

---

# Revisiting Ensembling in One-Shot Federated Learning

---

Youssef Allouah<sup>1</sup> Akash Dhasade<sup>1\*</sup> Rachid Guerraoui<sup>1</sup> Nirupam Gupta<sup>2</sup>

Anne-Marie Kermarrec<sup>1</sup> Rafael Pinot<sup>3</sup> Rafael Pires<sup>1</sup> Rishi Sharma<sup>1</sup>

<sup>1</sup>EPFL <sup>2</sup>University of Copenhagen

<sup>3</sup>Sorbonne Université and Université Paris Cité, CNRS, LPSM

## Abstract

Federated learning (FL) is an appealing approach to training machine learning models without sharing raw data. However, standard FL algorithms are iterative and thus induce a significant communication cost. One-shot federated learning (OFL) trades the iterative exchange of models between clients and the server with a single round of communication, thereby saving substantially on communication costs. Not surprisingly, OFL exhibits a performance gap in terms of accuracy with respect to FL, especially under high data heterogeneity. We introduce FENS, a novel federated ensembling scheme that approaches the accuracy of FL with the communication efficiency of OFL. Learning in FENS proceeds in two phases: first, clients train models locally and send them to the server, similar to OFL; second, clients collaboratively train a lightweight prediction aggregator model using FL. We showcase the effectiveness of FENS through exhaustive experiments spanning several datasets and heterogeneity levels. In the particular case of heterogeneously distributed CIFAR-10 dataset, FENS achieves up to a 26.9% higher accuracy over state-of-the-art (SOTA) OFL, being only 3.1% lower than FL. At the same time, FENS incurs at most  $4.3\times$  more communication than OFL, whereas FL is at least  $10.9\times$  more communication-intensive than FENS.

## 1 Introduction

FL is a widely adopted distributed machine learning (ML) approach, enabling clients to *collaboratively train* a common model over their collective data without sharing raw data with a central server [27]. Clients in FL engage in iterative parameter exchanges with the server over several communication rounds to train a model. While providing high accuracy, this process incurs substantial communication cost [19]. One-shot federated learning (OFL) [11] has been introduced to address the communication challenges in FL by reducing the exchange of models to a single round. Not surprisingly, this came with a loss of accuracy with respect to FL.

Typical OFL methods execute local training at the clients up to completion and form an ensemble of locally trained models at the server [7, 10, 46, 11]. The ensemble is distilled into a single model, through means of either auxiliary public dataset [10, 11] or synthetic data generated at the server [7, 14, 46]. While these OFL methods address communication challenges by reducing model exchanges to a single round, they often achieve lower accuracy compared to iterative FL. This is especially true when data distribution across clients is highly heterogeneous as OFL methods typically

---

\*Corresponding author <akash.dhasade@epfl.ch>

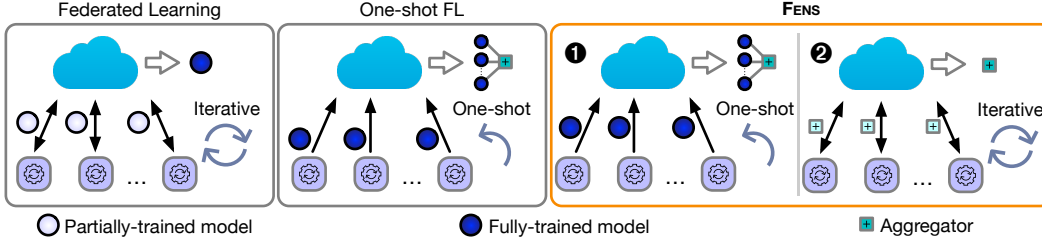


Figure 1: FENS in comparison to iterative and one-shot federated learning.

rely on simple prediction aggregation schemes such as averaging [11, 46], weighted averaging [7, 10] or voting [8].

We introduce FENS, a hybrid of OFL and standard FL. FENS aims to approach both the accuracy of iterative FL as well as the communication cost of OFL. Learning in FENS proceeds in two phases. In the first phase, similar to OFL, clients upload their locally-trained models to the server. Instead of using the traditional OFL aggregation, FENS employs a second phase of FL: the server constructs an ensemble with a prediction *aggregator model* stacked on top of the locally trained models. This advanced aggregation function is then trained by the clients in a lightweight FL training phase. The overall learning procedure is illustrated in Figure 1, alongside iterative and one-shot FL.

## 1.1 Our Contributions

We show for the first time, to the best of our knowledge, that a shallow neural network for the aggregator model suffices to satisfactorily bridge the gap between OFL and FL. Leveraging a shallow aggregator model enables two major benefits: first, it induces significantly lower communication cost in the iterative phase, and second, the iterative refinement of this aggregator model significantly improves accuracy over existing OFL methods. By utilizing elements from both OFL and FL in this novel ensembling scheme, FENS achieves the best of both worlds: accuracy of FL and communication efficiency of OFL.

Through extensive evaluations on several benchmark datasets (CIFAR-100, CIFAR-10, SVHN, and AG-News) across different heterogeneity levels, we demonstrate the efficacy of FENS in achieving FL-like accuracy at OFL-like communication cost. We extend our empirical evaluations to the FLamby benchmark [32], a realistic cross-silo FL dataset for healthcare applications. Our results show that in heterogeneous settings where even iterative FL algorithms struggle, FENS remains a strong competitor. We then conduct an extensive study of different aggregator models and highlight the accuracy vs. communication trade-off. Lastly, we show that FENS maintains high accuracy even with a comparable memory footprint.

To showcase FENS’s performance, we compare its accuracy and communication costs against Co-Boosting [7], a state-of-the-art OFL method, and FEDADAM [33], a popular iterative FL algorithm, as shown in Figure 2. These evaluations are performed on the CIFAR-10 dataset with 20 clients across three heterogeneity levels:  $\alpha = 0.01$  (very high),  $\alpha = 0.05$  (high), and  $\alpha = 0.1$  (moderate). Co-Boosting exhibits an accuracy gap of 13.7 – 26.9% compared to FEDADAM. FENS closes this accuracy gap, being only 0 – 3.1% lower than FEDADAM. To achieve this, FENS incurs only 3.8 – 4.3 $\times$  more communication than Co-Boosting whereas FEDADAM is 10.9 – 22.1 $\times$  more expensive than FENS.

## 1.2 Related Work

**One-shot Federated Learning.** Guha *et al.* [11] introduced one-shot FL, which limits communication to a single round. They proposed two main methods: (i) heuristic selection for final ensemble clients, and (ii) knowledge distillation (KD) for ensemble aggregation into a single model at the server using an auxiliary dataset. Subsequent methods based on KD [10, 22] require large, publicly available

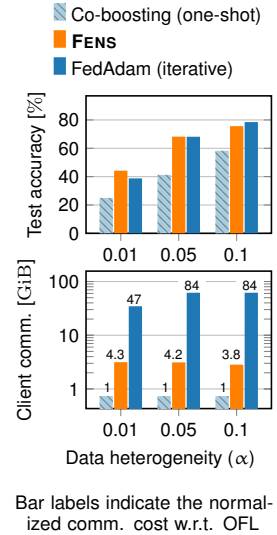


Figure 2: Test accuracy and communication cost of OFL, FENS and FL on CIFAR-10 dataset under high data heterogeneity.

datasets similar to local client data for good performance, which are often difficult to obtain [50]. To address this, synthetic data generation using generative adversarial networks (GAN)s has been proposed [7, 46]. The SOTA Co-Boosting algorithm [7] iteratively generates and refines synthetic data and the ensemble model. In FEDCVAE-ENS [14], clients train variational autoencoders (VAEs) locally and upload decoders to the server, which generates synthetic samples for classifier training. FEDOV [8] trains an open-set classifier at each client to predict “unknown” classes, with the server ensembling these models and using open-set voting for label prediction. Other OFL approaches either do not fully consider data heterogeneity [22, 37], or face difficulties under high data heterogeneity [48].

Another line of research in OFL focuses on aggregating fully trained client model parameters [42, 45]. PFNM [45] matches neurons across client models for fully-connected networks, while FEDMA [42] extends this to convolutional neural networks (CNNs) and LSTMs. However, the performance of these methods drops with more complex models. Few theoretical works exist, such as [18], which analyze global model loss for overparameterized ReLU networks. Despite the advances, OFL still exhibits accuracy gap with iterative FL. We show that FENS narrows this accuracy gap while preserving communication efficiency.

**Ensembles in Federated Learning.** Ensembles have been previously studied in FL for a variety of different goals. FEDDF [24] performs robust model fusion of client ensembles to support model heterogeneity. The FEDBE algorithm [5] uses Bayesian Model Ensemble to aggregate parameters in each global round, improving over traditional parameter averaging. Hamer *et al.* propose FEDBOOST [12] that constructs the ensemble using simple weighted averaging and analyze its optimality for density estimation tasks. However, these works are designed for standard FL and rely on substantial iterative communication. In the decentralized edge setting, [38] show that collaborative inference via neighbor averaging can achieve higher accuracy over local inference alone. However, they assume a setting where clients can exchange query data during inference and consider only IID data replicated on all edge devices. The idea of learning an aggregator model closely resembles late fusion techniques in multimodal deep learning [25]. The key difference is that FENS focuses on fusing single modality models trained on heterogeneous data under the communication constraints of federated settings.

## 2 Description of FENS

We consider a classification task represented by a pair of input and output spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively. The system comprises  $M$  clients, represented by  $[M] = \{1, \dots, M\}$  and a central server. Each client  $i$  holds a local dataset  $\mathcal{D}_i \subset \mathcal{X} \times \mathcal{Y}$ . For a model  $h_\theta : \mathcal{X} \rightarrow \mathcal{Z}$  parameterized by  $\theta \in \Theta \subseteq \mathbb{R}^d$ , each data point  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  incurs a loss of  $\ell(h_\theta(x), y)$  where  $\ell : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$ . Denoting by  $\mathcal{D} := \bigcup_{i \in [M]} \mathcal{D}_i$  the union of all local datasets, the objective is to solve the empirical risk minimization (ERM) problem:  $\min_{\theta \in \Theta} \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \ell(h_\theta(x), y)$ .

### 2.1 Federated Learning (FL) and One-shot FL (OFL)

FL algorithms, such as FedAvg [27], are iterative methods that enable the clients to solve the above ERM problem, without having to share their local data. In each iteration  $t$ , the server broadcasts the current model parameter  $\theta_t$  to a subset of clients  $S_t \subseteq [M]$ . Each client  $i \in S_t$  updates the parameter locally over its respective dataset  $\mathcal{D}_i$  using an optimization method, typically stochastic gradient descent (SGD). Clients send back to the server their locally updated model parameters  $\{\theta_t^{(i)}, i \in S_t\}$ . Lastly, the server updates the global model parameter to  $\theta_{t+1} := \frac{1}{|S_t|} \sum_{i \in S_t} \theta_t^{(i)}$ .

