

Introduction

In recent years, most Personalized Federated Learning (PFL) research has assumed that clients can afford tens to hundreds of communication rounds with the central server to adapt a global model to their local data distribution. In truly ad-hoc, one-shot settings, such as Low-Earth-Orbit (LEO) satellites and delivery drones, devices can only briefly connect to each other over highly constrained communication windows, rendering multi-round PFL infeasible.

Key Contributions

- We propose Federated Oriented Learning (FOL): a novel, four-stage one-shot personalization framework that enables clients to obtain fully personalized models in a single model exchange.
- We design an alignment-aware structured pruning mechanism: an approach that incorporates an alignment regularization term during pruning to retain only those filters and neurons in each neighbor's model that best match the client's own model and data distribution.
- We prove two theoretical guarantees: Upper bounds on the student-teacher risk discrepancy and convergence of the distillation process.

Related Work

One-Shot Federated Learning (OFL).

Methods such as DENSE and Co-Boosting can learn a single global model in one communication round, but that model remains generic rather than personalized. As a result, it often underperforms on individual clients' local datasets.

Dataset	Hurricane					
Satellite #	41	3	9	22	56	51
Methods	$\psi = 0.7$					
Local	90.45	82.35	88.63	90.67	86.01	91.18
FOL-A (E=1)	94.27	91.18	92.73	93.10	93.87	96.57
FOL-A (E=2)	95.54	94.12	93.64	93.68	95.16	97.06
FOL-A (E=3)	96.18	96.06	94.09	95.40	95.74	97.55
FOL (E=1)	93.11	85.29	90.02	91.95	89.81	93.63
FOL (E=2)	93.63	91.33	91.82	92.53	90.07	94.12
FOL (E=3)	94.27	93.04	92.27	94.25	91.92	95.59
DENSE	70.02	67.35	68.13	71.31	69.57	70.16
Co-Boosting	74.61	69.16	72.51	73.63	75.21	74.47

 Fig. 1 Performance of one-shot models on local data (Hurricane, $\psi=0.7$)

Personalized Federated Learning (PFL).

Methods like FedPer, FedRep and pFedMe yield client-specific models by decoupling or regularizing parameters, but they require tens to hundreds of server-client exchanges to converge, which is impractical in real-world settings with constrained communication, such as LEO satellites.

Key Challenges

The three key challenges to obtaining a fully personalized model under one-shot, server-free setting are:

Model Alignment under Heterogeneity.

Neighbor models may differ in architecture and are trained on non-IID data. How can a client, receiving these models a single communication round and without server coordination, adapt and prune each one to retain only the filters most relevant to its own architecture and data distribution?

Server-Free Ensemble Weighting.

How can a client compute optimal weights for the top-K adapted models to form a robust ensemble “teacher,” without any centralized coordination or additional communication?

Server-Free Compact Knowledge Distillation.

How can the ensemble's knowledge be efficiently distilled into a single, compact student model without any server-side orchestration or further communication?

Methodology

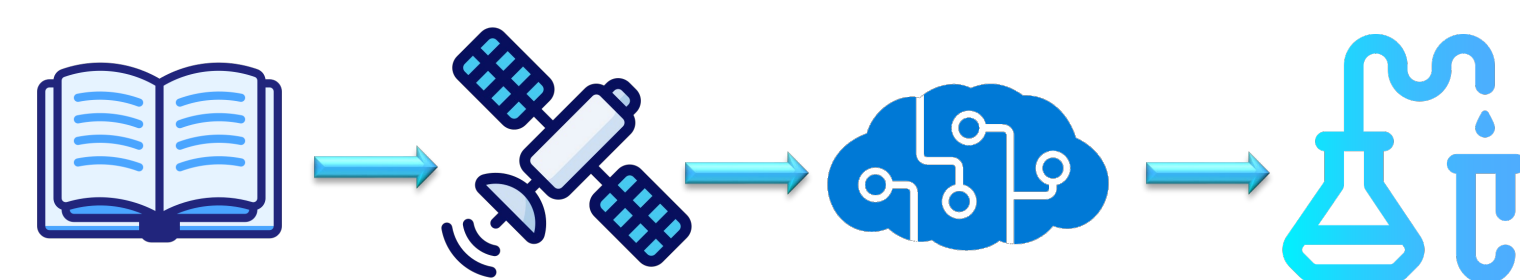


Fig. 2 Architecture Overview

1. Pretrain

- Each client trains an initial local model $\theta_k^{(1)-}$ on its own private dataset $\mathcal{D}_{\text{train}}^k$ using standard SGD:

