

# FedRN: Exploiting $k$ -Reliable Neighbors Towards Robust Federated Learning

SangMook Kim  
KAIST  
Republic of Korea  
sangmook.kim@kaist.ac.kr

Wonyoung Shin\*  
NAVER Shopping  
Republic of Korea  
wonyoung.shin@navercorp.com

Soohyuk Jang\*  
POSTECH  
Republic of Korea  
jang001@postech.ac.kr

Hwanjun Song<sup>†</sup>  
NAVER Corp.  
Republic of Korea  
hwanjun.song@navercorp.com

Se-Young Yun<sup>†</sup>  
KAIST  
Republic of Korea  
yunseyoung@kaist.ac.kr

## ABSTRACT

Robustness is becoming another important challenge of federated learning in that the data collection process in each client is naturally accompanied by noisy labels. However, it is far more complex and challenging owing to varying levels of data heterogeneity and noise over clients, which exacerbates the client-to-client performance discrepancy. In this work, we propose a robust federated learning method called FedRN, which exploits  $k$ -reliable neighbors with high data expertise or similarity. Our method helps mitigate the gap between low- and high-performance clients by training only with a selected set of clean examples, identified by a collaborative model that is built based on the reliability score over clients. We demonstrate the superiority of FedRN via extensive evaluations on three real-world or synthetic benchmark datasets. Compared with existing robust methods, the results show that FedRN significantly improves the test accuracy in the presence of noisy labels.

## CCS CONCEPTS

• **Computing methodologies** → **Supervised learning by classification**; • **Human-centered computing** → **Ubiquitous and mobile computing**.

## KEYWORDS

Federated Learning; Robust Learning; Label Noise

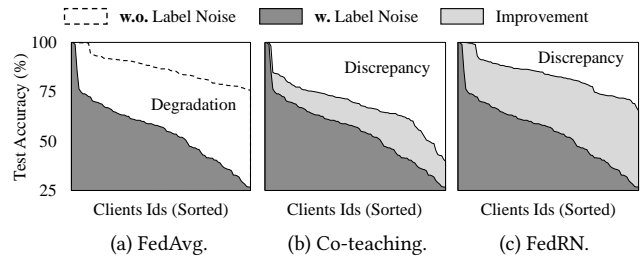
## ACM Reference Format:

SangMook Kim, Wonyoung Shin, Soohyuk Jang, Hwanjun Song, and Se-Young Yun. 2022. FedRN: Exploiting  $k$ -Reliable Neighbors Towards Robust Federated Learning. In *Proceedings of the 31st ACM International Conference on Information and Knowledge Management (CIKM '22)*, October 17–21, 2022, Atlanta, GA, USA. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3511808.3557322>

\*Both authors contributed equally to this research.

<sup>†</sup>Both authors are corresponding authors of this paper.

This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *Proceedings of the 31st ACM International Conference on Information and Knowledge Management (CIKM '22)*, October 17–21, 2022, Atlanta, GA, USA <https://doi.org/10.1145/3511808.3557322>.  
CIKM '22, October 17–21, 2022, Atlanta, GA, USA



**Figure 1: Performance difference of local models over 100 clients after training on CIFAR-10 with symmetric noise of 0–80% in the Non-IID setting [24]: (a) shows the performance degradation incurred by label noise in training data; (b) and (c) contrasts the improvement by Co-teaching and our proposed FedRN. Client Ids are sorted in ascending order by the noise ratio of their local data.**

## 1 INTRODUCTION

Federated learning (FL) is a privacy-preserving distributed learning technique, which handles a large corpus of decentralized data residing on multiple clients such as mobile or IoT devices [5, 34]. The main idea is to learn a joint model by alternating local update and model aggregation phases. Each client performs the local update phase by training the current model with its own training data without sharing any of the clients' raw training data. Then, the model aggregation phase is executed by a server that merges received local models [24]. As the need for on-device machine learning is rising, FL has attracted much attention in various fields [9, 17, 22, 25, 35, 41] and there have been numerous attempts to deploy it to on-device services, such as query suggestion for Google Keyboard [37] and question answering for Amazon Alexa [7].