In One-shot Federated Learning (OFL), the iterative exchanges in FL are replaced with a *one-shot* communication of local models. Specifically, each client  $i$  seeks a model  $\theta^{(i)}$  that approximately solves the ERM problem on their local data:  $\min_{\theta \in \Theta} \frac{1}{|\mathcal{D}_i|} \sum_{(x,y) \in \mathcal{D}_i} \ell(h_\theta(x), y)$ , and sends  $\theta^{(i)}$  to the server. Upon receiving the local parameter  $\theta^{(i)}$ , corresponding to parametric model  $\pi_i = h_{\theta^{(i)}}$ , the server builds an ensemble model of the form  $\pi(x) = \sum_{i \in [M]} w_i \pi_i(x)$ . This ensemble model is then distilled into a single global model at the server using either a public dataset or synthetic data (generated by the server). Existing OFL algorithms choose weights  $w_1, \dots, w_M$  in three different ways: (i) uniformly at random [46], (ii) based on local label distributions [10], and (iii) dynamically adjusted based on generated synthetic data [7].

## 2.2 FENS

In FENS, the server builds the ensemble model using a generic aggregator  $f_\lambda : \mathcal{Z}^M \rightarrow \mathcal{Z}$ , parameterized by  $\lambda \in \Lambda \subset \mathbb{R}^q$  to obtain a global model  $\pi : \mathcal{X} \rightarrow \mathcal{Z}$  defined to be

$$\pi(x) := f_\lambda(\pi_1(x), \dots, \pi_M(x)). \quad (1)$$

In case of standard aggregators such as weighted averaging,  $q = M$  and  $\lambda \in (w \in \mathbb{R}_+^M \mid \sum_{i=1}^M w_i = 1)$ , and  $f_\lambda(\pi_1(x), \dots, \pi_M(x)) := \sum_{i \in [M]} \lambda_i \pi_i(x)$ . In general,  $f_\lambda$  can be a non-linear trainable model such as a neural network. The overall learning procedure in FENS comprises two phases:

1. **Local training and one-shot communication:** Each client  $i$  does local training to compute  $\theta^{(i)}$ , identical to OFL, and sends it to the server.
2. **Iterative aggregator training:** Upon receiving the local parameters  $\theta^{(i)}$ , the server reconstructs  $\pi_i := h_{\theta^{(i)}}$ , and obtains  $\hat{\lambda}$  that approximately solves the following ERM problem:

$$\min_{\lambda \in \Lambda} \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \ell(f_\lambda(\pi_1(x), \dots, \pi_M(x)), y). \quad (2)$$

The above ERM problem is solved using an iterative FL scheme (described above). For doing so, the server transmits the set of local models  $\{\pi_1, \dots, \pi_M\}$  to all the clients. The final model is given by  $\pi(x) := f_{\hat{\lambda}}(\pi_1(x), \dots, \pi_M(x))$ .

When solving for (2) using iterative FL, only the aggregator parameters are transferred between the server and the clients. As we show through experiments, in the subsequent section, training an aggregator model is much simpler than training the local models  $\pi_i$ , and a shallow neural network suffices for  $f_\lambda$ . Algorithms 1 and 2 (Appendix C) provide the pseudo for FENS.

**Connection with stacked generalization.** The use of a trainable aggregator corresponds to an instance of stacked generalization [44] in the centralized ensemble literature, wherein the aggregation function is regarded as *level 1* generalizer, while the clients' models are regarded as *level 0* generalizers. It has been shown that level 1 generalizer serves the role of correcting the biases of level 0 generalizers, thereby improving the overall learning performance of the ensemble [44]. While stacked generalization has been primarily studied in centralized settings, through FENS we show that this scheme can be efficiently extended to an FL setting.

## 3 Experiments

We split our evaluation into the following sections: (i) FENS vs OFL in Section 3.2; (ii) FENS vs FL and analysis of when FENS can match FL in Section 3.3; (iii) FENS on real-world cross-silo FLamby benchmark [32] in Section 3.4; (iv) FENS on language dataset in Section 3.5; (v) dissecting components of FENS in Section 3.6; and (vi) enhancing FENS efficiency in Section 3.7.

### 3.1 Experimental setup

**Datasets.** We consider three standard vision datasets with varying level of difficulty, including SVHN [30], CIFAR-10 [30] and CIFAR-100 [30], commonly used in several OFL works [7, 10, 46] as well as one language dataset AG-News [47]. Vision experiments involve 20 clients, except in the scalability study, where client numbers vary; and AG-News uses 10 clients. The original training splits of these datasets are partitioned across clients using the Dirichlet distribution  $\text{Dir}_{20}(\alpha)$ , in line with previous works [7, 10, 14]. The parameter  $\alpha$  determines the degree of heterogeneity, with lower values leading to more heterogeneous distributions (see Appendix A, Figure 9). For our experiments involving the realistic healthcare FLamby benchmark, we experiment with 3 datasets: Fed-Camelyon16, Fed-Heart-Disease, and Fed-ISIC2019. Table 5 (Appendix A) presents an overview of the selected tasks. The datasets consist of a natural non independent and identically distributed (non-IID) partitioning across clients. In FENS, each client performs local training using 90% of their local training data while reserving 10% for the iterative aggregator training. We observed that by splitting the datasets, we achieve better performance than reusing the training data for aggregator training. For fairness, OFL and FL baselines run with each client using 100% of their dataset for

local training. The testing set of each dataset is split (50-50%) for validation and testing. We use the validation set to tune hyperparameters and always report the accuracy on the testing split.

**One-shot FL baselines.** We compare FENS against 6 one-shot baselines: (i) one-round FEDAVG [27]; (ii) FEDENS [11], the first one-shot method constituting an ensemble with uniform weights; (iii) FEDKD [10], based on auxiliary dataset; (iv) one-shot version of FED-ET [6]; (v) the data-free FEDCVAE-ENS [14]; and (vi) Co-Boosting [7], based on synthetic data generation. We use the best-reported hyperparameters in each work for the respective datasets wherever applicable or tune them. Appendix B.1 provides additional details regarding the one-shot baselines.

**Iterative FL baselines.** For comparison with FL, we consider 6 algorithms: (i) FEDAVG [27]; (ii) FEDPROX [23]; (iii) FEDNOVA [43]; (iv) SCAFFOLD [20]; (v) FEDYOGI [33]; and (vi) FEDADAM [33]. We tune learning rates for each algorithm. In addition to these baselines, we implement gradient compression with FEDAVG STC, following the sparsification and quantization schemes of STC [29]. In particular, we set the quantization precision to 16-bit and sparsity level to 50%, to reduce the communication cost of FEDAVG by  $4\times$  and keep the remaining setup to the same as above baselines. For the FLamby benchmark experiments, we use the reported hyperparameters which were obtained after extensive tuning, except with one difference. The authors purposefully restricted the number of rounds to be approximately the same as the number of epochs required to train on pooled data (see [32]). Since this might not reflect true FL performance, we rerun all FL strategies to convergence using the reported tuned parameters. Precisely, we run up to  $10\times$  more communication rounds than evaluated in the FLamby benchmark. We include more details on FL baselines in Appendix B.2.

**FENS.**<sup>2</sup> For the CIFAR-10 and CIFAR-100 datasets, each client conducts local training for 500 epochs utilizing SGD as the local optimizer with an initial learning rate of 0.0025. For the SVHN and AG-News datasets, local training extends to 50 and 20 epochs respectively with a learning rate of 0.01. The learning rate is decayed using Cosine Annealing across all datasets. For the FLamby benchmark experiments, each client in FENS performs local training with the same hyperparameters as the client local baselines of FLamby. We experiment with two aggregator models, a 2-layer perceptron with ReLU activations and another that learns per-client per-class weights. We train the aggregator model using the FEDADAM algorithm where the learning rate is separately tuned for each dataset (Table 8, Appendix B.3). To reduce the communication costs corresponding to the ensemble download, we employ post-training model quantization at the server from FP32 to INT8. Appendix B.3 provides more details on FENS.

**Configurations.** In line with related work [10, 24, 35], we use ResNet-8 [13] as the client local model for our vision tasks and fine-tune DistilBert [34] for our language task. Our FLamby experiments use the same models as defined in the benchmark for each task (see Table 5, Appendix A). We report the average results across at least 3 random seeds. For iterative FL baselines, the communication cost corresponds to the round in which the best accuracy is achieved.

### 3.2 FENS vs OFL

To assess FENS’s efficacy, we experiment in non-IID settings, varying  $\alpha \in \{0.01, 0.05, 0.1\}$ , and present results across datasets and baselines in Table 1. Our observations reveal challenges for one-shot methods under high heterogeneity, with the optimal baseline differing across settings. FEDAVG with one round exhibits the poorest performance. While FEDCVAE-ENS maintains consistent accuracy across various heterogeneity levels and datasets, it struggles particularly with CIFAR-10 and CIFAR-100, indicating challenges in learning effective local decoder models. Regarding distillation-based methods, FEDKD and Co-Boosting demonstrate similar performance on SVHN and CIFAR-10. However, FEDKD outperforms Co-Boosting on CIFAR-100, facilitated by the auxiliary public dataset for KD while Co-Boosting arduously generates its synthetic transfer dataset. FED-ET improves over FEDENS and is also competitive to FEDKD. Notably, FENS consistently outperforms the best baseline in each scenario, except for SVHN at  $\alpha = 0.01$ , where FEDCVAE-ENS excels. FENS achieves significant accuracy gains, surpassing the best baseline by 11.4 – 26.9% on CIFAR-10 and 8.7 – 15.4% on CIFAR-100, attributed to its advanced aggregator model.

We chart the client communication costs incurred by all algorithms in Figure 3. The clients in FENS expend 3.6 –  $4.3\times$  more than one-shot algorithms owing to the ensemble download and iterative

<sup>2</sup>Source code available at: <https://github.com/sacs-epfl/fens>.

Table 1: FENS vs one-shot FL for various heterogeneity levels across datasets. The highest achieved accuracy is presented as bold and the top-performing baseline is underlined. The rightmost column presents the performance difference between FENS and the top-performing baseline.