$$\theta_k^{(1)-} \leftarrow \arg \min_{\theta_k^0} \frac{1}{|\mathcal{D}_{\text{train}}^k|} \sum_{(x_i, y_i) \in \mathcal{D}_{\text{train}}^k} \ell(f_k(x_i; \theta_k^0), y_i).$$

2. Collect and Align

- Receive up to Q neighboring models $\{\phi_j^{(e)}\}_{j=1}^Q$ in each model collection round.
- Fine-tune each received model on the client's local training set:

$$\phi_{j \rightarrow k}^{(e)} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_{\text{train}}^k|} \sum_{(x_i, y_i) \in \mathcal{D}_{\text{train}}^k} \ell(f_j(x_i; \phi), y_i),$$

where ϕ is initialized by $\phi \leftarrow \phi_j^{(e)}$.

- Apply alignment-aware structured pruning by solving following joint objective function:

$$\begin{aligned} \min_{\substack{\{\alpha_l\}_{(l, l') \in \mathbb{L}_{\text{shared}}(k, j)}, \\ \{\alpha_u\}_{u \in \mathbb{L}_{\text{unshared}}(k, j)}}} & \underbrace{\frac{1}{|\mathcal{D}_{\text{train}}^k|} \sum_{(x_i, y_i) \in \mathcal{D}_{\text{train}}^k} \ell(f_j(x_i; \tilde{\phi}_{j \rightarrow k}^{(e)}, \{\alpha_l\}, \{\alpha_u\}), y_i)}_{\text{1. Task Loss on Local Data}} \\ & + \underbrace{\lambda_p \sum_{(l, l') \in \mathbb{L}_{\text{shared}}(k, j)} \sum_{i=1}^{m_l} \left\| \alpha_{l, i} \mathbf{w}_{l, i}^{(j \rightarrow k)} - \mathbf{w}_{l', i}^k \right\|_2}_{\text{2. Alignment Regularization (Shared Layers Only)}} \\ & + \underbrace{\gamma_{\text{shared}} \sum_{l \in \mathbb{L}_{\text{shared}}(k, j)} \left\| \alpha_l \odot \mathbf{w}_l^{(j \rightarrow k)} \right\|_{2,1}}_{\text{3. Group-Lasso for Shared Layers}} \\ & + \underbrace{\gamma_{\text{unshared}} \sum_{u \in \mathbb{L}_{\text{unshared}}(k, j)} \left\| \alpha_u \odot \mathbf{w}_u^{(j \rightarrow k)} \right\|_{2,1}}_{\text{4. Group-Lasso for Unshared Layers}}, \end{aligned}$$

where λ_p and γ are hyperparameters controlling the strength of the alignment regularization and the structured pruning, respectively. $\|\cdot\|_{2,1}$ represents the group-lasso norm. \odot denotes element-wise multiplication.

- Apply a post fine-tuning on each pruned model to restore any lost accuracy:

$$\phi_{j \rightarrow k}^{(e)} \leftarrow \arg \min_{\phi} \frac{1}{|\mathcal{D}_{\text{train}}^k|} \sum_{(x_i, y_i) \in \mathcal{D}_{\text{train}}^k} \ell(f_j(x_i; \phi), y_i),$$

where ϕ is initialized by $\tilde{\phi}_{j \rightarrow k}^{(e)}$.

