  $\bar{H} = [M] \backslash H$ . For  $u \in H$ ,

$$\Pr\left(\tilde{g}_u^{(t)} < \mu_{(K)} - \frac{\Delta}{2}\right) \le \exp\left(-\frac{\Delta^2}{8\,\sigma(\beta)^2}\right).$$

A union bound over the K elements of H yields  $\Pr(\min_{u \in H} \tilde{g}_u^{(t)} < \mu_{(K)} - \Delta/2) \leq K \exp(-\Delta^2/(8\sigma(\beta)^2))$ . Similarly, for  $v \in \bar{H}$ ,  $\Pr(\tilde{g}_v^{(t)} > \mu_{(K+1)} + \Delta/2) \leq \exp(-\Delta^2/(8\sigma(\beta)^2))$ , and union over M - K indices gives  $\Pr(\max_{v \in \bar{H}} \tilde{g}_v^{(t)} > \mu_{(K+1)} + \Delta/2) \leq (M - K) \exp(-\Delta^2/(8\sigma(\beta)^2))$ . Union over these two bad events proves Theorem 1.

#### C.3 EXPECTED TEMPORAL SELF-OVERLAP

Let  $\mathcal B$  be the subset of K-boundary positions whose mean gaps exceed  $c\,\sigma(\beta)\sqrt{\log M}$ ; assume  $|\mathcal B|/K\geq \rho$ . On the event guaranteed by Theorem 1 at times t-1 and t, the indices in  $\mathcal B$  persist in the Top-K. Taking expectation over the complement yields  $\mathbb E[\Gamma^{(t)}]\geq \rho-O(M^{1-c^2/8})$ .

# D MORE THEORETICAL ANALYSIS

## D.1 ROBUSTNESS TO RE-SCALING AND DIAGONAL PRECONDITIONING

**Setting (global Top-**K**).** We select a single global set  $(g, K) \subseteq [M]$  over the concatenated parameter vector; *no per-layer budgets* are used. The results below distinguish (i) invariances under positive scalings, (ii) effects of blockwise (layer-constant) rescaling, and (iii) diagonal preconditioners.

**Lemma C.1 (Positive scalar invariance; blockwise order preservation).** Let importance be monotone in |g| or  $g^2$ . Then, for any a > 0,

$$(ag, K) = (g, K).$$

Moreover, if D rescales each layer  $\ell$  by a positive constant  $d_{\ell} > 0$ , the *within-layer* ranking of coordinates is unchanged, although the *global* Top-K membership may change via cross-layer swaps. *Proof.* Positive scalar scaling preserves all pairwise orders; blockwise scaling preserves orders within blocks (layers).

Lemma C.2 (Cross-layer stability under bounded block rescaling). Let  $D = \operatorname{diag}(d_u)$  with  $d_u = d_\ell$  for all u in layer  $\ell$ , and block scales  $d_\ell \in [1/\chi, \chi]$ . Let  $I^\star = (g, K)$  and suppose g has dispersion margin  $\kappa > 1$  at the boundary: for every pair (u, v) with  $u \in I^\star$ ,  $v \notin I^\star$ , we have  $g_u/g_v \ge \kappa$ . Then

$$\frac{\left|\,(Dg,K)\,\triangle\,I^\star\,\right|}{K}\;\leq\;\eta(\kappa,\chi,K),$$

where  $\eta(\kappa,\chi,K)$  counts boundary pairs whose ratio lies in  $[1,\chi^2)$  (i.e., near-ties that block rescaling can invert). In particular,  $\eta(\kappa,\chi,K)\downarrow 0$  as  $\chi\downarrow 1$  or as  $\kappa$  increases. *Proof.* For any boundary pair (u,v), the post-rescaling ratio is  $(d_{\ell(u)}/d_{\ell(v)})\cdot (g_u/g_v)\in [\kappa/\chi^2,\kappa\chi^2]$ . If  $\kappa>\chi^2$ , the order is preserved. Violations can only arise from near-ties, which bounds the symmetric difference.  $\square$ 

Corollary C.3 (Diagonal preconditioners). Let  $D_t = \operatorname{diag}(d_{t,u})$  be a diagonal preconditioner (e.g., Adam's  $\hat{v}_t^{-1/2}$ ) with overall condition number  $\chi_t = \frac{\max_u d_{t,u}}{\min_u d_{t,u}} \leq \bar{\chi}$ . Then, across a round,

$$\frac{\left| (D_t g, K) \triangle (g, K) \right|}{K} \leq \eta(\kappa, \bar{\chi}, K).$$

*Proof.* For any boundary pair, the scaled ratio lies in  $[\kappa/\bar{\chi}^2,\kappa\bar{\chi}^2]$ ; the same near-tie argument as in Lemma C.2 applies.