Method	$\alpha$	FEDAVG	FEDENS	FEDKD	FED-ET	FEDCVAE	Co-Boosting	FENS	$\Delta$
CF-100	0.01	04.22 $\pm$ 0.33	16.59 $\pm$ 2.07	<u>28.98<math>\pm</math>4.55</u>	20.37 $\pm$ 1.53	07.84 $\pm$ 0.98	14.76 $\pm$ 2.14	<b>44.46<math>\pm</math>0.31</b>	+15.4
	0.05	05.53 $\pm$ 0.13	20.56 $\pm$ 3.51	<u>39.01<math>\pm</math>1.11</u>	33.20 $\pm$ 1.34	07.90 $\pm$ 0.83	20.28 $\pm$ 1.94	<b>49.70<math>\pm</math>0.86</b>	+10.6
	0.1	06.04 $\pm$ 0.92	27.41 $\pm$ 2.71	<u>42.38<math>\pm</math>0.78</u>	38.53 $\pm$ 1.04	08.05 $\pm$ 0.69	25.50 $\pm$ 0.65	<b>51.11<math>\pm</math>0.37</b>	+8.7
CF-10	0.01	10.35 $\pm$ 0.29	15.66 $\pm$ 6.11	18.59 $\pm$ 2.92	16.94 $\pm$ 7.88	<u>29.98<math>\pm</math>0.88</u>	24.97 $\pm$ 4.72	<b>44.20<math>\pm</math>3.29</b>	+14.2
	0.05	13.56 $\pm$ 3.33	39.56 $\pm$ 6.33	<u>38.84<math>\pm</math>6.03</u>	37.51 $\pm$ 2.87	29.40 $\pm$ 2.53	<u>41.25<math>\pm</math>5.88</u>	<b>68.22<math>\pm</math>4.19</b>	+26.9
	0.1	17.38 $\pm$ 0.22	48.40 $\pm$ 9.01	<u>64.14<math>\pm</math>5.17</u>	47.06 $\pm$ 2.31	32.10 $\pm$ 1.79	58.24 $\pm$ 3.54	<b>75.61<math>\pm</math>1.85</b>	+11.4
SVHN	0.01	11.35 $\pm$ 7.05	20.31 $\pm$ 3.49	23.62 $\pm$ 10.1	12.63 $\pm$ 6.23	<b>69.71<math>\pm</math>0.21</b>	25.32 $\pm$ 8.94	57.35 $\pm$ 12.6	-12.3
	0.05	12.85 $\pm$ 5.34	38.91 $\pm$ 7.28	37.41 $\pm$ 9.62	41.14 $\pm$ 6.66	<u>70.63<math>\pm</math>1.59</u>	46.36 $\pm$ 3.29	<b>76.76<math>\pm</math>2.98</b>	+6.1
	0.1	27.80 $\pm$ 9.83	51.99 $\pm$ 7.85	61.38 $\pm$ 3.90	58.91 $\pm$ 2.81	<u>72.38<math>\pm</math>0.28</u>	59.19 $\pm$ 7.06	<b>83.64<math>\pm</math>0.75</b>	+11.2

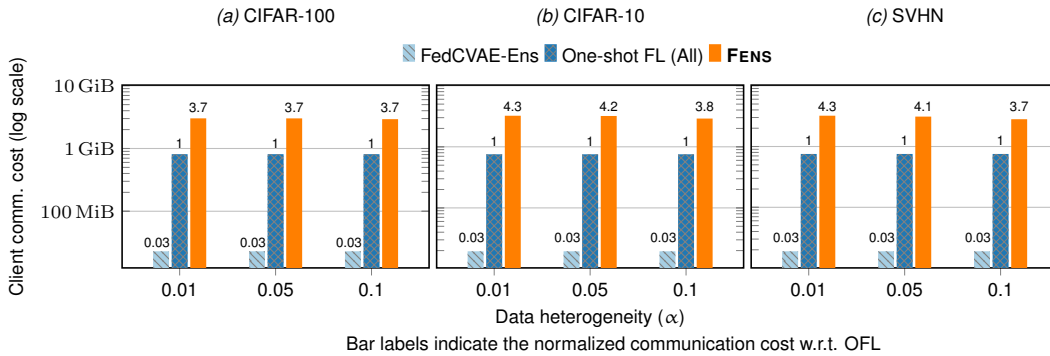


Figure 3: Total communication cost of FENS against OFL baselines. The clients in FENS expend roughly 3.7 – 4.3 $\times$  more than OFL in communication costs.

aggregator training. While these costs are greater than OFL, they are significantly lower than the costs incurred by iterative FL baselines as shown in Section 3.3. FEDCVAE-ENS has the lowest cost since the clients only upload the decoder component of their VAE model to the server. However, it also suffers from a significant performance gap with respect to other OFL baselines and FENS on the CIFAR-10 and CIFAR-100 datasets.

**Varying number of clients.** We also assess the performance of different baselines by varying the number of clients. We consider the CIFAR-10 dataset with  $\alpha = 0.1$  and present the results in Table 2. FENS achieves the best accuracy surpassing the best-performing baseline FEDKD by 5.9 – 11.4% in accuracy points. This again demonstrates the benefits of utilizing an advanced ensemble model with a trainable aggregator function.

Table 2: FENS vs one-shot FL on CIFAR-10 with varying number of clients.

$M$	FEDAVG	FEDENS	FEDKD	FED-ET	FEDCVAE	Co-Boosting	FENS	$\Delta$
5	32.00 $\pm$ 0.69	62.70 $\pm$ 8.45	<u>71.74<math>\pm</math>4.11</u>	47.83 $\pm$ 2.71	27.73 $\pm$ 1.50	60.86 $\pm$ 4.33	<b>77.70<math>\pm</math>2.62</b>	+5.9
10	24.00 $\pm$ 5.68	46.80 $\pm$ 1.28	<u>66.85<math>\pm</math>5.18</u>	49.22 $\pm$ 4.81	30.00 $\pm$ 2.21	54.10 $\pm$ 4.19	<b>76.44<math>\pm</math>2.64</b>	+9.5
20	17.38 $\pm$ 0.22	48.40 $\pm$ 9.01	<u>64.14<math>\pm</math>5.17</u>	47.06 $\pm$ 2.31	32.10 $\pm$ 1.79	58.24 $\pm$ 3.54	<b>75.61<math>\pm</math>1.85</b>	+11.4
50	18.57 $\pm$ 8.64	49.89 $\pm$ 3.82	<u>59.26<math>\pm</math>2.90</u>	51.59 $\pm$ 2.31	34.59 $\pm$ 1.26	51.38 $\pm$ 4.10	<b>70.53<math>\pm</math>1.10</b>	+11.2

### 3.3 FENS vs Iterative FL

We now compare the accuracy and communication cost of FENS against iterative FL baselines. After our extensive evaluation of all 6 iterative FL baselines (Table 10, Appendix D) on the CIFAR-10 dataset across various heterogeneity levels, we find that FEDADAM and FEDYOGI consistently perform the best. Hence, we show FEDADAM in our remaining evaluations. Figure 4 presents the results, showing FEDADAM as the representative of FL, FEDAVG STC as the representative of gradient compression, FEDKD as the representative of OFL, and FENS. Moreover, we show two versions of FEDKD, one additional with multi-round support to match the communication cost

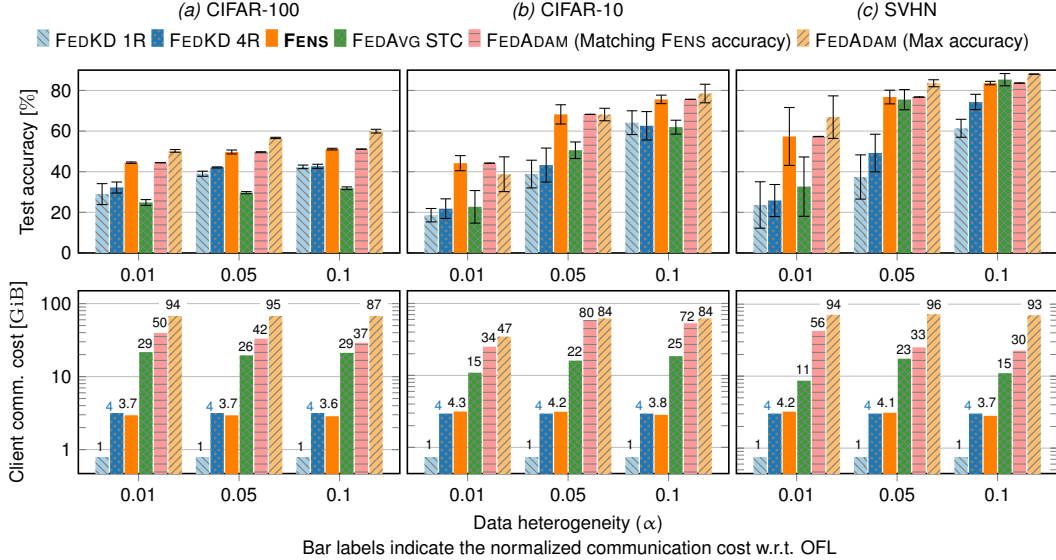


Figure 4: **FENS against iterative FL.** The R indicates the number of global rounds, signifying the multi-round version of the OFL baseline. FENS achieves accuracy properties of iterative FL (FEDADAM) with a modest increase in communication cost compared to OFL (FEDKD). Numerical accuracy results are included in Table 11 (Appendix D).

of FENS (details in Appendix B.2). We also show two versions of FEDADAM, one achieving the accuracy of FENS and the other with its maximum accuracy to facilitate effective comparison of communication costs.

We observe that FENS with its iteratively trained aggregator significantly closes the accuracy gap between OFL (FEDKD) and FL (FEDADAM) across all datasets and heterogeneity levels. Remarkably, the boost achieved is sufficient to match FEDADAM’s accuracy at  $\alpha = \{0.01, 0.05\}$  on the CIFAR-10 dataset. This comes at only a modest increase in communication costs which are  $\approx 4\times$  that of OFL across all cases. We observe that FEDADAM incurs  $30 - 80\times$  more communication than OFL to reach the same accuracy as FENS. Even adding multi-round support to FEDKD only marginally improves its performance. While the best accuracy achieved by FEDADAM still remains higher, it also comes at significant communication costs of  $47 - 96\times$  that of OFL. Furthermore, we observe that communication compression (FEDAVG STC) fails to preserve the accuracy of FL under high heterogeneity. Thus FENS achieves the best accuracy vs. communication trade-off, demonstrating accuracy properties of iterative FL while retaining communication efficiency of OFL.

### 3.3.1 When can FENS match iterative FL?

In this section, we aim to understand when FENS can match iterative FL. The performance of ensembles depends upon (i) the quality of local models; and (ii) data heterogeneity. The quality of local models in turn depends on the amount of local data held by the clients. As local models improve at generalizing locally, the overall performance of the ensemble is enhanced. In contrast, FL struggles to generate a good global model when the local datasets of clients significantly differ. Thus more volume of data does not analogously benefit FL due to high data heterogeneity. However, as heterogeneity reduces, FL excels and benefits significantly from collaborative updates. This suggests that FENS under *sufficiently large local datasets* and *high heterogeneity* can match iterative FL’s performance. We confirm this intuition through the following experiments on the SVHN dataset.

**Setup.** We study the performance of FL and FENS by progressively increasing the volume of data held by the clients. To this end, we consider the classification task on the SVHN dataset due to the availability of an extended training set of 604 388 samples, *i.e.*,  $\approx 10\times$  bigger than the default SVHN dataset. We then experiment with fractions ranging from 10 to 100% of the total training set. Each client locally utilizes 90% for one-shot local model training and reserves 10% for iterative aggregator training, similar to previous sections. We then compare FENS with FEDAVG (FL baseline) on three

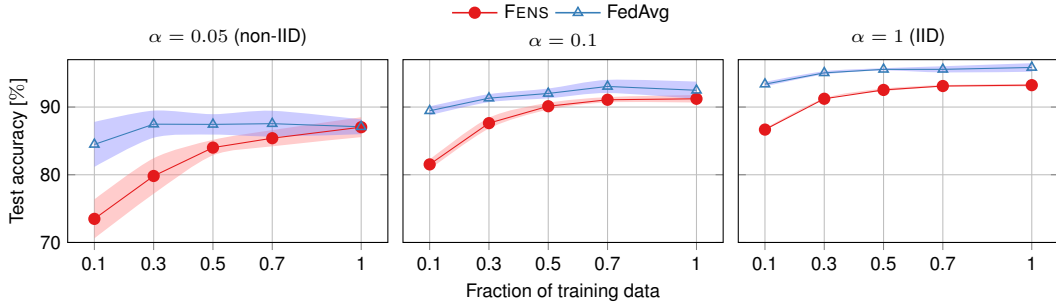


Figure 5: Accuracy of FENS for increasing dataset size. Performance of FENS rapidly increases as the data volume increases. At high data heterogeneity, FENS matches iterative FL’s accuracy.

levels of heterogeneity:  $\alpha = \{0.05, 0.1, 1\}$ , varying from highly non-IID to IID. We tune the learning rate for FEDAVG (details in Appendix B.2) and keep the remaining setup as in previous experiments.