- Compute a validation score for each post-tuned neighbor model $\phi_{j \rightarrow k}^{(e)}$, and its own local model $\theta_k^{(e)-}$ on $\mathcal{D}_{\text{val}}^k$:

$$\text{score}_k^{(e)}(\theta) = \frac{1}{|\mathcal{D}_{\text{val}}^k|} \sum_{(x_i, y_i) \in \mathcal{D}_{\text{val}}^k} \mathbb{1}(\arg \max f(x_i; \theta) = y_i),$$

where $\mathbb{1}(\cdot)$ is the indicator function.

- Rank all candidates by their validation scores (breaking ties by cosine similarity) and choose the Top-K models for the ensuing ensemble stage.

3. Top-K Ensemble

- Form the optimal weighted ensemble “teacher”:

$$A_{\mathbf{w}_k^{(e)}}(x; \{s_i^{(e)}\}_{i=1}^K) = \sum_{i=1}^K w_i^{(e)} \cdot f_i(x; s_i^{(e)}),$$

where the optimal weights $\mathbf{w}_k^{(e)}$ is computed by minimizing the following loss:

$$\mathbf{w}_k^{(e)} = \arg \min_{\mathbf{w}_k^0} \frac{1}{|\mathcal{D}_{\text{train}}^k|} \sum_{(x_i, y_i) \in \mathcal{D}_{\text{train}}^k} \ell(A_{\mathbf{w}_k^0}(x_i; \{s_i^{(e)}\}_{i=1}^K), y_i).$$

4. Regularization-based Knowledge Distillation.

- Distill the weighted ensemble $A_w^{(e)}$ into the client's student model $\theta_k^{(e)+}$ by minimizing following KL-based distillation loss:

$$\begin{aligned} \mathcal{L}_{\text{KD}}(\theta_k^{(e)+}) = & \frac{1}{|\mathcal{D}_{\text{train}}^k|} \sum_{x_i \in \mathcal{D}_{\text{train}}^k} \text{KL}\left(\text{softmax}\left(\frac{A_{\mathbf{w}_k^{(e)}}(x_i)}{T}\right) \parallel \right. \\ & \left. \text{softmax}\left(\frac{f_k(x_i; \theta_k^{(e)+})}{T}\right)\right) + \lambda \|\theta_k^{(e)+} - \theta_k^{(e)-}\|^2, \end{aligned}$$

where $T > 0$ controls the smoothness of the softmax distributions applied to the logits.

Theoretical Analysis

Theorem 1. Risk Discrepancy Bound.

- Let $\theta_k^{(e)}$ be the student model obtained by minimizing the distillation loss $\mathcal{L}_{\text{KD}}(\theta_k^{(e)})$ on $\mathcal{D}_{\text{train}}^k$. Then, for a C-class problem with L-Lipschitz cross-entropy loss, $T > 0$, and softmax outputs in $(\alpha, 1-\alpha)$, the empirical risk discrepancy between the student and teacher models is bounded as follows:

$$|R_S(\theta_k^{(e)}) - R(A_{\mathbf{w}_k^{(e)}})| \leq \frac{L \cdot CT}{\alpha(1-\alpha)} \cdot \left(\frac{\mathcal{L}_{\text{KD}}(\theta_k^{(e)})}{2} + \frac{1}{8} \right).$$

Theorem 2. Convergence of Knowledge Distillation.

- Suppose $\{\theta_k^r\}_{r=0}^R$ are generated by $\theta_k^{r+1} = \theta_k^r - \eta \nabla \mathcal{L}_{\text{KD},k}(\theta_k^r, \xi_k^r)$, under standard assumptions that the distillation loss $\mathcal{L}_{\text{KD},k}$ is L_s -smooth and μ -strongly convex, and that the variance of the stochastic gradient is bounded by σ^2 , then for $r \geq 0$, and any step size $0 < \eta < 1/L_s$, the following bound holds:

$$\mathbb{E}[\|\theta_k^r - \theta_k^*\|^2] \leq \gamma^r \|\theta_k^0 - \theta_k^*\|^2 + \sum_{\tau=0}^{r-1} \gamma^\tau \beta,$$

Where $\gamma = (1 - 2\eta\mu + \frac{L_s^3}{\mu}\eta^2)$, $\beta = \eta^2\sigma^2$, and θ_k^* is the minimizer of $\mathcal{L}_{\text{KD},k}$.