Since the FL system is deployed over heterogeneous networks, *robustness* has become an important challenge of federated learning, because the data collection process in each client is naturally accompanied by *noisy labels*, which may be corrupted from the ground-truth labels [19]. In the presence of noisy labels, deep neural networks (DNNs) easily overfit to the noisy labels, thereby leading to poor generalization on unseen data [40]. Extensive efforts have been devoted to overcoming noisy labels in the centralized scenario [20, 32], but have yet to be studied widely in the federated learning (decentralized) scenario. Handling noisy labels in the FL



**Figure 2: Overview of FedRN ( $k = 1$ ):** (a)  $k$ -reliable neighbors are retrieved to fit GMMs to the loss values of all training examples in the target client. Their ensembled GMMs are leveraged to identify clean examples, which are used to update the target model in the local update phase; (b) The weights of updated local models are averaged in the model aggregation phase.

scenario is far more challenging than in the centralized scenario, particularly on the following two difficulties:

- (1) **Data Heterogeneity:** Each client has non-independent and identically distributed (Non-IID) data with respect to the number of training examples for each class, which differ between clients.
- (2) **Varying Label Noise:** Noise ratios and types of data flaws vary between clients depending on how their data was collected, mainly attributed to the malfunction of data collectors, crowd-sourcing, and adversarial attacks.

As shown in Figure 1(a), these two difficulties exacerbate the performance discrepancy over clients because local data in each client varies in terms of data heterogeneity and label noise. Such a large performance discrepancy arguably results in poor performance of the global model, because client models are aggregated regardless of whether their models are corrupted by label noise [8].

The *small-loss trick*, which updates a DNN with a certain number of small-loss examples every training iteration [14, 18, 28, 38], could be a prominent direction to provide robustness; small-loss examples are typically regarded as the clean set because DNNs tend to first learn from clean data and then gradually overfit to noisy data [3]. However, as can be seen from Figure 1(b), the performance discrepancy still remains even when all the local models were trained using a popular robust learning method, Co-teaching [14], that uses the small-loss trick for sample selection. Although the overall performance is higher than that of FedAvg, the performance gap between the high- and low-performance model is considerably large. Therefore, it is essential to not only (i) identify clean examples from noisy data, but also (ii) improve the low-performance models to increase the overall robustness for federated learning with noisy labels.

In this paper, we propose a simple yet effective robust approach called **FedRN** (Fed-erated learning with Reliable Neighbors), which exploits  $k$ -reliable neighbors in the client pool to help identify clean examples even when the target client is unreliable. We introduce the notion of  $k$ -reliable neighbors, which are  $k$  clients that have *similar* data class distributions to the target client (data similarity), or *minimal* label noise in their data (data expertise). Using  $k$ -reliable neighbors, FedRN improves the performance of every client’s local model; the low-performance models benefit more from their  $k$ -reliable neighbors, as shown in Figure 1(c).

In the local update phase shown in Figure 2(a), FedRN first retrieves  $k$ -reliable neighbors from the server. “Neighbor 2” has a

relatively clean dataset or a similar class distribution to the target client’s, so the server transmits the model of “Neighbor 2” to the target client. The retrieved  $k$ -reliable models are used to identify clean examples from the client’s dataset collaboratively with the target model; hence, with clear guidance from reliable neighbors, the target model can be improved significantly even if its performance is poor due to heavy label noise. In detail, the target and  $k$ -reliable models fit bi-modal univariate Gaussian mixture models (GMMs) to the loss distribution of the target’s training data, respectively. Next, we build joint mixture models by aggregating all the GMMs based on the reliability scores of neighbor clients. Since the loss distributions of clean and noisy examples are bi-modal [3], the training examples with a higher probability of belonging to the clean (*i.e.*, small-loss) modality are treated as clean examples and used to robustly update the target model. In this process, the low-performance model hardly hinders the accuracy of sample selection due to its low reliability score.

After the local update phase, the weights of updated models are averaged to build a global model in the model aggregation phase in Figure 2(b). These two phases are repeated until the global model converges following the standard federated learning pipeline [24].

Our main contributions are summarized as follows:

- This is a new FL framework that exploits reliable neighbors to tackle the challenge of data heterogeneity and varying label noise.
- We introduce and examine the two indicators of data expertise and similarity, measuring the reliability score of neighbor clients without infringing on data privacy.
- FedRN remarkably improves the performance of underperforming client models by mitigating the client-to-client discrepancy, thereby boosting overall robustness.
- FedRN significantly outperforms state-of-the-art methods on three real-world or synthetic benchmark datasets with varying levels of data heterogeneity and label noise.

## 2 RELATED WORK

**Federated Learning.** Federated learning aims to learn a strong global model by exploiting all of the client data without infringing on data privacy. However, due to the heterogeneity in training data, it typically suffers from performance degradation of the global model [44] or weight divergence [42]. In this regard, extensive efforts have been made to address these issues. Duan et al.

[13] addressed the problem of long-tailed distribution in training data by performing data augmentation. Zhao et al. [42] let all of the clients share the same public data to mitigate the data heterogeneity. Meanwhile, Li et al. [21] and Acar et al. [1] added a regularization term that prevents the local model from diverging to the global model due to its own skewed Non-IID data. Although this family of methods contributes to improving the effectiveness of federated learning, they simply assume that all the labels in training data are clean, which is not a realistic scenario.