**Results.** Figure 5 shows the results and confirms our prior insight behind the effective performance of FENS. Specifically, we observe that the growing volume of training data benefits FENS much more than FEDAVG. When the data distribution is homogeneous ( $\alpha = 1$ ), the performance of FENS improves faster than FEDAVG, but still remains behind. On the other hand, under high heterogeneity ( $\alpha = 0.01$ ), FENS quickly catches up with the performance of FEDAVG, matching the same accuracy when using the full training set. We conclude that under regimes of high heterogeneity and sufficiently large local datasets, FENS presents a practical alternative to communication expensive iterative FL.

### 3.4 Performance on real-world datasets

In this section, we evaluate the performance of FENS on the real-world cross-silo FLamby benchmark [32]. Specifically, we present 5 iterative baselines and 2 one-shot baselines along with the client local baselines (Figure 6). For the one-shot FEDAVG and FEDPROX OFL baselines, we additionally tune the number of local updates. FEDKD is infeasible in these settings since it requires a public dataset for distillation, unavailable in the medical setting. FEDCVAE-ENS and Co-Boosting are also infeasible due to the difficulty in learning good decoder models or synthetic dataset generators for medical input data, a conclusion supported by their poor performance on the comparatively simpler CIFAR-100 task (Table 1). Figure 6 shows the results with the first row comparing FENS against iterative FL algorithms and the second row against one-shot FL and the client local baselines.

When comparing to iterative FL, we observe that FENS is on par for the Fed-Heart-Disease dataset and performs better for Fed-Camelyon16. The iterative FL performance is affected by high heterogeneity [32] where the deterioration is more significant for Fed-Camelyon16, which learns on large breast slides ( $10000 \times 2048$ ) than for Fed-Heart-Disease, which learns on tabular data. In such scenarios of heterogeneity, FENS can harness diverse local classifiers through the aggregator model to attain good performance. On the Fed-ISIC2019 dataset, however, FENS does not reach the accuracy of iterative FL. Clients in Fed-ISIC2019 exhibit high variance in local data amounts and model performance, with the largest client having 12k samples and the smallest only 225 (Table 5, Appendix A). We thus speculate that the Fed-ISIC2019 dataset falls within the low local training fraction regime depicted in Figure 5, exhibiting a larger

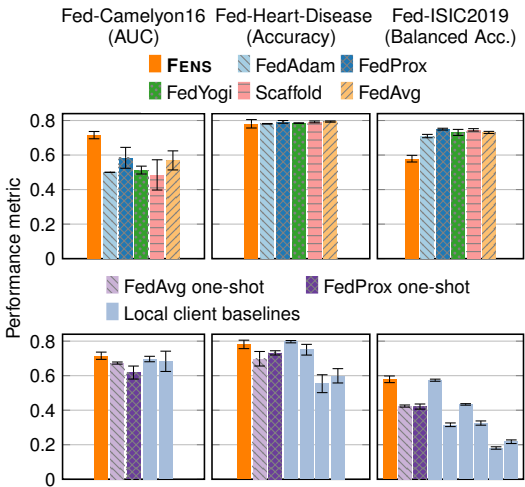


Figure 6: **FENS in FLamby.** FENS is on par with iterative FL (row-1), except when local models are weak (Fed-ISIC2019) while remaining superior in the one-shot setting (row-2). Numerical results included in Tables 12 to 17 (Appendix D).



accuracy gap compared to iterative FL. However, we note that FENS achieves superior performance over one-shot FEDAVG and one-shot FEDPROX, while performing at least as well as the best client local baseline across all datasets. Overall, these observations for FL algorithms have spurred new interest in developing a better understanding of performance on heterogeneous cross-silo datasets [32]. We show that FENS remains a strong competitor in such settings.

### 3.5 Performance on language dataset

We now study the performance of FENS on the AG-News dataset, comparing it against top-performing baselines FEDADAM and FEDKD in the iterative and one-shot categories, respectively. Figure 7 shows the results: at  $\alpha = 0.1$ , FEDKD achieves 71.5% accuracy, leaving a gap to FEDADAM at 82.3%. FENS effectively closes this gap, reaching 78.8%. As heterogeneity reduces at  $\alpha = 0.3$ , all algorithms achieve higher accuracies. FENS improves upon FEDKD from 79.3% to 84.5% while FEDADAM achieves 88.3%. Thus we observe consistent results on the language task as our vision benchmarks.

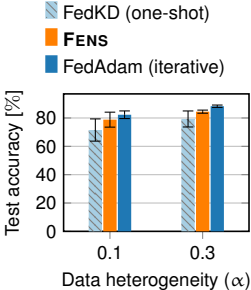


Figure 7: FENS on the AG-News dataset.

### 3.6 Dissecting FENS

We extensively evaluate various aggregation functions (details in Appendix B.4) on the CIFAR-10 dataset across diverse heterogeneity levels.

In particular, we assess static aggregation rules including averaging and weighted averaging, parametric aggregator models including a linear model, and a shallow neural network. We also evaluate an advanced version of voting [2] which involves computing competency matrices to reach a collective decision. In addition, we evaluate the Mixture-of-Experts (MoE) aggregation rule [36] where only the gating function is trained via federation. Figure 8 illustrates the accuracy, communication costs, and breakdown for all aggregations. Trained aggregator models outperform static aggregations, incurring additional costs for ensemble download and iterative training. The NN aggregator emerges as the top performer, offering the best accuracy vs. communication trade-off. Notably, the iterative training cost of the NN aggregator model for several rounds is lower than the OFL phase itself. Regarding accuracy, only MoE outperforms NN at  $\alpha = 0.01$ , where extreme heterogeneity induces expert specialization, while the trained gating network accurately predicts the right expert. However, MoE’s performance declines as heterogeneity decreases while its communication costs remain higher due to the size of the gating network.

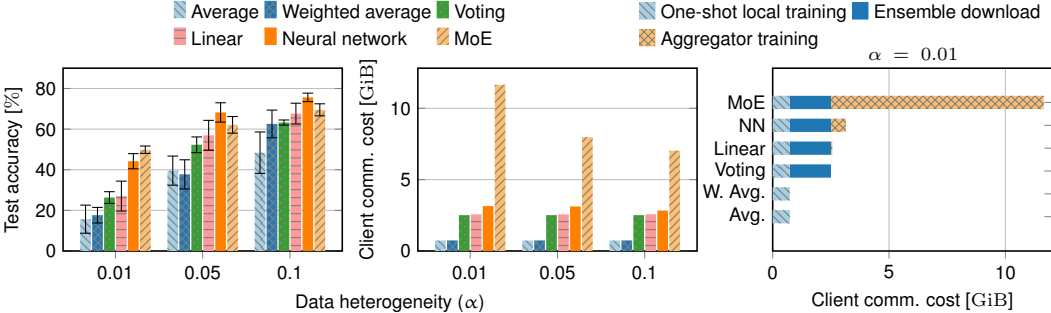


Figure 8: Accuracy of different aggregation functions on the CIFAR-10 dataset. NN offers the best accuracy vs. communication trade-off, with its iterative training taking up only a fraction of the total cost. Numerical accuracy values are included in Table 9 (Appendix D).

### 3.7 Enhancing FENS efficiency

The FENS global model comprises the aggregator model stacked atop the ensemble of client local models. Although FENS achieves strong accuracy, the ensemble model can be computationally and memory intensive. We used FP32 to INT8 quantization in our previous experiments which reduces the memory costs by  $4\times$  (Appendix B.3). In this section, we explore two additional approaches to reduce FENS’s overheads.

**What if we distill FENS into a single model?** To enable efficient inference, we can distill the FENS global model into a single model at the server using KD once the training is completed. Specifically, we distill the FENS ensemble model comprising 20 ResNet-8 client models and the shallow aggregator neural network into a single ResNet-8 model. Table 3 presents the results on the CIFAR-10 dataset for  $\alpha = \{0.01, 0.05, 0.1\}$  distilled using CIFAR-100 as the auxiliary dataset. We observe a slight accuracy drop arising from the process of distillation, which is standard behavior [39]. While we distill using the standard distillation algorithm [15], we note that this accuracy gap can be further reduced through the use of more advanced distillation methods [16, 49].

Table 3: Accuracy of FENS after distillation on the CIFAR-10 dataset.

Algorithm	$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.1$
FENS	44.20 $\pm$ 3.29	68.22 $\pm$ 4.19	75.61 $\pm$ 1.85
FENS distilled	43.81 $\pm$ 2.58	65.56 $\pm$ 3.25	71.59 $\pm$ 1.21

Table 4: FENS vs FEDADAM under similar memory footprint on CIFAR-10. DS stands for downsized.

Algorithm	$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.1$	Memory (MiB)
FENS	43.32 $\pm$ 3.92	66.99 $\pm$ 4.01	74.14 $\pm$ 1.63	377.03
FENS-DS	43.08 $\pm$ 3.41	64.59 $\pm$ 3.74	72.43 $\pm$ 2.16	17.95
FEDADAM	39.32 $\pm$ 7.85	68.74 $\pm$ 2.76	78.73 $\pm$ 3.55	18.85
FEDADAM-DS	29.69 $\pm$ 3.62	63.77 $\pm$ 0.23	72.25 $\pm$ 3.22	0.88

**What if we match the memory footprint of FENS to FEDADAM?** While the above approach enables efficient inference, we still need to mitigate training time costs. We now consider a downsized (DS) version of ResNet-8, where the width of a few layers is reduced so that the total size of the FENS downsized (FENS-DS) model approximately matches the size of the single model in FEDADAM. Table 4 presents the results on the CIFAR-10 dataset for various heterogeneity levels. Note that no quantization is considered for FENS and FENS-DS, hence the contrast in values with Table 3. Under a comparable memory footprint, FENS-DS remains competitive with the original FEDADAM, with only a slight drop in accuracy compared to FENS. On the other hand, using the downsized model as the global model in FEDADAM-DS results in a significant accuracy drop (from 39.32% to 29.69%) under high data heterogeneity ( $\alpha = 0.01$ ). Thus, the memory overhead of FENS can be alleviated while retaining its communication benefits without too much impact on accuracy.

## 4 Discussion and Conclusion

**Limitations.** One limitation is the memory required on client devices to store the ensemble model for aggregator training. We explored quantization and downsizing to mitigate this issue. Future work could investigate aggregator models that do not require access to all client models in the ensemble. This memory issue is only present during training; after training, FENS can be distilled into a single global model on the server, enabling efficient inference as shown in Section 3.7. Another limitation is the increased vulnerability to attacks during iterative aggregator training, unlike OFL, which limits the attack surface to one round. However, this only affects the aggregator model, since the client local models are still uploaded in one shot. Privacy can be further enhanced in FENS through techniques such as differential privacy [9] or trusted execution environments [28]. Specifically, clients can use differentially private SGD [1] for local training, providing a differentially private local model for the ensemble, while the aggregator training could leverage a differentially private FL algorithm [31].