Experimental Results

 Table 1. Test accuracies (%) on Wildfire and Hurricane ($\psi = 0.7$), reported as mean \pm std.

Dataset	Wildfire			Hurricane		
Satellite #	13	28	48	35	32	44
Methods	$\psi = 0.7$					
Local	94.23 \pm 1.84	94.12 \pm 1.80	90.53 \pm 1.57	86.93 \pm 1.56	87.34 \pm 1.60	89.82 \pm 1.82
FOL-A (E=1)	97.19 \pm 1.53	97.16 \pm 1.24	95.97 \pm 1.55	95.34 \pm 1.42	96.18 \pm 1.02	97.61 \pm 1.68
FOL-A (E=2)	97.50 \pm 1.12	97.52 \pm 1.17	97.33 \pm 1.23	96.59 \pm 1.76	96.97 \pm 1.41	97.87 \pm 1.22
FOL-A (E=3)	97.53 \pm 0.76	97.70 \pm 0.98	97.99 \pm 0.93	96.90 \pm 1.09	97.47 \pm 1.11	98.20 \pm 1.03
FOL (E=1)	94.94 \pm 1.38	95.21 \pm 1.32	91.26 \pm 1.62	90.09 \pm 1.55	89.87 \pm 0.69	91.62 \pm 0.58
FOL (E=2)	95.23 \pm 1.35	95.57 \pm 0.72	91.60 \pm 1.29	91.23 \pm 1.57	91.77 \pm 0.83	95.21 \pm 1.49
FOL (E=3)	96.32 \pm 0.96	95.75 \pm 1.39	91.95 \pm 1.31	92.26 \pm 1.05	92.41 \pm 1.68	95.81 \pm 1.88
FOL-AN (E=1)	94.38 \pm 1.86	94.86 \pm 1.67	91.28 \pm 1.82	88.24 \pm 1.82	91.14 \pm 1.13	92.81 \pm 1.10
FOL-AN (E=2)	95.63 \pm 1.40	95.04 \pm 1.43	93.29 \pm 1.51	90.09 \pm 0.64	92.47 \pm 1.86	94.01 \pm 1.70
FOL-AN (E=3)	95.94 \pm 0.71	96.45 \pm 0.65	95.97 \pm 1.43	93.19 \pm 1.23	93.04 \pm 1.19	96.41 \pm 1.26
FOL-N (E=1)	93.44 \pm 1.68	94.68 \pm 1.79	88.59 \pm 2.31	85.76 \pm 1.85	89.22 \pm 0.93	91.62 \pm 1.19
FOL-N (E=2)	94.69 \pm 0.53	94.86 \pm 0.88	90.60 \pm 1.01	89.16 \pm 1.31	90.21 \pm 1.28	92.22 \pm 1.65
FOL-N (E=3)	95.31 \pm 1.49	95.21 \pm 0.98	91.95 \pm 0.97	90.71 \pm 0.59	90.51 \pm 1.21	94.61 \pm 0.73
DENSE	88.75 \pm 1.91	87.41 \pm 1.63	83.22 \pm 1.57	67.49 \pm 1.81	69.95 \pm 1.70	73.05 \pm 1.62
Co-Boosting	90.31 \pm 1.26	89.19 \pm 1.13	88.02 \pm 1.25	72.14 \pm 1.52	74.45 \pm 1.72	74.04 \pm 1.54
FedAvg (E=1)	73.19 \pm 1.73	73.94 \pm 1.96	68.18 \pm 2.02	60.21 \pm 1.73	62.03 \pm 1.95	66.26 \pm 1.62
FedAvg (E=2)	73.13 \pm 1.91	72.29 \pm 1.74	66.92 \pm 1.55	59.44 \pm 1.64	64.33 \pm 1.33	69.88 \pm 1.57
FedAvg (E=3)	74.61 \pm 1.54	71.58 \pm 1.16	68.48 \pm 1.23	63.70 \pm 0.71	65.16 \pm 1.14	67.82 \pm 0.92

 Table 2. Test accuracies (%) on Wildfire and Hurricane ($\psi \in \{0.5, 0.3, 0.1\}$), reported as mean \pm std.