**Learning with Noisy Labels.** Extensive studies have been conducted to overcome noisy labels in the centralized scenario. A prominent direction for handling noisy labels is sample selection, which trains DNNs for a possibly clean subset of noisy training data. For example, Co-teaching [14] trains two DNNs where each DNN selects small-loss examples in a mini-batch and then feeds them to its peer network for further training. MORPH [28] divides the training process into two learning periods and employs two different criteria for sample selection. Another possible direction is to modify the loss of training examples based on importance reweighting or label correction [15, 26, 30, 43].

Most recently, to leverage all the training data, sample selection methods have been combined with other approaches, such as loss correction and semi-supervised learning. SELFIE [27] uses relabeled noisy examples in conjunction with the selected clean ones. DivideMix [20] treats selected clean examples as labeled data and applies MixMatch [4] for semi-supervised learning. Despite these methods' success, they do not perform well in the federated learning setting due to the client-to-client performance discrepancy.

**Robust Federated Learning.** To address the noisy labels in FL, most studies assume that trustworthy data exists either on the client- or server-side. Chen et al. [11] estimated the credibility of each client with small clean validation data, and aggregated models based on the estimated credibilities of clients. Tuor et al. [31] proposed a data filtering method, which identifies highly relevant examples for the given specific task using a reference model trained on clean benchmark data. In another direction, a few robust methods using label correction have been proposed [33, 36, 39]. The most representative Robust FL [36] shares the central representations of clients' local data to maintain the consistent decision boundary over clients, and performed global-guided label correction assuming the IID scenario. However, existing methods often rely on unrealistic assumptions including the existence of clean validation data or the IID scenario. The recent label correction approaches including Robust FL still suffer from false label corrections produced by incorrect models under heavy label noise or severe Non-IID scenarios.

**Difference from Existing Work.** We clarify *why FL with noisy labels is problematic*; the client-to-client performance discrepancy hinders the use of existing robust methods in the label noise community into FL. Although a few robust methods have been proposed, this problem has been overlooked. In this paper, we thus bridge the two different topics to improve overall robustness by introducing *k*-reliable neighbors such that they deliver credible information to the clients. Moreover, compared to all aforementioned robust methods, FedRN mainly focuses on identifying reliable neighbor clients for sample selection, which does not require any unrealistic

supervision, such as a clean validation set or knowledge of true noise rates per client.

### 3 PRELIMINARIES

A multi-class classification problem requires training data  $\mathcal{D} = \{(x_i, y_i^*)\}_{i=1}^{|\mathcal{D}|}$ , where  $x$  is a training example with its ground-truth label  $y^*$ . However, noisy training data  $\tilde{\mathcal{D}} = \{(x_i, \tilde{y}_i)\}_{i=1}^{|\mathcal{D}|}$  may have corrupted labels such that  $\tilde{y}_i \neq y_i^*$  for some  $i$ . We herein briefly summarize the conventional learning pipeline for (i) standard federated learning and (ii) robust learning with sample selection.

- **Federated Learning (FedAvg):** It is assumed that each client  $i$  has its own Non-IID dataset  $\mathcal{D}_i$  (i.e.,  $\mathcal{D} = \cup_i \{\mathcal{D}_i\}$ ), and cannot directly access another client's data. Therefore, a global model  $\Theta_{\text{global}}$  is trained by alternating the local update and model aggregation phase.

During the local update phase, a fixed number of target clients  $\mathcal{M}$  are selected from the entire client pool. For each target client  $c \in \mathcal{M}$ , all the other clients are defined as its neighbors  $n \in \mathcal{N}_c$ . Next, the target client's local model  $\Theta_c$  is trained on its own data  $\mathcal{D}_c$  for certain local epochs with the standard stochastic gradient descent [24]. All the trained models for  $\mathcal{M}$  are then received and averaged by the server in the model aggregation phase,

$$\Theta_{\text{global}} \leftarrow \sum_{c \in \mathcal{M}} w_c \Theta_c, \text{ where } w_c = \frac{|\mathcal{D}_c|}{\sum_{c' \in \mathcal{M}} |\mathcal{D}_{c'}|} \quad (1)$$

where  $w_c$  is the aggregation weight that is based on the data size of each client. Lastly, the updated global model is broadcasted to the clients. This training round is called a transmission round and is repeated until convergence.