**Benefits.** In addition to low communication costs and good accuracy, FENS provides three important advantages. First, it supports model heterogeneity, allowing different model architectures across federated clients [21]. Second, FENS enables rapid client unlearning [4], towards the goal of the *right to be forgotten* in GDPR [26]. In FENS, unlearning a client can be achieved by simply re-executing the lightweight aggregator training by excluding the requested clients’ model from the ensemble. This is more efficient than traditional FL, where disincorporating knowledge from a single global model can be costly. Lastly, if a small server-side dataset is available, such as a proxy dataset for bootstrapping FL [3, 19], FENS can train the aggregator model on the server. This makes FENS applicable in model market scenarios of OFL [7, 41] where clients primarily offer pre-trained models.

To conclude, we introduce FENS, a hybrid approach combining OFL and FL. FENS emphasizes local training and one-shot model sharing, similar to OFL, which limits communication costs. It then performs lightweight aggregator training in an iterative FL-like fashion. Our experiments on diverse tasks demonstrated that FENS is highly effective in settings with high data heterogeneity, nearly achieving FL accuracy while maintaining the communication efficiency of OFL. Additionally, FENS supports model heterogeneity, rapid unlearning, and is applicable to model markets.

## Acknowledgments

Nirupam is partly supported by Swiss National Science Foundation (SNSF) project 200021\_200477, “Controlling The Spread of Epidemics: A Computing Perspective”. The authors are thankful to Milos Vujanovic and Sayan Biswas for their helpful discussions, and to the anonymous reviewers of NeurIPS 2024 for their valuable time and constructive inputs that shaped the final version of this work.

## References

- [1] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318, 2016.
- [2] Ruth Ben-Yashar and Jacob Paroush. Optimal decision rules for fixed-size committees in polychotomous choice situations. *Social Choice and Welfare*, 18(4):737–746, 2001.
- [3] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmityr Huba, Alex Ingerman, Vladimir Ivanov, Chloé Kiddon, Jakub Konečný, Stefano Mazzocchi, Brendan McMahan, Timon Van Overveldt, David Petrou, Daniel Ramage, and Jason Roselander. Towards federated learning at scale: System design. In *MLSys*, 2019.
- [4] Lucas Bourtole, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *2021 IEEE Symposium on Security and Privacy (SP)*, pages 141–159. IEEE, 2021.
- [5] Hong-You Chen and Wei-Lun Chao. Fed{be}: Making bayesian model ensemble applicable to federated learning. In *International Conference on Learning Representations*, 2021.
- [6] Yae Jee Cho, Andre Manoel, Gauri Joshi, Robert Sim, and Dimitrios Dimitriadis. Heterogeneous ensemble knowledge transfer for training large models in federated learning. In Lud De Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 2881–2887. International Joint Conferences on Artificial Intelligence Organization, 7 2022. Main Track.
- [7] Rong Dai, Yonggang Zhang, Ang Li, Tongliang Liu, Xun Yang, and Bo Han. Enhancing one-shot federated learning through data and ensemble co-boosting. In *The Twelfth International Conference on Learning Representations*, 2024.
- [8] Yiqun Diao, Qinbin Li, and Bingsheng He. Towards addressing label skews in one-shot federated learning. In *International Conference on Learning Representations*, 2023.
- [9] Robin C Geyer, Tassilo Klein, and Moin Nabi. Differentially private federated learning: A client level perspective. *arXiv preprint arXiv:1712.07557*, 2017.
- [10] Xuan Gong, Abhishek Sharma, Srikrishna Karanam, Ziyang Wu, Terrence Chen, David Doermann, and Arun Innanjan. Preserving privacy in federated learning with ensemble cross-domain knowledge distillation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(11):11891–11899, Jun. 2022.
- [11] Neel Guha, Ameet Talwalkar, and Virginia Smith. One-shot federated learning. *arXiv preprint arXiv:1902.11175*, 2019.
- [12] Jenny Hamer, Mehryar Mohri, and Ananda Theertha Suresh. Fedboost: A communication-efficient algorithm for federated learning. In *International Conference on Machine Learning*, pages 3973–3983. PMLR, 2020.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

- [14] Clare Elizabeth Heinbaugh, Emilio Luz-Ricca, and Huajie Shao. Data-free one-shot federated learning under very high statistical heterogeneity. In *The Eleventh International Conference on Learning Representations*, 2023.
- [15] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*, 2015.
- [16] Fotis Iliopoulos, Vasilis Kontonis, Cenk Baykal, Gaurav Menghani, Khoa Trinh, and Erik Vee. Weighted distillation with unlabeled examples. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 7024–7037. Curran Associates, Inc., 2022.
- [17] Maximilian Ilse, Jakub Tomczak, and Max Welling. Attention-based deep multiple instance learning. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2127–2136. PMLR, 10–15 Jul 2018.
- [18] Divyansh Jhunjhunwala, Shiqiang Wang, and Gauri Joshi. Towards a theoretical and practical understanding of one-shot federated learning with fisher information. In *Federated Learning and Analytics in Practice: Algorithms, Systems, Applications, and Opportunities*, 2023.
- [19] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.
- [20] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020.
- [21] Daliang Li and Junpu Wang. Fedmd: Heterogenous federated learning via model distillation. *NeurIPS Workshop on Federated Learning for Data Privacy and Confidentiality*, 2019.
- [22] Qinbin Li, Bingsheng He, and Dawn Song. Practical one-shot federated learning for cross-silo setting. In Zhi-Hua Zhou, editor, *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 1484–1490. International Joint Conferences on Artificial Intelligence Organization, 8 2021. Main Track.
- [23] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*, 2:429–450, 2020.
- [24] Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust model fusion in federated learning. *Advances in Neural Information Processing Systems*, 33:2351–2363, 2020.
- [25] Kuan Liu, Yanen Li, Ning Xu, and Prem Natarajan. Learn to combine modalities in multimodal deep learning. *arXiv preprint arXiv:1805.11730*, 2018.
- [26] Alessandro Mantelero. The eu proposal for a general data protection regulation and the roots of the ‘right to be forgotten’. *Computer Law & Security Review*, 29(3):229–235, 2013.
- [27] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pages 1273–1282. PMLR, 2017.
- [28] Aghiles Ait Messaoud, Sonia Ben Mokhtar, Vlad Nitu, and Valerio Schiavoni. Shielding federated learning systems against inference attacks with arm trustzone. In *Proceedings of the 23rd ACM/IFIP International Middleware Conference*, pages 335–348, 2022.
- [29] Alessio Mora, Luca Foschini, and Paolo Bellavista. Structured sparse ternary compression for convolutional layers in federated learning. In *2022 IEEE 95th Vehicular Technology Conference: (VTC2022-Spring)*, pages 1–5, 2022.

- [30] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [31] Maxence Noble, Aurélien Bellet, and Aymeric Dieuleveut. Differentially private federated learning on heterogeneous data. In *International Conference on Artificial Intelligence and Statistics*, pages 10110–10145. PMLR, 2022.
- [32] Jean Ogier du Terrail, Samy-Safwan Ayed, Edwige Cyffers, Felix Grimberg, Chaoyang He, Regis Loeb, Paul Mangold, Tanguy Marchand, Othmane Marfoq, Erum Mushtaq, Boris Muzellec, Constantin Philippenko, Santiago Silva, Maria Teleńczuk, Shadi Albarqouni, Salman Avestimehr, Aurélien Bellet, Aymeric Dieuleveut, Martin Jaggi, Sai Praneeth Karimireddy, Marco Lorenzi, Giovanni Neglia, Marc Tommasi, and Mathieu Andreux. Flamby: Datasets and benchmarks for cross-silo federated learning in realistic healthcare settings. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 5315–5334. Curran Associates, Inc., 2022.
- [33] Sashank J. Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and Hugh Brendan McMahan. Adaptive federated optimization. In *International Conference on Learning Representations*, 2021.
- [34] Victor Sanh, L Debut, J Chaumond, and T Wolf. Distilbert, a distilled version of bert: Smaller, faster, cheaper and lighter. arxiv 2019. *arXiv preprint arXiv:1910.01108*, 2019.
- [35] Felix Sattler, Arturo Marban, Roman Rischke, and Wojciech Samek. Cfd: Communication-efficient federated distillation via soft-label quantization and delta coding. *IEEE Transactions on Network Science and Engineering*, 9(4):2025–2038, 2022.
- [36] Noam Shazeer, \*Azalia Mirhoseini, \*Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference on Learning Representations*, 2017.
- [37] MyungJae Shin, Chihoon Hwang, Joongheon Kim, Jihong Park, Mehdi Bennis, and Seong-Lyun Kim. Xor mixup: Privacy-preserving data augmentation for one-shot federated learning. *arXiv preprint arXiv:2006.05148*, 2020.
- [38] Nir Shlezinger, Erez Farhan, Hai Morgenstern, and Yonina C. Eldar. Collaborative inference via ensembles on the edge. In *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8478–8482, 2021.
- [39] Samuel Don Stanton, Pavel Izmailov, Polina Kirichenko, Alexander A Alemi, and Andrew Gordon Wilson. Does knowledge distillation really work? In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
- [40] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- [41] Manasi Vartak, Harihar Subramanyam, Wei-En Lee, Srinidhi Viswanathan, Saadiyah Husnoo, Samuel Madden, and Matei Zaharia. Modeldb: a system for machine learning model management. In *Proceedings of the Workshop on Human-In-the-Loop Data Analytics, HILDA '16*, New York, NY, USA, 2016. Association for Computing Machinery.
- [42] Hongyi Wang, Mikhail Yurochkin, Yuekai Sun, Dimitris Papailiopoulos, and Yasaman Khazaeni. Federated learning with matched averaging. In *International Conference on Learning Representations*, 2020.
- [43] Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. *Advances in neural information processing systems*, 33:7611–7623, 2020.
- [44] David H. Wolpert. Stacked generalization. *Neural Networks*, 5(2):241–259, 1992.
- [45] Mikhail Yurochkin, Mayank Agarwal, Soumya Ghosh, Kristjan Greenewald, Nghia Hoang, and Yasaman Khazaeni. Bayesian nonparametric federated learning of neural networks. In *International conference on machine learning*, pages 7252–7261. PMLR, 2019.

- [46] Jie Zhang, Chen Chen, Bo Li, Lingjuan Lyu, Shuang Wu, Shouhong Ding, Chunhua Shen, and Chao Wu. DENSE: Data-free one-shot federated learning. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [47] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.
- [48] Yanlin Zhou, George Pu, Xiyao Ma, Xiaolin Li, and Dapeng Wu. Distilled one-shot federated learning. *arXiv preprint arXiv:2009.07999*, 2020.
- [49] Yichen Zhu, Ning Liu, Zhiyuan Xu, Xin Liu, Weibin Meng, Louis Wang, Zhicai Ou, and Jian Tang. Teach less, learn more: On the undistillable classes in knowledge distillation. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [50] Zhuangdi Zhu, Junyuan Hong, and Jiayu Zhou. Data-free knowledge distillation for heterogeneous federated learning. In *International Conference on Machine Learning*, pages 12878–12889. PMLR, 2021.