Dataset	Wildfire			Hurricane		
Satellite #	32	43	48	8	26	44
Methods	$\psi = 0.5$	$\psi = 0.3$	$\psi = 0.1$	$\psi = 0.5$	$\psi = 0.3$	$\psi = 0.1$
Local	79.07 \pm 1.71	90.37 \pm 1.76	85.50 \pm 2.16	86.77 \pm 1.90	57.14 \pm 2.87	77.78 \pm 1.35
FOL-A (E=1)	95.35 \pm 1.42	94.07 \pm 1.89	90.63 \pm 1.92	95.04 \pm 1.70	90.48 \pm 1.57	88.89 \pm 1.92
FOL-A (E=2)	96.52 \pm 1.02	94.92 \pm 1.25	96.14 \pm 1.16	95.34 \pm 1.16	91.72 \pm 1.26	91.67 \pm 1.26
FOL-A (E=3)	97.67 \pm 0.71	95.76 \pm 0.85	96.88 \pm 1.01	95.87 \pm 1.03	93.65 \pm 1.14	94.44 \pm 0.87
FOL (E=1)	90.70 \pm 1.75	90.68 \pm 1.01	88.46 \pm 1.99	89.26 \pm 1.25	84.13 \pm 1.57	83.33 \pm 1.69
FOL (E=2)	91.96 \pm 1.09	91.53 \pm 1.78	90.63 \pm 1.77	90.08 \pm 1.74	85.71 \pm 1.38	84.43 \pm 1.92
FOL (E=3)	93.02 \pm 1.22	92.37 \pm 1.27	93.75 \pm 1.40	90.91 \pm 1.38	87.30 \pm 1.07	86.11 \pm 1.18
FOL-AN (E=1)	90.77 \pm 1.38	91.53 \pm 1.26	87.51 \pm 2.32	91.34 \pm 1.70	87.47 \pm 2.55	86.73 \pm 1.94
FOL-AN (E=2)	93.22 \pm 1.85	93.22 \pm 1.17	90.63 \pm 1.69	92.56 \pm 1.18	88.89 \pm 1.91	88.67 \pm 1.75
FOL-AN (E=3)	95.35 \pm 1.25	94.07 \pm 1.21	90.94 \pm 1.14	93.39 \pm 1.37	90.48 \pm 1.55	91.39 \pm 1.26
FOL-N (E=1)	86.05 \pm 1.96	88.14 \pm 1.67	85.13 \pm 1.92	85.95 \pm 1.95	76.19 \pm 1.73	80.56 \pm 2.11
FOL-N (E=2)	87.35 \pm 1.41	89.83 \pm 1.76	86.38 \pm 2.07	86.74 \pm 1.83	80.95 \pm 1.94	81.94 \pm 1.38
FOL-N (E=3)	90.54 \pm 1.51	90.06 \pm 1.59	88.47 \pm 1.37	87.60 \pm 1.49	82.54 \pm 1.76	83.37 \pm 1.56
DENSE	79.91 \pm 1.73	78.63 \pm 1.98	52.08 \pm 2.03	61.10 \pm 1.51	58.73 \pm 1.43	46.14 \pm 1.81
Co-Boosting	86.05 \pm 1.68	85.59 \pm 1.65	54.51 \pm 1.85	72.29 \pm 1.68	52.38 \pm 1.85	48.78 \pm 1.50
FedAvg (E=1)	53.11 \pm 1.82	63.25 \pm 1.87	35.33 \pm 2.76	66.12 \pm 1.50	41.27 \pm 1.99	46.14 \pm 1.72
FedAvg (E=2)	56.03 \pm 2.53	67.52 \pm 1.92	45.16 \pm 1.97	58.79 \pm 1.86	45.16 \pm 1.26	42.61 \pm 1.86
FedAvg (E=3)	51.07 \pm 1.93	66.10 \pm 2.05	42.86 \pm 1.53	60.33 \pm 1.24	44.44 \pm 1.76	43.33 \pm 1.46