- **Robust Learning with Sample Selection:** Many recent robust approaches have adopted sample selection, which treats a certain number of small-loss examples as the clean set [14, 20, 28, 38]. In the centralized scenario, given a noisy mini-batch  $\mathcal{B}(c, \tilde{\mathcal{D}})$ , the clean subset  $\mathcal{S}$  is identified and used to directly update the global model  $\Theta_{\text{global}}$ ,

$$\Theta_{\text{global}} \leftarrow \Theta_{\text{global}} - \eta \nabla \frac{1}{|\mathcal{S}|} \sum_{x \in \mathcal{S}} \ell(x, \tilde{y}; \Theta_{\text{global}}), \quad (2)$$

where  $\eta$  is the learning rate and  $\ell$  is a certain loss function. The remaining examples (mostly large-loss examples) in the mini-batch are excluded to pursue robust learning.

Note that training data is typically assumed to be *clean* in federated learning, while a global model is assumed to have access to *all* of the training data in robust learning. However, our objective is to robustly train a global model  $\Theta_{\text{global}}$  even when noisy labels exist in the federated learning scenario.

## 4 PROPOSED METHOD: FEDRN

We introduce the notion of *k*-reliable neighbors and describe how to leverage them for robust federated learning with noisy labels.

### 4.1 Main Concept: k-Reliable Neighbors

The key idea of FedRN is to exploit useful information in other clients without violating privacy. Hence, the main challenge is

to find the most helpful neighbor clients using limited information about each client. Intuitively speaking, these helpful neighbor clients have (1) high expertise on their data, or (2) data distributions similar to the target client's:

- **Data Expertise:** The training data should be sufficiently *clean*. The neighbor client with heavy noise could hinder identifying clean examples from the client's noisy data.
- **Data Similarity:** Data and class distributions should be *similar* to the target client's. The distribution shift problem could lead to model incompatibility between target and neighbor clients.

We investigate how these two considerations affect selecting reliable neighbor clients for robust federated learning. Given a target client  $c$ , a reliable neighbor client is defined to be one of the  $k$ -reliable neighbors  $\mathcal{R}_c(k)$ , as described in Definition 4.1. Satisfying the condition,  $k$ -reliable neighbors can deliver clean and useful information even when the training data of the target client is unreliable due to the heavy label noise.

*Definition 4.1.* Let  $\text{Exp}(\cdot)$  and  $\text{Sim}(\cdot, \cdot)$  be the score function for data expertise of a client and data similarity between two clients, respectively. Given a target client  $c$  and available neighbor clients  $\mathcal{N}_c = \{n_i\}_{i=1}^{|\mathcal{N}_c|}$ , the  $k$ -reliable neighbors  $\mathcal{R}_c(k)$  is a subset of  $\mathcal{N}_c$  with size  $k$  satisfying,

$$\begin{aligned} \mathcal{R}_c(k) &= \text{argmax}_{|\mathcal{R}' : n \in \mathcal{R}'| = k} R(c, n), \text{ where} \\ R(c, n) &= \alpha \cdot \text{Exp}(n) + (1 - \alpha) \cdot \text{Sim}(c, n) \end{aligned} \quad (3)$$

and  $\alpha$  is the balancing term that determines the contribution of each score. The data similarity to itself is 1, i.e.,  $\text{Sim}(c, c) = 1$ .  $\square$

To compute the reliability score  $R(c, n)$  in Eq. (3), we propose two indicators of data expertise and similarity using client's local models. Since it is infeasible to access the raw data in neighbor clients, we leverage the server's copies of local models, which are sent to the server during the model aggregation phase. By doing so, our approach does not compromise the privacy-preserving property in federated learning.

With the intention of selecting clients with cleaner examples, we define the data expertise score motivated by the memorization effect of DNNs. In general, DNNs memorize clean examples faster than noise examples during training, which has been observed consistently even at extreme noise ratios [28]. In this sense, the early memorization of clean examples leads to higher training accuracy if the client has cleaner local data. Consequently, during the model aggregation phase, the server receives the training accuracy of each client along with their local models, then computes and normalizes the data expertise  $\text{Exp}(\cdot)$  per client,

$$\text{Exp}(c) = \frac{\text{Acc}(c) - \min \text{Acc}(\{c\} \cup \mathcal{N}_c)}{\max \text{Acc}(\{c\} \cup \mathcal{N}_c) - \min \text{Acc}(\{c\} \cup \mathcal{N}_c)}, \quad (4)$$

where  $\text{Acc}$  denotes the training accuracy of the local model for the client  $c$ ;  $\max \text{Acc}$  and  $\min \text{Acc}$  are the maximum and minimum training accuracy of the given client set.