## Organization of the Appendix

- A Datasets
- B Additional Experimental Details
  - B.1 One-shot FL baselines
  - B.2 Iterative FL baselines
  - B.3 FENS
  - B.4 Aggregation rules
- C FENS pseudo code
- D Numerical results
- E Compute resources
- F Broader impact

## A Datasets

As mentioned in Section 3.1, we focus on classification tasks and experiment with 3 datasets in FLamby benchmark including Fed-Camelyon16, Fed-Heart-Disease and Fed-ISIC2019. Table 5 overviews the selected tasks in this work.

Table 5: Overview of selected datasets and tasks in FLamby. We defer additional details to [32].

DATASET	INPUT ( $x$ )	PREDICTION ( $y$ )	TASK TYPE	# CLIENTS	# EXAMPLES PER CLIENT	MODEL	METRIC
Fed-Camelyon16	Slides	Tumor on Slide	Binary Classification	2	239, 150	DeepMIL [17]	AUC
Fed-Heart-Disease	Patient Info.	Heart Disease	Binary Classification	4	303, 261, 46, 130	Logistic Regression	Accuracy
Fed-ISIC2019	Dermoscopy	Melanoma Class	Multi-class Classification	6	12413, 3954, 3363, 225, 819, 439	EfficientNet [40] + Linear layer	Balanced Accuracy

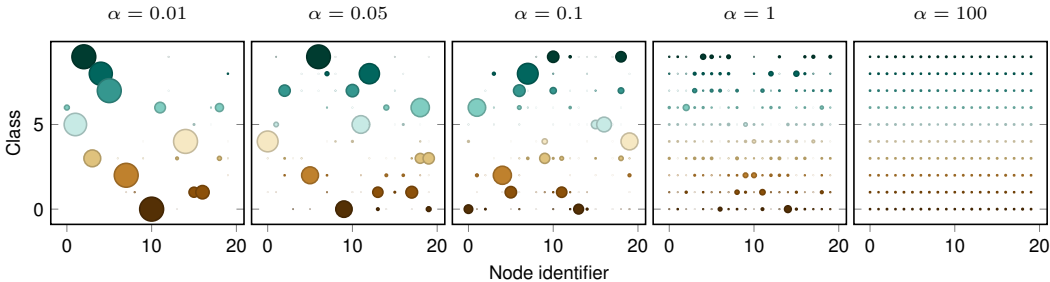


Figure 9: Visualizing the effect of changing  $\alpha$  on the CIFAR-10 dataset. Dot size corresponds to the number of samples of a given class in a given node.

## B Additional Experimental Details

### B.1 One-shot FL baselines

The client local training for all OFL baselines (except FEDCVAE-ENS) as well as FENS is conducted alike, using the parameters reported in [10] for the CIFAR-10 and CIFAR-100 datasets. For the SVHN and AG-News datasets, local training is conducted using the SGD optimizer with a learning rate of 0.01 for 50 and 20 local epochs respectively and decayed using Cosine Annealing. All vision datasets use a batch size of 16 while AG-News uses a batch-size of 8 for local training. For FEDCVAE-ENS, we use the same CVAE architecture and local training parameters as reported by the authors [14]. However, for the distillation at the server, we use ResNet-8 as the classifier model at the server to maintain fairness with other baselines. For the distillation phase of FEDKD,

we use the setup described by the authors [10] without inference quantization. For the one-shot version of FED-ET, we set the diversity regularization parameter to the best value of  $\lambda = 0.05$ . For both FEDKD and FED-ET, we use CIFAR-10, CIFAR-100, and TinyImageNet as the auxiliary datasets for distillation for SVHN, CIFAR-10, and CIFAR-100 datasets respectively. We consider a 60 – 40% split for the AG-News dataset where local training is conducted on the 60% split while the remaining 40% is treated as the auxiliary dataset for distillation at the server in FEDKD. Distillation in Co-Boosting using synthetically generated data also uses the best-reported hyperparameters [7]. For one-round FEDAVG, we additionally tune the number of local epochs performed before aggregation by considering  $\{1, 2, 5, 10, 15, 20\}$  epochs.

## B.2 Iterative FL baselines

We experiment with 6 different baselines including FEDAVG, FEDPROX, FEDNOVA, FEDADAM, FEDYOGI, and SCAFFOLD. We perform extensive hyperparameter tuning on the CIFAR-10 dataset for all levels of heterogeneity as detailed below. The server assumes full client participation, *i.e.*, all clients participate in each round. For our vision benchmarks, each client performs 2 local epochs per round using a batch size of 16. For our language task, clients train for 50 local steps in each round using a batch-size of 8. We run FL training until convergence and report the best accuracy achieved. We found all algorithms to converge in less than 100 communication rounds on the vision tasks and in 150 rounds on the language task.

**Hyperparameter tuning.** Below we describe our tuning procedure derived from several previous works [23, 32, 33]. For the FEDAVG algorithm, we tune the client learning rate ( $\eta_l$ ) over the values  $\{0.1, 0.01, 0.001, 0.0001\}$  separately for every  $\alpha \in \{0.01, 0.05, 0.1\}$ . For the FEDPROX algorithm, our grid space was  $\{0.1, 0.01\}$  and  $\{1, 0.1, 0.01\}$  for the client learning rate ( $\eta_l$ ) and the proximal parameter ( $\mu$ ) respectively. This was again separately tuned for every value of  $\alpha \in \{0.01, 0.05, 0.1\}$ . For the FEDYOGI and the FEDADAM algorithm, we consider the grid space of  $\{0.1, 0.01, 0.001, 0.0001\}$  and  $\{10, 1, 0.1, 0.01, 0.001\}$  for the client learning rate ( $\eta_l$ ) and the server learning rate ( $\eta_s$ ) respectively. This explodes the search significantly when tuning for every value of  $\alpha$ . From our tuning results for the FEDAVG and the FEDPROX algorithm, we noticed that the optimal parameter values were the same within the following two subgroups of  $\alpha - \{0.01, 0.05\}$  and  $\{0.1\}$  (Table 6). Hence, to keep the tuning tractable, we tune only for one  $\alpha$  in each subgroup and reuse the values for other alphas within the same subgroup. For the FEDNOVA algorithm, we use the version with both client and global momentum which was reported to perform the best [43]. We consider the search space  $\{0.005, 0.01, 0.02, 0.05, 0.08\}$  for the client learning rate ( $\eta_l$ ) as done by the authors [43] and tune separately for every value of  $\alpha$ . Finally, for the SCAFFOLD algorithm, we consider the search spaces  $\{0.1, 0.01, 0.001, 0.0001\}$  and  $\{1.0, 0.1, 0.01\}$  for  $\eta_l$  and  $\eta_s$  respectively and also tune separately for every  $\alpha$ . We first conduct the tuning procedure on the CIFAR-10 dataset and report the obtained hyperparameters in Table 6.

Table 6: Best hyperparameters obtained for the different algorithms on the CIFAR-10 dataset.

$\alpha$	FEDAVG	FEDPROX		FEDYOGI		FEDADAM		FEDNOVA	SCAFFOLD	
	$\eta_l$	$\eta_l$	$\mu$	$\eta_l$	$\eta_s$	$\eta_l$	$\eta_s$	$\eta_l$	$\eta_l$	$\eta_s$
0.01	0.01	0.01	0.01	0.01	0.01	-	-	0.0001	0.0001	1.0
0.05	0.01	0.01	0.01	-	-	0.01	0.01	0.02	0.01	0.1
0.1	0.1	0.1	0.01	0.01	0.01	0.01	0.01	0.005	0.01	1.0

We run all iterative FL baselines using the above parameters and present the results in Table 10. Based on our results in the Table 10 for CIFAR-10, we observe that FEDADAM and FEDYOGI consistently perform the best. Hence, to keep the experiments tractable, we tune and present just FEDADAM as a representative of the iterative FL family for our evaluations on the SVHN, CIFAR-100 datasets in Section 3.3 and AG-News in Section 3.5. For FEDADAM on the AG-News dataset, we note that the training is conducted using only the 60% split (see Appendix B.1) to achieve a fair comparison with FEDKD. For our experiments involving the extended SVHN dataset in Section 3.3.1, we again tune the client learning rate ( $\eta_l$ ) for the FEDAVG algorithm over the search space  $\{0.1, 0.01, 0.001, 0.0001\}$  separately for every  $\alpha \in \{0.01, 0.05, 0.1\}$ .



**FEDAVG with gradient compression.** To implement FEDAVG with gradient compression, we followed the sparsification and quantization schemes of STC [29]. We use the quantization level of 16-bit and sparsity of 50%. This results in a communication cost reduction of  $4\times$  against standard FEDAVG in every round. For each dataset and heterogeneity level, we tune the learning rate over the search space  $\{0.1, 0.05, 0.01, 0.001\}$ . We keep the remaining setup the same as FEDAVG.

**Multi-round FEDKD.** Since FENS incurs a communication cost four times that of OFL, we also evaluate FEDKD with multi-round support. We explore two approaches: *i*) pre-training for 3 rounds using FEDAVG, then applying FEDKD, and *ii*) using FEDKD followed by 3 fine-tuning rounds with FEDAVG. In the first case, each FEDKD client begins training from the global model produced by FEDAVG, while in the second, FEDAVG starts from the FEDKD model. We observe that pre-training with FEDAVG offers little improvement, likely due to the forgetting effect from multiple local training epochs, whereas fine-tuning with FEDAVG boosts FEDKD performance. For our experiments in Section 3.3, we thus present the multi-round version of FEDKD with fine-tuning. We remark that this multi-round support is still insufficient to match the performance of FENS, which achieves significantly higher accuracy as shown in Table 7 while incurring similar communication costs.

Table 7: FEDKD under multi-round support on the CIFAR-10 dataset.

Dataset	$\alpha$	FEDKD	3 rounds FEDAVG + FEDKD	FEDKD + 3 rounds FEDAVG	FENS
CF-10	0.01	18.59 $\pm$ 2.92	19.31 $\pm$ 3.16	21.81 $\pm$ 4.27	<b>44.20</b> $\pm$ 3.29
	0.05	38.84 $\pm$ 6.03	37.92 $\pm$ 7.04	43.26 $\pm$ 7.39	<b>68.22</b> $\pm$ 4.19
	0.1	64.14 $\pm$ 5.17	63.25 $\pm$ 6.22	62.61 $\pm$ 6.17	<b>75.61</b> $\pm$ 1.85

### B.3 FENS

Clients in FENS perform local training similar to OFL baselines as described in Section 3.1 and appendix B.1. Let  $\mathbf{z}_j \in \mathbb{R}^C$  denote the logits obtained from each client model  $\pi_j$  for all  $j \in [M]$  on a given input where  $C$  is the number of classes. We use one of the two aggregator models as follows. The first one is a multilayer perceptron using ReLU activations and a final classifier head as follows:  $f = \mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{z})$  where  $\mathbf{W}_1 \in \mathbb{R}^{MC \times k}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{k \times C}$ ,  $\mathbf{z} = \text{concat}(\mathbf{z}_1, \dots, \mathbf{z}_M) \in \mathbb{R}^{MC}$  is the concatenated logit vector and  $\sigma(x) = \max\{x, 0\}$  is the ReLU function. The parameter  $k$  determines the number of units in the hidden layer of this perceptron model. The second one is  $f = \sum_{i=1}^M \lambda_i \odot \mathbf{z}_i$  where  $\lambda_1, \dots, \lambda_M \in \mathbb{R}^C$  are weight vectors and  $\odot$  denotes coordinate-wise product. This model learns per-class per-client weights as the model parameters. For all datasets, the aggregator model is trained using the FEDADAM algorithm where the learning rate is separately tuned for each dataset. Table 8 presents the tuned learning rate and tuned training parameters per dataset.