### D.2 CHOOSING K: COVERAGE, STABILITY, AND SAMPLE COMPLEXITY

**Proposition D.1 (Monotone coverage; selection complexity).** Let  $h_i$  denote the head of length K after ordering coordinates of  $g_i$  by decreasing importance. Then  $C_i(K) = \|h_i\|_2^2/\|g_i\|_2^2$  is non-decreasing in K. The smallest  $K \leq K_{\max}$  with  $C_i(K) \geq \tau$  can be found by binary search in  $O(\log K_{\max})$  iterations, each using a selection step that is O(M) expected time (via Quickselect) or  $O(M\log M)$  with sorting. *Proof.* Adding indices cannot reduce head energy; complexity follows from selection/sorting costs.

**Proposition D.2 (Stability estimator concentration).** Let  $\hat{\rho}_i(K) = \frac{1}{r} \sum_{s=1}^r Z_s$  where  $Z_s \in \{0,1\}$  indicates whether the Top-K set is unchanged between two adjacent EMA slices (or light bootstrap resamples). Then, for any  $\delta \in (0,1)$ ,

$$\Pr(|\hat{\rho}_i(K) - \rho_i(K)| > \delta) \le 2\exp(-2r\delta^2).$$

Thus  $r \geq \frac{1}{2\delta^2}\log\frac{2}{\gamma}$  samples suffice for accuracy  $\delta$  with confidence  $1-\gamma$ . *Proof.* Hoeffding's inequality.

**Remark D.3 (End-to-end frontier).** Choosing  $(\tau, \rho_0, r, \beta)$  traces a utility–privacy–bandwidth curve. Larger  $\beta$  improves  $\hat{\rho}$  but slows drift response; r trades runtime for confidence.

#### D.3 MINHASH-LSH FACTS AND END-TO-END COMPLEXITY

**Lemma E.1 (MinHash unbiasedness and variance).** Let  $S_1, S_2 \subseteq [M]$  and  $s = J(S_1, S_2) = \frac{|S_1 \cap S_2|}{|S_1 \cup S_2|}$ . With h independent MinHash functions, the estimator  $\hat{s} = \frac{1}{h} \sum_{u=1}^{h} \mathbf{1}\{ \mathrm{mh}_u(S_1) = \mathrm{mh}_u(S_2) \}$  satisfies

$$\mathbb{E}[\hat{s}] = s, \quad \operatorname{Var}[\hat{s}] = \frac{s(1-s)}{h}.$$

*Proof.* Each indicator is Bernoulli(s) by MinHash collision equivalence.

**Lemma E.2 (Banding retrieval probability).** Arrange h = br MinHash values into b bands of r rows; two sets with similarity s collide in a band with probability  $s^r$ , and are retrieved with probability  $1 - (1 - s^r)^b$ . *Proof.* Standard LSH banding analysis with independent bands.

**Theorem E.3 (Candidate set size and total cost).** Let n clients each submit a K-subset (the TopK indices). Using h = br MinHashes and b hash tables: expected candidate set size per query is

$$\mathbb{E}[C] = 1 + \sum_{q \neq i} \left( 1 - (1 - s_{iq}^r)^b \right),$$

where  $s_{iq}$  are pairwise similarities. The total server cost per round is

$$O(nh) + O\left(\sum_{i=1}^n \mathbb{E}[C_i]\right) + O\left(\sum_{(i,q) \in \text{cands}} |I_i \cap I_q|\right),$$

i.e., signature build + candidate lookups + final exact overlaps, with the last term bounded by  $O(\sum_i \mathbb{E}[C_i] \cdot K)$ . Signature memory is O(nh). Proof. Linearity of expectation and that exact overlap is  $O(\min(K_i, K_g)) = O(K)$ .