Regarding the data similarity score, we propose to use *predictive difference*, which provides a hint to the similarity between two heterogeneous data distributions without the use of raw data. Hence, the data similarity between two clients is approximated by the prediction difference for the *same* random input, where a smaller difference indicates a higher similarity. Accordingly, after the local

update phase, each client generates the same input  $\tilde{x}$  of Gaussian random noise [16] and transmits its softmax outputs  $p(\tilde{x}; \Theta_n)$  into the server, where  $p(x; \Theta_n)$  denotes the softmax output of the client  $n$  for input  $x$ . The data similarity  $\text{Sim}(\cdot, \cdot)$  between the target client  $c$  and neighbor client  $n$  is then computed by the server,

$$\text{Sim}(c, n) = \text{Cosine}(p(\tilde{x}; \Theta_c), p(\tilde{x}; \Theta_n)). \quad (5)$$

We use the cosine similarity denoted by  $\text{Cosine}$  to measure the predictive difference, because the *symmetry* of similarity is guaranteed, i.e.,  $\text{Sim}(c, n) = \text{Sim}(n, c)$ . The min-max normalization is also applied similarly to data expertise. We provide an in-depth analysis of how the two proposed indicators work precisely in Section 5.4.

## 4.2 Robust Learning with $k$ -Reliable Neighbors

For every communication round, each client participating in the upcoming round receives its  $k$ -reliable neighbors from the server. When updating the target client's model, we follow the general philosophy of *sample selection* for robust learning by identifying clean examples from the noisy training data.

Different to robust methods in centralized scenarios,  $k + 1$  models, comprising  $k$ -reliable neighbors and the target client's model, cooperate to mitigate the performance discrepancy among clients and obtain more reliable results. For each of the  $k + 1$  models, we fit two-component GMMs to model the loss distributions of clean and noisy examples in view of loss distributions being *bi-modal* [2, 20]. Based on the ensembled results of these GMMs, the target client's model is only updated with examples selected as clean.

To be specific, our aim is to estimate the probability of each training example being clean with  $k + 1$  models with different reliability scores. At the beginning of the local update phase, the GMM is fitted to the loss of all available training examples in the target client by using the Expectation-Maximization (EM) algorithm. Given a noisy example  $x$ , the probability of being clean is obtained through its posterior probability for the clean (small-loss) modality,

$$p(g|\ell(x, \tilde{y}; \Theta)) = p(g)p(\ell(x, \tilde{y}; \Theta)|g)/p(\ell(x, \tilde{y}; \Theta)), \quad (6)$$

where  $g$  denotes the Gaussian modality with a smaller mean (i.e., smaller loss) and  $\ell(\cdot)$  is the loss function. Subsequently, the probability of each example to be clean is ensembled over  $k + 1$  models with their reliability scores  $R(\cdot, \cdot)$  in Eq. (3). Given a target client  $c$ , the clean probability of joint mixtures can be formulated as:

$$\begin{aligned} p(\text{clean}|x; \mathcal{R}_c(k)) &= \sum_{n \in \{c\} \cup \mathcal{R}_c(k)} R'(c, n) \times p(g|\ell(x, \tilde{y}; \Theta_n)), \\ \text{where } R'(c, n) &= R(c, n) / \sum_{n' \in \{c\} \cup \mathcal{R}_c(k)} R(c, n'). \end{aligned} \quad (7)$$

The ensemble with the reliability scores give higher weight to the GMMs of trusted neighbors with high data expertise or similarity.

Lastly, FedRN constructs a clean set  $\mathcal{S}_c$  such that each training example in the set has a clean probability greater than 0.5,

$$\mathcal{S}_c = \{x \in \tilde{\mathcal{D}}_c : p(\text{clean}|x; \mathcal{R}_c(k)) > 0.5\}, \quad (8)$$

where  $\tilde{\mathcal{D}}_c$  is the noisy data of the target client  $c$ . The model is finally updated with only the clean set for a specified number of local epochs in the local update phase, while the complement of the set, which are likely to be noisy, are discarded for robust learning.