Table 8: Aggregator training in FENS. We use FEDADAM as the FL algorithm with the following client ( $\eta_l$ ) and server ( $\eta_s$ ) learning rates. The parameter  $k$  corresponds to the weight matrices  $\mathbf{W}_1$  and  $\mathbf{W}_2$ .

Dataset	Aggregator Model	$k$	$\eta_l$	$\eta_s$	Batch Size	Local Steps	Global Rounds
CIFAR-10	$f = \mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{z})$	40	1.0	0.001	128	1	500
CIFAR-100	$f = \sum_{i=1}^M \lambda_i \odot \mathbf{z}_i$	—	1.0	0.003	128	1	1000
SVHN	$f = \mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{z})$	40	0.1	0.01	128	1	500
AG-News	$f = \mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{z})$	40	1.0	0.001	128	1	500
Fed-Heart-Disease	$f = \sum_{i=1}^M \lambda_i \odot \mathbf{z}_i$	—	0.1	0.1	2	5	50
Fed-Camelyon16	$f = \sum_{i=1}^M \lambda_i \odot \mathbf{z}_i$	—	1.0	0.0005	64	1	2000
Fed-ISIC2019	$f = \mathbf{W}_2^T \sigma(\mathbf{W}_1^T \mathbf{z})$	24	1.0	0.001	16	1	2500

**Model quantization in FENS.** Clients in FENS incur a critical cost of downloading the ensemble model from the server to initiate the aggregator training. To reduce the communication burden on the clients, the server employs post-training model quantization of all received client local models from FP32 to INT8, reducing the download costs by  $4\times$ . Alternatively, the quantization can also be executed on the client side. The quantization results in a drop of  $\approx 1 - 2\%$  test accuracy for every client model compared to the corresponding non-quantized model. However, the subsequent

aggregator training phase in FENS provides resilience to this drop in the accuracy of client models in the ensemble. In fact, we observe that the final accuracy achieved after aggregator training is slightly higher when using the quantized models as compared to unquantized models due to the regularising effect of quantization on generated logits. The model quantization also provides reduced memory usage on client devices during aggregator training. We use the standard PyTorch quantization library to implement quantization in FENS.

#### B.4 Aggregation rules

**Averaging.** Averaging corresponds to the following static aggregation rule:  $f = \sum_{i=1}^M \lambda_i \odot z_i$  where  $\lambda_i = [\frac{1}{M}, \dots, \frac{1}{M}]$  and  $\odot$  denotes coordinate-wise product.

**Weighted Averaging.** In weighted averaging, the  $\lambda_1, \dots, \lambda_M \in \mathbb{R}^C$  are typically assigned based on local training dataset statistics. In FEDKD [10],  $\lambda_i = [\frac{n_i^1}{\sum_{i \in [M]} n_i^1}, \frac{n_i^2}{\sum_{i \in [M]} n_i^2}, \dots, \frac{n_i^C}{\sum_{i \in [M]} n_i^C}]$  where  $n_i^j$  corresponds to the number of samples of class  $j$  with client  $i$ .

**Linear.** This aggregation corresponds to having a single learnable scalar weight for each client  $f = \sum_{i=1}^M w_i z_i$ . The learnable parameters in this case consist of the vector  $[w_1, w_2, \dots, w_M]^T$ .

**Neural network (NN).** Let  $z_j \in \mathbb{R}^C$  denote the logits obtained from each client model  $\pi_j$  for all  $j \in [M]$  on a given input where  $C$  is the number of classes. The NN aggregation corresponds to any neural network-based model  $f : \mathcal{Z}^M \rightarrow \mathcal{Z}$  that operates on the logits produced by the client models. Denoting  $z = \text{concat}(z_1, \dots, z_M) \in \mathbb{R}^{MC}$  as the concatenated vector of logits,  $f$  corresponds to the following 2 layer neural network  $f = \mathbf{W}_2^T \sigma(\mathbf{W}_1^T z)$  where  $\mathbf{W}_1 \in \mathbb{R}^{MC \times k}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{k \times C}$  and  $\sigma(x) = \max\{x, 0\}$  is the ReLU function. Here  $k$  determines the number of units in the hidden layer and controls the expressivity of the network. This aggregation is much more powerful than the previously mentioned aggregations, owing to its ability to discern complex patterns across all  $M \times C$  logits. The learnable parameters comprise the weight matrices  $\{\mathbf{W}_1, \mathbf{W}_2\}$ .

**Polychotomous Voting.** Polychotomous voting [2], was originally developed in social choice theory to reach a collective decision when offered  $C$  alternatives (classification labels in our case) in a committee of  $M$  experts (clients in our case). This method requires as input: (i) classwise ‘‘competency’’ scores of each client:  $P_i^c(r)$  indicating the probability of  $i^{\text{th}}$ -client to vote for label  $c$  when the ground truth is  $r$ ; (ii)  $p_{\text{prior}}(r)$ : prior probability distribution over correct alternatives  $r$ ; and (iii) the ‘‘benefit’’ vector of the committee:  $B(c|r)$  indicating the committee’s benefit in choosing label  $c$  when the correct class is  $r$ . Given this information, Ben-Yashar and Paroush [2] derive a criterion for the optimal decision that maximizes expected utility. This criterion is not computed using a closed-form expression, and we generically express it as

$$f(\pi_1, \dots, \pi_M; \mathbf{P}_1, \dots, \mathbf{P}_M; \mathbf{p}_{\text{prior}}; B) \quad (3)$$

where  $f$  corresponds to a procedure that evaluates and compares benefits for each choice  $c \in [C]$  given the competency matrices  $\{\mathbf{P}_i\}_{i=1}^M$ , the priors  $\mathbf{p}_{\text{prior}}$  and the benefit function  $B$ . Since the competency matrices for each client model are not directly available in our distributed setting, we learn them by federation in the network. More specifically, each client computes the competency matrix for every client model in the ensemble on its local data and transfers them to the server. The server then aggregates the received competency matrix to produce the final competency matrix per client to be used in decision-making. We further use a simple benefit function for our experiments that assigns  $B(c/r) = 1$  when  $c = r$  (correct choice) and  $B(c/r) = 0$  when  $c \neq r$  (incorrect choice) and set the prior to be uniform over the set of labels.

**Mixture-of-Experts (MoE).** The MoE aggregation [36] considers both the input and the logits in the following form:  $f = \sum_{i \in [M]} G(x)_i \pi_i(x)$ . Thus, MoE aggregation  $f : (\mathcal{X}, \mathcal{Z}^M) \rightarrow \mathcal{Z}$  is more expressive compared to other aggregations which only consider logits  $f : \mathcal{Z}^M \rightarrow \mathcal{Z}$ . Here,  $G : \mathcal{X} \rightarrow [0, 1]^M$  is called a gating network that generates scalar weights for every expert  $\pi_1, \dots, \pi_M$  based on the input  $x \in \mathcal{X}$ . In FENS with MoE aggregation, only the gating network is trained via the federation in the network while  $\{\pi_i\}_{i=1}^M$  correspond to the locally trained client models. We use a simple CNN with two convolutional blocks comprising ReLU activation and max pooling layers followed by 2 FC layers with ReLU activation, which in turn is followed by the final classification head. Despite its expressivity, learning a good gating network incurs significant communication costs and remains difficult under heterogeneous data in federated settings.

---

**Algorithm 1:** FENS from server perspective

---

**Require :**  $M$  clients, boolean *quantize*

```
1 Procedure FENS_SERVER():
2   Initialize model  $\pi$  with parameters  $\theta$  for local training
3   Send  $\theta$  to all  $M$  clients
4   Receive parameters of locally trained models  $\{\theta^{(i)}, i \in [M]\}$ 
5   if quantize then
6      $\theta^{(i)} \leftarrow \text{quantization\_alg}(\theta^{(i)}, \text{FP32}, \text{INT8}) \forall i \in [M]$ 
7   Send  $\{\theta^{(i)}, i \in [M]\}$  to all  $M$  clients1
8   Initialize aggregator model  $f_\lambda$  with parameters  $\lambda_0$ 
9   for  $t = 0, 1, \dots$  until convergence do
10    Select  $S_t \subseteq [M]$  and send them the aggregator parameters  $\lambda_t$ 
11    Receive updated parameters  $\{\lambda_t^{(i)}, i \in S_t\}$ 
12    Update global aggregator model  $\lambda_{t+1} := \frac{1}{|S_t|} \sum_{i \in S_t} \lambda_t^{(i)}$ 2
```

---

---

**Algorithm 2:** FENS from the clients perspective

---

**Require :** Local dataset  $\mathcal{D}_i$ , loss function  $\ell$ , local steps  $K$  and client learning rate  $\eta_l$ 

```
1 Procedure FENS_CLIENT():
2   Split  $\mathcal{D}_i$  randomly into 90%  $\mathcal{D}_{i1}$  and 10%  $\mathcal{D}_{i2}$ 
3   Receive parameters  $\theta$  from the server
4   Obtain  $\theta^{(i)}$  through local training of  $\theta$  on  $\mathcal{D}_{i1}$ 
5   Send converged model parameters  $\theta^{(i)}$  to server (one-shot)
6   Receive  $\{\theta^{(i)}, i \in [M]\}$  from the server
7   while Receive aggregator model parameters  $\lambda_t$  from the server do
8     Initialize  $\lambda_t^{(i)} \leftarrow \lambda_t$ 
9     for  $k = 0, 1, \dots, K$  do
10      Sample mini-batch  $b \in \mathcal{D}_{i2}$ 
11       $\ell_b \leftarrow \frac{1}{|b|} \sum_{(x,y) \in b} \ell \left( f_{\lambda_t^{(i)}}(\pi_1(x), \dots, \pi_M(x)), y \right)$ 
12       $\lambda_t^{(i)} \leftarrow \lambda_t^{(i)} - \eta_l \nabla \ell_b$ 
13    Send  $\lambda_t^{(i)}$  back to the server
```

---

## C FENS Algorithm

Algorithm 1 outlines the role of the server in FENS. The process begins with the server initializing the parameter  $\theta$  corresponding to the parametric model  $\pi = h_\theta$ , which it sends to all clients for local training (lines 2-3). Once clients complete their local training, they return their updated models to the server (line 4). If quantization is enabled, the server quantizes all the local models from FP32 to INT8 using a quantization algorithm (line 6). The server then redistributes all models back to the clients (line 7), enabling each to contribute to the aggregation process that follows. In the final stage, the server iteratively trains an aggregator model in FL fashion, which is designed to combine client models into a single, improved global model (lines 9-12). During each round, the server selects a subset of clients to update the aggregator model (line 10), refining it further with each iteration until convergence.