## D.4 Complexity and Communication Analysis ( $K \ll M$ )

The client-side overhead to accumulate importances adds one O(M) elementwise update per minibatch (squared-gradient or second moment) on top of backprop; across E local epochs and  $|D_i|/B$ batches this is  $O(\frac{E|D_i|}{B}M)$  arithmetic with one extra M-vector in memory. Extracting supports uses a partial selection in O(M) expected time (e.g., nth\_element) or  $O(M\log M)$  for a full sort; storing indices costs  $O(K_i)$  integers (or O(M/word) words for a bitset when many intersections are reused). Server-side exact similarity with sorted lists performs one merge-style intersect per pair:  $O(\sum_{i < j} \min(K_i, K_j))$  time (uniform  $K: O(n^2K)$ ) and  $O(\sum_i K_i)$  storage; bitset AND+popcount is  $O(n^2M/\text{word})$  time and O(nM/word) storage and is preferable only when K is not tiny or the same bitsets are intersected many times. With MinHash-LSH, each client builds a signature in  $O(hK_i)$  time and inserts b band keys; per query, computing the signature and probing b buckets is O(h+b), and exact re-scoring over the retrieved candidate set  $C_q$  costs  $O(\sum_{j \in C_q} \min(K_q, K_j))$  (uniform  $K: O(|C_q|K)$ ), with memory O(nh) for signatures and O(nb)bucket pointers. Per-round communication is  $K_i \lceil \log_2 M \rceil$  bits if indices are sent (or  $h \cdot w$  bits for signatures only, with on-demand index fetch for candidates); compared to value sharing  $(M \cdot q)$  bits for q-bit weights/gradients), IOS reduces bytes by  $\Omega(M/K)$  (or  $\Omega(M/h)$  for signature-only discovery). End to end, an exact all-pairs round is dominated by  $O(n^2K)$  server time and  $O(\sum_i K_i \lceil \log_2 M \rceil)$ uplink, whereas the LSH pipeline is O(nh) build plus  $O(\sum_q |C_q|K)$  rescoring with  $|C_q| \ll n$  when banding is tuned so the selectivity threshold matches the target neighborhood; amortizing stable supports across rounds (EMA + delta-encoding or incremental rehash) further lowers both build time and bytes.

## E EXTENDED EXPERIMENTAL RESULTS

# F COMPARATIVE TEST ACCURACY OF MODEL BASED ON FL ROUNDS

We reported the aggregate test accuracies for Clustered Personalized FL (CPFL) application in Table1 of the main text. Here we complement those results with the *training dynamics*: Figures 1–3

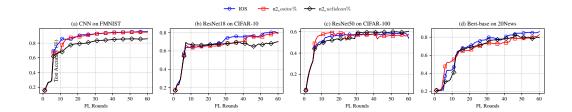


Figure 1: Comparative test accuracy of model based on FL rounds (Patho(20%)).

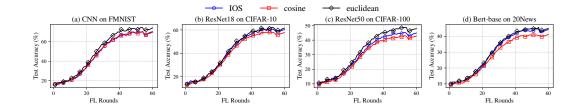


Figure 2: Comparative test accuracy of model based on FL rounds (Dir(0.1)).

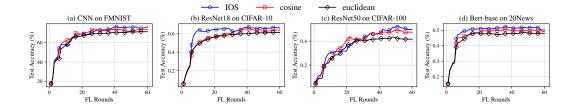


Figure 3: Comparative test accuracy of model based on FL rounds (Dir(0.3)).

plot test accuracy versus FL rounds under different heterogeneity regimes (Path(20%), Dir(0.1) and Dir(0.3)) for different model/dataset pairs: CNN on FMNIST, ResNet18 on CIFAR-10, ResNet50 on CIFAR-100, and BERT-base on 20News. The CPFL protocol is identical across methods—same partitions, architectures, optimizers, and schedules—and only the clustering similarity changes among IOS (Top-K index), cosine, and Euclidean. IOS exchanges indices only (no weights, gradients, logits, or activations).