**Algorithm 1** FedRN Algorithm

---

INPUT:  $\tilde{\mathcal{D}}_c$ : noisy data of client  $c$ , *rounds*: # total rounds,  
 $e$ : # local epochs,  $k$ : # reliable neighbors  
OUTPUT:  $\Theta_{\text{global}}^t$ : a global model

- 1:  $t \leftarrow 1$ ,  $\Theta_{\text{global}}^t \leftarrow$  Initialize network parameters
- 2: **for** *round* = 1 **to** *rounds* **do**
- 3:    $\mathcal{M} \leftarrow$  Randomly select  $m$  clients ( $m > k$ )
- 4:   /\* Find  $k$ -reliable neighbors in Definition 4.1 \*/
- 5:   Compute  $\forall_{c \in \mathcal{M}} \forall_{n \in \mathcal{N}_c} \text{Sim}(c, n)$  by Eq. (5)
- 6:   Retrieve  $\mathcal{R}_c(k)$  for  $c \in \mathcal{M}$
- 7:   Broadcast  $\Theta_{\text{global}}^t$  and  $\mathcal{R}_c(k)$  to clients  $\in \mathcal{M}$
- 8:   /\* I. Local Update Phase \*/
- 9:   **for**  $c \in \mathcal{M}$  **in parallel do**
- 10:      $\mathcal{R}_c(k) \leftarrow$  Receive  $k$ -reliable neighbors
- 11:      $\mathcal{R}_c(k) \leftarrow$  Fine-tune for  $\mathcal{S}_c^{\text{aux}}$  in Eq. (9).
- 12:     /\* Sample selection with  $\mathcal{R}_c(k)$  by Eq. (8) \*/
- 13:      $\mathcal{S}_c \leftarrow \{x \in \tilde{\mathcal{D}}_c : p(\text{clean}|x; \mathcal{R}_c(k)) > 0.5\}$
- 14:      $\Theta_c \leftarrow \Theta_{\text{global}}^t$  /\* Initialization \*/
- 15:      $\Theta_c \leftarrow$  Train for  $e$  epochs with  $\mathcal{S}_c$
- 16:      $\text{Acc}(c) \leftarrow$  Compute data expertise by Eq. (4)
- 17:      $p(\tilde{x}; \Theta_c) \leftarrow$  Compute the softmax output for  $\tilde{x}$
- 18:     Send  $\Theta_c$ ,  $\text{Acc}(c)$ , and  $p(\tilde{x}; \Theta_c)$  to the server
- 19:     /\* II. Model Aggregation Phase \*/
- 20:      $\Theta_{\text{global}}^{t+1} \leftarrow \sum_{c \in \mathcal{M}} w_c \Theta_c$  s.t.  $w_c = \frac{|\tilde{\mathcal{D}}_c|}{\sum_{c' \in \mathcal{M}} |\tilde{\mathcal{D}}_{c'}|}$
- 21:      $t \leftarrow t + 1$
- 22: **return**  $\Theta_{\text{global}}^t$

---

This collaborated approach with  $k$ -reliable neighbors considerably increases the robustness of *low*-performance clients with the aid of neighbors with high data expertise or similarity. Therefore, FedRN reduces the performance gap between clients and enhances the overall robustness, thus leading to a stronger global model.

### 4.3 Fine-tuning with k-Reliable Neighbors

Due to the data heterogeneity in federated learning, the problem of distribution shift between clients' local models is another challenge when ensembling their predictions for sample selection. Hence, FedRN is integrated with the fine-tuning technique that helps quickly adapt to the local data distribution [12].

Before constructing the clean set in Eq. (8), we identify an auxiliary set of clean examples using *only* the target client's model,

$$\mathcal{S}_c^{\text{aux}} = \{x \in \tilde{\mathcal{D}}_c : p(\text{clean}|x; \{c\}) > 0.5\} \quad (9)$$

and fine-tune all the retrieved  $k$ -reliable neighbors for the auxiliary set. The auxiliary set could, however, be very noisy especially when the target client has heavy noise in its training data. To prevent feature extractors from being corrupted by such client-side noise, we only fine-tune the classification head. In Section 5.5, we examine the effect of the fine-tuning technique on data expertise and similarity used in neighbor search, and provide insights on their use for robust federated learning with noisy labels.

	# of Train	# of Val.	# of Classes	Noise Ratio
CIFAR-10	50,000	10,000	10	$\approx 0\%$
CIFAR-100	50,000	10,000	100	$\approx 0\%$
mini-WebVision	65,944	2500	50	$\approx 20\%$

**Table 1: Summary of datasets.**

### 4.4 Algorithm Pseudocode

Algorithm 1 describes the overall procedure of FedRN. We perform standard federated learning, FedAvg, for warm-up epochs before applying FedRN, because robust training methods are generally applied after the warm-up phase in the literature [14, 27]. When training with FedRN, retrieved neighbor models are fine-tuned for 1 epoch before sample selection (Line 11). Next, the target client identifies clean examples by leveraging  $k$ -reliable neighbor models (Line 13). The received global model is updated with estimated clean samples (Lines 14–15). In the aggregation phase, the central server aggregates the updated local models (Line 20).