Algorithm 2 explains the client-side process in FENS. Each client starts by splitting its local dataset into two parts: one for one-shot local training and a smaller part for the iterative aggregator training (line 2). Using the received model, clients first train on  $\mathcal{D}_{i1}$  (line 4) and send their converged local model back to the server (line 5). In subsequent rounds, clients receive from the server the aggregator model (line 7) which they refine it locally using  $\mathcal{D}_{i2}$  (lines 8-12). Finally, clients send the updated aggregator parameters back to the server (line 13), contributing to the global aggregation process.

<sup>1</sup>The server can potentially send all parameters except the client’s own to reduce costs.

<sup>2</sup>The server can use any FL algorithm for aggregator training, FEDAVG shown for simplicity.

## D Numerical Results

In this section, we include the numerical values in Tables 9 to 17 corresponding to the plots presented in Sections 3.3, 3.4 and 3.6 for a complete reference.

Table 9: FENS aggregation methods on CIFAR-10. Results of Figure 8.

Algorithm	$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.1$
Averaging	15.66 $\pm$ 6.11	39.56 $\pm$ 6.33	48.40 $\pm$ 9.01
Weighted Averaging	17.62 $\pm$ 3.38	37.72 $\pm$ 6.35	62.55 $\pm$ 6.03
Polychotomous Voting	26.32 $\pm$ 2.54	52.28 $\pm$ 3.42	63.23 $\pm$ 1.12
Linear Aggregator	27.03 $\pm$ 6.50	56.94 $\pm$ 6.50	67.64 $\pm$ 4.52
NN Aggregator	44.20 $\pm$ 3.29	<b>68.22</b> $\pm$ 4.19	<b>75.61</b> $\pm$ 1.85
MoE (Gating)	<b>49.86</b> $\pm$ 1.61	62.11 $\pm$ 3.62	69.50 $\pm$ 2.61

Table 10: FENS vs SOTA FL algorithms on the CIFAR-10 dataset.

Algorithm	$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.1$
FEDADAM	39.324 $\pm$ 7.855	<b>68.748</b> $\pm$ 2.762	78.736 $\pm$ 3.552
FEDAVG	37.600 $\pm$ 6.428	60.344 $\pm$ 1.705	77.062 $\pm$ 3.678
FEDNOVA	32.316 $\pm$ 3.844	60.280 $\pm$ 2.300	78.732 $\pm$ 3.252
FEDPROX	37.344 $\pm$ 6.001	59.772 $\pm$ 1.413	75.250 $\pm$ 4.705
FEDYOGI	39.788 $\pm$ 7.726	67.168 $\pm$ 2.793	<b>78.980</b> $\pm$ 3.047
SCAFFOLD	21.824 $\pm$ 3.168	27.308 $\pm$ 8.303	65.892 $\pm$ 23.98
FENS	<b>44.200</b> $\pm$ 3.297	68.220 $\pm$ 4.197	75.613 $\pm$ 1.858

Table 11: FENS vs. iterative FL. Results from Figure 4.

Dataset	$\alpha$	FEDKD 1R	FEDAVG STC	FEDKD 4R	FENS	FEDADAM (Max acc.)
CF-100	0.01	28.98 $\pm$ 4.55	24.83 $\pm$ 1.27	32.22 $\pm$ 2.39	44.46 $\pm$ 0.31	50.25 $\pm$ 0.58
	0.05	39.01 $\pm$ 1.11	29.67 $\pm$ 0.46	42.12 $\pm$ 0.23	49.70 $\pm$ 0.86	56.62 $\pm$ 0.28
	0.1	42.38 $\pm$ 0.78	31.87 $\pm$ 0.58	42.66 $\pm$ 0.89	51.11 $\pm$ 0.37	59.96 $\pm$ 0.80
CF-10	0.01	18.59 $\pm$ 2.92	22.70 $\pm$ 7.08	21.81 $\pm$ 4.27	44.20 $\pm$ 3.29	39.32 $\pm$ 7.85
	0.05	38.84 $\pm$ 6.03	50.56 $\pm$ 3.60	43.26 $\pm$ 7.39	68.22 $\pm$ 4.19	68.74 $\pm$ 2.76
	0.1	64.14 $\pm$ 5.17	61.93 $\pm$ 3.01	62.61 $\pm$ 6.17	75.61 $\pm$ 1.85	78.73 $\pm$ 3.55
SVHN	0.01	23.62 $\pm$ 10.1	32.64 $\pm$ 12.9	25.82 $\pm$ 6.96	57.35 $\pm$ 12.6	66.85 $\pm$ 9.26
	0.05	37.41 $\pm$ 9.62	75.48 $\pm$ 4.39	49.21 $\pm$ 8.17	76.76 $\pm$ 2.98	83.55 $\pm$ 1.52
	0.1	61.38 $\pm$ 3.90	85.32 $\pm$ 2.66	74.33 $\pm$ 3.34	83.64 $\pm$ 0.75	88.04 $\pm$ 0.06

Table 12: Figure 6 results. FENS vs. iterative FL – Fed-Camelyon16 (row 1).

Algorithm	AUC
FEDADAM	0.500 $\pm$ 0.000
FEDAVG	0.569 $\pm$ 0.063
FEDPROX	0.584 $\pm$ 0.069
FEDYOGI	0.513 $\pm$ 0.026
SCAFFOLD	0.485 $\pm$ 0.100
FENS	<b>0.715</b> $\pm$ 0.024

Table 13: Figure 6 results. FENS vs. one-shot FL – Fed-Camelyon16 (row 2).

Algorithm	AUC
Client 0	0.696 $\pm$ 0.018
Client 1	0.683 $\pm$ 0.067
FEDAVG-OS	0.673 $\pm$ 0.007
FEDPROX-OS	0.618 $\pm$ 0.043
FENS	<b>0.715</b> $\pm$ 0.024

Table 14: Figure 6 results. FENS vs. iterative FL – Fed-Heart-Disease (row 1).

Algorithm	Accuracy
FEDADAM	0.781 $\pm$ 0.002
FEDAVG	<b>0.794</b> $\pm$ 0.004
FEDPROX	0.792 $\pm$ 0.009
FEDYOGI	0.785 $\pm$ 0.002
SCAFFOLD	0.792 $\pm$ 0.006
FENS	0.781 $\pm$ 0.027

Table 15: Figure 6 results. FENS vs. one-shot FL – Fed-Heart-Disease (row 2).

Algorithm	Accuracy
Client 0	<b>0.796</b> $\pm$ 0.007
Client 1	0.750 $\pm$ 0.035
Client 2	0.553 $\pm$ 0.059
Client 3	0.599 $\pm$ 0.047
FEDAVG-OS	0.698 $\pm$ 0.048
FEDPROX-OS	0.732 $\pm$ 0.014
FENS	0.781 $\pm$ 0.027

Table 16: Figure 6 results. FENS vs. iterative FL – Fed-ISIC2019 (row 1).

Algorithm	Balanced Accuracy
FEDADAM	0.710 $\pm$ 0.011
FEDAVG	0.731 $\pm$ 0.007
FEDPROX	0.750 $\pm$ 0.006
FEDYOGI	0.731 $\pm$ 0.020
SCAFFOLD	<b>0.746</b> $\pm$ 0.008
FENS	0.579 $\pm$ 0.022

Table 17: Figure 6 results. FENS vs. one-shot FL – Fed-ISIC2019 (row 2).

Algorithm	Balanced Accuracy
Client 0	0.573 $\pm$ 0.007
Client 1	0.315 $\pm$ 0.012
Client 2	0.433 $\pm$ 0.005
Client 3	0.325 $\pm$ 0.014
Client 4	0.181 $\pm$ 0.008
Client 5	0.217 $\pm$ 0.012
FEDAVG-OS	0.424 $\pm$ 0.007
FEDPROX-OS	0.421 $\pm$ 0.017
FENS	<b>0.579</b> $\pm$ 0.022

## E Compute Resources

We use a cluster comprising a mix of 2x Intel Xeon Gold 6240 @ 2.5 GHz of 36 hyper-threaded cores and 2x AMD EPYC 7302 @ 3 GHz of 64 hyper-threaded cores, equipped with 4x NVIDIA Tesla V100 32G and 8x NVIDIA Tesla A100 40G GPU respectively. Training of local models can take up to 2.5 hours in wall-clock time depending on the dataset, while FENS aggregator model training executes in under 30 minutes in wall-clock time. The time required for executing the baselines varies significantly depending on the baseline, with up to 24 hours in wall clock time for Co-Boosting. Across all experiments that are presented in this article, including different seeds and hyperparameter tuning, the total virtual CPU and GPU time is approximately 3500 and 8000 hours respectively.

## F Broader Impact

Federated Learning (FL) has significantly advanced privacy-preserving machine learning, particularly in sensitive domains like healthcare and finance, by facilitating collaborative model training without sharing raw data. The development of FENS, which combines FL’s accuracy with the communication efficiency of One-shot FL (OFL), carries numerous positive implications. By simultaneously reducing communication costs and maintaining high accuracy, FENS enhances the accessibility and practicality of FL, thereby promoting wider adoption, especially in resource-constrained environments. This has the potential to catalyze advancements in privacy-sensitive sectors like healthcare and finance, where FL is extensively utilized.

## NeurIPS Paper Checklist

### 1. Claims

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

Answer: [Yes]

Justification: Our claims in the abstract and the introductions are appropriately scoped and well supported through our extensive empirical assessments spanning multiple datasets, baselines and settings.

Guidelines:

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

### 2. Limitations

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

Answer: [Yes]

Justification: We have discussed the limitations of our work in Section 4.

Guidelines:

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

### 3. Theory Assumptions and Proofs

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

Answer: [NA]

Justification: The paper does not include theoretical results.

Guidelines:

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

#### 4. Experimental Result Reproducibility

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

Answer: [Yes]

Justification: We have included elaborate descriptions regarding the setup and hyperparameters needed to reproduce the paper in Appendix B, complementing our description in Section 3.1.

Guidelines:

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

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We open-source our code at <https://github.com/sacs-epfl/fens>.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

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

Answer: [Yes]

Justification: Yes, we have included all details in Section 3.1 and Appendix B.

Guidelines:

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

## 7. Experiment Statistical Significance

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

Answer: [Yes]

Justification: Each of our experiments is repeated with at least three random seeds. We report standard deviations in all tables and display the 95% confidence intervals in our plots.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).



- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

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

Answer: [Yes]

Justification: We have included the details on compute resources in Appendix E.

Guidelines:

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

## 9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

Answer: [Yes]

Justification: We have reviewed the Code of Ethics and believe our work adheres to these guidelines.

Guidelines:

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

## 10. Broader Impacts

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

Answer: [Yes]

Justification: We have discussed the broad impacts of our work in Appendix F.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

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

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

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

#### 12. Licenses for existing assets

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

Answer: [Yes]

Justification: We affirm that the medical datasets employed in our experiments within the FLamby benchmark were acquired and utilized in strict accordance with their respective licensing agreements and ethical guidelines. We obtained the necessary permissions and approvals from the appropriate authorities and/or institutions responsible for data collection, and we adhered to all relevant ethical standards and regulations. The owners of the original assets used in this paper were properly cited and attributed.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. **New Assets**

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

Answer: [Yes]

Justification: We open-source our code at: <https://github.com/sacs-epfl/fens>. Our repository is well documented with all the instructions to run the code.

Guidelines:

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

### 14. **Crowdsourcing and Research with Human Subjects**

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

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

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

### 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

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

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.