Across all settings, *IOS* mirrors the learning trajectory of cosine and Euclidean, rising and saturating at comparable *round counts*, while achieving slightly higher final accuracy in most cases (consistent with Table X). The advantage is most noticeable on FMNIST and CIFAR-10, with parity or minor gaps on the harder CIFAR-100 and 20News tasks. The pattern persists when moving from Dirichlet to more pathological splits, indicating that overlap on salient indices provides a stable relatedness signal even when value-based metrics are sensitive to scale, normalization, or noise. For CPFL, *IOS* is a robust proxy for inter-client similarity: it reproduces the clustering behavior and convergence profile of cosine/Euclidean and typically yields modest accuracy gains, while *preserving privacy* by avoiding any transfer of real-valued parameters or gradients.

# F.1 Effect of the support size K for IOS

For each model/dataset we compare three *IOS* support sizes: the task-specific  $K^{\star}$  (selected by our coverage+stability rule; typically  $\approx$ 7–12% and concretely 12.2% for CNN/FMNIST, 7.3% for ResNet18/CIFAR-10, 9.6% for ResNet50/CIFAR-100, and 8.2% for BERT-base/20News), plus two

Table 6: **CPFL Accuracy** (%) with *IOS* at three K settings. Within each distribution block, columns are  $K^*$ , 15%, and 25%. Bold marks the best within each block.

	1	Dir(0.3	3)	1	Dir(0.1	.)	Pa	tho(20	0%)	Patho(30%)			
Model / Dataset	$K^{\star}$	15%	25%	$K^{\star}$	15%	25%	$K^{\star}$	15%	25%	$K^{\star}$	15%	25%	
CNN / FMNIST	76.23	73.23	70.03	74.79	71.29	67.39	95.18	92.98	90.38	83.15	80.45	77.65	
ResNet18 / CIFAR-10	66.13	63.53	60.13	62.10	58.70	54.90	81.44	78.54	74.84	74.81	72.41	69.61	
ResNet50 / CIFAR-100	49.71	45.81	40.71	48.30	44.20	38.50	58.07	54.87	49.77	52.98	49.48	44.08	
BERT-base / 20News	50.30	48.10	45.60	45.41	43.41	40.31	58.26	55.86	51.96	54.25	52.15	49.05	

Table 7: **Neighbor Selection Quality** (*Recall@8* vs. oracle) with *IOS* at three K settings. Higher is better. Bold marks the best within each block.

	I	Dir(0.3)	3)	L	Dir(0.1)	L)	Pa	tho(20	0%)	Pat	Patho(30%)			
Model / Dataset	$K^{\star}$	15%	25%	$K^{\star}$	$K^{\star}$ 15% 25%		$K^{\star}$	15%	25%	$K^{\star}$	15%	25%		
CNN / FMNIST ResNet18 / CIFAR-10 ResNet50 / CIFAR-100 BERT-base / 20News	0.80 0.70	0.85 0.77 0.66 0.80	0.73 0.61	0.73	0.74 0.69 0.63 0.81	0.70 0.65 0.57 0.77	0.74 0.74 0.63 0.82	0.60	0.67 0.65 0.55 0.77	0.73 0.64	0.75 0.70 0.61 0.81	0.66 0.56		

larger supports 15% and 25%. All other settings (non-IID regimes, data splits, training budgets, and pipelines) match §6.

Clustered Personalized FL (CPFL) accuracy vs. K. We run IFCA/Clustered-FL style training with the *same* clustering/training pipeline while varying only K in IOS to quantify how support size impacts downstream CPFL test accuracy.

Table 6 reports CPFL accuracy (%) across four heterogeneity regimes; within each block, columns are  $K^*$ , 15%, and 25%.

 $K^{\star}$  is best across all models/regimes. Moving to 15% costs roughly 2–6 pp depending on task/heterogeneity, while 25% introduces larger losses (typically 4–10 pp, up to  $\sim$ 10 pp on the hardest cases). Deeper models under stronger heterogeneity suffer more. On ResNet50/CIFAR-100 under Dir(0.1), accuracy falls from 48.30 at  $K^{\star}$  to 44.20 at 15% (-4.10 pp) and 38.50 at 25% (-9.80 pp). For BERT-base/20News under Patho(20%), the score decreases from 58.26 to 55.86 at 15% (-2.40 pp) and to 51.96 at 25% (-6.30 pp). On the lighter CNN/FMNIST regime Dir(0.3), accuracy drops from 76.23 to 73.23 at 15% (-3.00 pp) and to 70.03 at 25% (-6.20 pp). The results shows that Small, stable supports near the coverage/stability knee ( $K^{\star}$ ) give the best downstream accuracy; enlarging K dilutes salience, destabilizes Top-K boundaries, and consistently reduces CPFL performance, especially in harder regimes and deeper nets.