The main additional cost of FedRN is the communication overhead from sending  $k$ -reliable models from the server to the selected client (see Appendix 5.9 for detailed cost analysis). However, a small number of reliable neighbors is sufficient (see Section 5.6) and the cost can be drastically reduced by over a factor of 32 times in modern federated learning by sending the difference of parameters [16] or a compressed model [23]. Furthermore, there is no issues in model convergence because we applied a sample selection approach, which don't has convergence problems, in a decentralized setting using FedAvg [24], which has proven model convergence.

## 5 EVALUATION

Our evaluation was conducted to support the following:

- FedRN is **more robust** than five state-of-the-art methods for federated learning with noisy labels (Section 5.2).
- FedRN consistently identifies clean examples from noisy data with **high precision** and **recall** (Section 5.3).
- The reliability metric is **effective** in finding  $k$ -reliable neighbors for robust learning (Section 5.4).
- The use of fine-tuning is **necessary** owing to the distribution shift in local data (Section 5.5).
- A small number ( $k \leq 2$ ) of reliable neighbors is **sufficient** to achieve high robustness (Section 5.6).

### 5.1 Experiment Configuration

We verified the superiority of FedRN compared with five robust learning methods on three real-world or synthetic benchmark datasets in a Non-IID federated learning setting, where the Non-IID setting is more realistic and difficult than the IID setting [6, 42, 44].

**Datasets.** We performed image classification on synthetic and real-world benchmark data: CIFAR-10 and CIFAR-100; mini-WebVision, a subset of real-world noisy data consisting of large-scale web images. We used the first 50 classes in WebVision V1 following the literature [20]. The statistics of them are summarized in Table 1.

For the federated learning setup with noisy labels, we merged the standard setting in the two research communities:

- **Non-IID Data:** Two popular ways of client data partitioning were applied [44]; (i) Sharding: training data is sorted by its class labels

Non-IID Type	Shard ( $S = 2$ )				Shard ( $S = 5$ )				Dirichlet ( $\beta = 0.5$ )			
	Symmetric		Asym	Mixed	Symmetric		Asym	Mixed	Symmetric		Asym	Mixed
Noise Type	0.0–0.4	0.0–0.8	0.0–0.4	0.0–0.4	0.0–0.4	0.0–0.8	0.0–0.4	0.0–0.4	0.0–0.4	0.0–0.8	0.0–0.4	0.0–0.4
Oracle <sup>1</sup>	69.10	67.19	69.10	68.89	78.00	75.94	78.00	77.7	81.79	80.38	81.79	81.39
FedAvg [24]	46.39	38.69	49.09	46.82	66.44	56.46	67.23	67.95	75.92	70.60	77.77	77.17
Co-teaching [14]	63.20	52.72	62.48	61.68	74.66	67.17	74.18	75.18	79.97	75.39	79.53	79.94
Joint-optm [30]	50.44	42.23	38.04	47.28	69.22	64.02	67.29	68.84	75.46	70.43	74.81	75.43
SELFIE [27]	62.64	53.69	64.21	62.63	74.57	66.90	74.77	74.44	78.57	72.92	78.58	79.02
DivideMix [20]	62.35	58.18	62.07	63.38	68.73	65.82	68.95	69.32	74.26	72.25	73.29	73.47
Robust FL [36]	56.25	45.59	55.52	57.58	70.30	62.89	69.04	70.02	75.75	70.63	74.00	75.49
<b>FedRN (<math>k=1</math>)</b>	67.33	60.33	67.51	67.92	76.37	72.30	76.74	76.92	79.99	75.92	80.05	79.79
<b>FedRN (<math>k=2</math>)</b>	<b>67.62</b>	<b>62.94</b>	<b>68.33</b>	<b>68.11</b>	<b>76.81</b>	<b>72.85</b>	<b>77.33</b>	<b>76.99</b>	<b>80.34</b>	<b>76.49</b>	<b>80.38</b>	<b>80.28</b>

**Table 2: Test accuracy (%) on CIFAR-10 with symmetric, asymmetric (Asym), and mixed noise (Mixed).**

Non-IID Type	Shard ( $S = 20$ )			
	Symmetric		Asym	Mixed
Noise Type	0.0–0.4	0.0–0.8	0.0–0.4	0.0–0.4
Oracle	47.83	44.82	47.83	47.33
FedAvg [24]	33.36	26.15	34.90	34.02
Co-teaching [14]	43.43	37.05	42.90	43.83
Joint-optm [30]	35.89	30.92	33.15	35.29
SELFIE [27]	44.14	37.65	43.25	43.78
DivideMix [20]	46.21	42.83	45.59	46.94
Robust FL [36]	35.33	29.20	33.08	34.27
<b>FedRN (<math>k=1</math>)</b>	47.30	42.27	<b>47.62</b>	46.62
<b>FedRN (<math>k=2</math>)</b>	<b>47.46</b>	<b>43.07</b>	47.52	<b>47.17</b>

**Table 3: Test accuracy (%) on CIFAR-100.**

and divided into “ $S \times (\# \text{ clients})$ ” shards, which are randomly assigned to each client with an equal number of  $S$ . **(ii)** Dirichlet Distribution: each client is assigned with training data drawn from the Dirichlet distribution with a concentration parameter  $\beta$ . For the experiments,  $S$  was set to 2, 5, 10, or 20, and  $\beta$  was set to 0.5.