**Neighbor selection quality (Recall@8) vs.** K. We build the similarity graph with IOS at each K and perform k-NN (k=8) neighbor selection. We then compute Recall@8 against the oracle neighbors (defined by Wasserstein distance over true per-client label histograms; used only for evaluation).

Table 7 reports Recall@8 (higher is better) across the four heterogeneity regimes; within each block, columns are  $K^*$ , 15%, and 25%. The results indicate that  $K^*$  yields the highest Recall@8 throughout. Increasing to 15% reduces Recall by roughly  $0.02{\text -}0.05$ ; to 25% by about  $0.05{\text -}0.09$ . Losses are more pronounced in harder regimes and for deeper models (e.g., ResNet50 under Dir(0.1)). For  $ResNet50/CIFAR{\text -}100$  under Dir(0.1), Recall@8 declines from 0.67 at  $K^*$  to 0.63 at 15% (-0.04) and 0.57 at 25% (-0.10). On  $ResNet18/CIFAR{\text -}10$  with Patho(20%), it goes from 0.74 to 0.70 at 15% (-0.04) and 0.65 at 25% (-0.09). For  $BERT{\text -}base/20News$  under Dir(0.3), Recall@8 moves from 0.82 to 0.80 at 15% (-0.02) and 0.77 at 25% (-0.05). Neighbor retrieval quality tracks the same pattern as downstream accuracy: supports near  $K^*$  best preserve salient head overlap across clients, while larger supports introduce tail coordinates that dilute overlap and make Top-K less stable—lowering Recall and, downstream, CPFL performance.

Table 8: **CPFL Accuracy** (%) with *IOS* at  $K^*$  across local epochs. Within each distribution block, columns are 2/4/6/8/10 epochs. Bold marks the best within each block.

Dir(0.3)							Dir(0.1)					Pa	tho(20	0%)		Patho(30%)					
Model / Dataset	2	4	9	∞	10	2	4	9	∞	10	2	4	9	∞	10	2	4	9	∞	10	
CNN / FMNIST ResNet18 / CIFAR-10																72.00 63.20					
ResNet50 / CIFAR-100 BERT-base / 20News																40.90 44.50					

Table 9: **Neighbor Selection Quality** (Recall@8 vs. oracle) with IOS at  $K^*$  across epochs. Within each distribution block, columns are 2/4/6/8/10 epochs. Higher is better; bold marks the best.

		I	Dir(0.3)	3)		Dir(0.1)						Pa	tho(20	0%)		Patho(30%)				
Model / Dataset	2	4	9	∞	10	2	4	9	∞	10	2	4	9	∞	10	7	4	9	∞	10
CNN / FMNIST ResNet18 / CIFAR-10			0.84		0.89 0.81				0.78 0.73	0.79 0.74			0.70 0.70						0.78 0.73	
			0.65 0.79	0.70		0.52		0.61	0.67	0.68	0.48	0.53	0.58 0.78	0.63	0.64	0.49	0.54	0.59	0.64 0.83	0.65

# F.2 TRAINING EPOCHS ABLATION FOR IOS

**Protocol.** Unless stated otherwise, our study uses 8 local epochs. Here we ablate the number of local epochs  $\{2,4,6,8,10\}$  while *fixing IOS* at the task-specific  $K^*$  (coverage+stability rule) and keeping all other settings (non-IID regimes, data splits, budgets, optimizers, and pipelines) identical to §6. The goal is to test how many epochs are sufficient for clients to reflect their underlying distributions in both downstream CPFL accuracy and neighbor retrieval.

**Clustered Personalized FL (CPFL) accuracy vs. epochs.** We run IFCA/Clustered-FL with the same clustering/training pipeline and vary only the number of local epochs. We report CPFL test accuracy (%) across all four heterogeneity regimes.

Table 8 shows CPFL accuracy for IOS at  $K^*$  under epochs 2,4,6,8,10. The results show that accuracy increases steadily with more local training and saturates by 8 epochs: moving from  $6 \rightarrow 8$  yields small gains (typically +1–7 pp over 6), and 10 epochs bring only marginal improvements (often  $\leq 0.3$  pp) under the same budget. Eight local epochs are enough: they capture the client distribution well, while additional epochs yield diminishing returns; fewer epochs underfit and fail to expose sufficient salience structure for clustering to exploit. In contrast, under-training markedly hurts: 2 epochs are about 10–15 pp worse than 8 depending on model/regime, and 4 epochs lag by 6–10 pp. The effect is most pronounced in harder regimes (Dir(0.1), Patho) and for deeper networks (ResNet50), where more local steps are needed to shape client-specific heads and stabilize overlap. On ResNet50/CIFAR-100 with Dir(0.1), accuracy climbs from 35.10 (2 ep) to 39.90 (4 ep), 43.40 (6 ep), and 48.30 (8 ep), with only a negligible rise to 48.50 at 10 epochs; similarly, for BERT-base/20News under Patho(30%), scores go  $44.50 \rightarrow 48.90 \rightarrow 52.10 \rightarrow 54.25$  with a tiny lift to 54.34 at 10 epochs, while CNN/FMNIST under Patho(20%) jumps from 83.40 (2 ep) to 95.18 (8 ep) and only nudges to 95.30 at 10.

**Neighbor Selection: Retrieval Quality vs. Oracle.** We build the similarity graph with *IOS* at  $K^*$  for each epoch setting and perform k-NN (k=8) neighbor selection. We then compute Recall@8 against oracle neighbors (defined via Wasserstein distance over true per-client label histograms; used only for evaluation).

Table 9 reports Recall@8 across epochs 2, 4, 6, 8, 10. Neighbor retrieval improves smoothly with more local training and stabilizes by 8 epochs; moving to 10 brings at most a +0.01 gain in Recall@8. Under-training substantially lowers recall: at 2 epochs we see about 0.10-0.15 absolute drops relative to 8, at 4 epochs about 0.06-0.10, and at 6 epochs about 0.01-0.07. The pattern mirrors accuracy: deeper models and stronger heterogeneity require more local updates before the importance head is sharp enough to recover oracle-like neighborhoods. For *ResNet50/CIFAR-100* (Dir(0.1)), Recall@8 increases from 0.52 (2 ep) to 0.57 (4 ep), 0.61 (6 ep), and 0.67 (8 ep), with only a minimal change to 0.68 at 10; on *ResNet18/CIFAR-10* with Patho(20%), the trajectory  $0.59 \rightarrow 0.65 \rightarrow 0.70 \rightarrow 0.74 \rightarrow 0.75$  exhibits the same saturation; and *BERT-base/20News* under Dir(0.3) moves from 0.70 (2 ep) to 0.82 (8 ep) with a tiny step to 0.83 at 10.

USE OF LARGE LANGUAGE MODELS (LLMS) **Scope and intent.** LLMs were used *only* to aid and polish writing (grammar, clarity, concision, tone, and LaTeX hygiene). They were *not* used to design experiments, analyze data, generate results, choose hyperparameters, or create figures/tables. All technical contributions, algorithms, proofs, and empirical results originate from the authors. **Tools.** We used ChatGPT (GPT-5 Thinking) in an editorial capacity. Typical operations included: sentence rephrasing for clarity, reducing redundancy, harmonizing terminology/notation, improving caption phrasing, fixing cross-references, and standardizing style (e.g., capitalization, punctuation, hyphenation). When requested, it proposed concise alternatives that the authors reviewed and edited. Content provenance and verification. No passages were accepted verbatim without author re-view. The model was not allowed to introduce new claims, citations, equations, or numbers. All references and quantitative values in the paper were produced by our code/experiments and cross-checked by the authors. **Data and privacy.** We did *not* upload raw datasets, private code, or proprietary results. Shared text was limited to draft paragraphs, captions, and LaTeX snippets necessary for stylistic edits. No confidential or personal data were provided to the LLM. Bias and accountability. LLMs may reflect stylistic or cultural biases. Final wording, framing, and interpretations are the authors' responsibility. Any errors remain our own. **Reproducibility note.** The use of LLMs does not affect the reproducibility of results. All experi-ments can be reproduced from the released code, configurations, and seeds; LLM involvement was purely editorial. **Authorship.** All authors reviewed and approved the final text. The LLM is not listed as an author and did not meet authorship criteria.