- **Label Noise:** We injected artificial noise by following typical protocols [29]; **(i)** Symmetric Noise: a true label is flipped into all possible labels with equal probability. **(ii)** Asymmetric Noise: a true label is flipped into only a certain label. Note that we injected **varying levels** of label noise into the clients, so each client had different noise rates. For instance, a noise of 0.0–0.8 means that the noise rate increases linearly from 0.0 to 0.8 as the client’s index increases; thus, the average noise rate of clients’ training data is 0.4. We further tested for **mixed noise** where the two different types of label noise are injected into the clients. That is, symmetric noise is applied to half of the clients, while asymmetric noise is applied to the remaining clients.

We did not inject any label noise into mini-WebVision because it contains real label noise whose rate is estimated at 20.0% [29].

**Implementation Details.** We compared FedRN with a standard federated learning method, FedAvg [24], and five state-of-the-art robust methods for handling noisy labels, namely Co-teaching [14], Joint-optimization [30], SELFIE [27], DivideMix [20], and Robust FL [36]. All the robust methods were combined with FedAvg to support decentralized learning. The total number of clients was set to 100 with a fixed participation rate of 0.1. The training (transmission) round was set to be 500 for CIFAR-10 and 1000 for CIFAR-100 and

Shard ( $S = 10$ ) Method	mini-WebVision	
	Top-1 Accuracy	Top-5 Accuracy
FedAvg [24]	20.80	45.00
Co-teaching [14]	21.76 (+0.96)	46.64 (+1.64)
Joint-optm [30]	13.56 (–7.24)	35.56 (–9.44)
SELFIE [27]	22.20 (+1.40)	48.76 (+3.76)
DivideMix [20]	20.84 (+0.04)	45.72 (+0.72)
Robust FL [36]	14.12 (–6.68)	35.24 (–9.76)
<b>FedRN (<math>k=1</math>)</b>	22.52 (+1.72)	48.48 (+3.48)
<b>FedRN (<math>k=2</math>)</b>	<b>22.76 (+1.96)</b>	<b>49.16 (+4.16)</b>

**Table 4: Validation accuracy (%) on mini-WebVision. The values in parentheses are the improvement over FedAvg.**

WebVision, ensuring all the models’ training convergence. The SGD optimizer was used with a learning rate of 0.01 and a momentum of 0.5. The number of local epochs was set to be 5. We used the model for learning consists of 4 convolution layers and 1 linear classifier for CIFAR-10 and MobileNet for CIFAR-100 and mini-WebVision.

As for hyperparameters related to our method, FedRN requires the balancing term  $\alpha$  and number of reliable neighbors  $k$ . We used  $\alpha = 0.6$ , which was obtained via a grid search (see the section 5.7). We used  $k = 1$  and  $k = 2$  since the performance gain of FedRN was consistently high as long as  $k \geq 1$ , which is detailed in Section 5.6.

The hyperparameters of the baseline robust methods were configured favorably in line with the literature. We clarify the hyperparameter setup of all baseline algorithms, as follows:

- Joint Optimization [30] requires the two coefficients for its two regularization losses for robust learning. These two coefficients were set to be 1.2 and 0.8, respectively.
- SELFIE [27] requires the uncertainty threshold and the length of label history as its hyperparameters. As suggested by the authors, they were set to be 0.05 and 15, respectively.
- DivideMix [20] applies MixMatch for semi-supervised learning, which requires the sharpening temperature, the unsupervised loss coefficient, and the Beta distribution parameter. They were set to be 0.5, 25, and 4, respectively.
- Robust FL [36] requires two coefficients for its two additional regularization losses similar to Joint Optimization. They were set to be 1.0 and 0.8, respectively.

<sup>1</sup>Oracle trains the network using FedAvg with all the true-labeled examples, i.e.,  $(1 - \text{noise rate}) \times 100\%$  of all the examples.



All of the algorithms were implemented using PyTorch 1.7.0 and executed using four NVIDIA RTX 2080Ti GPUs. In support of reliable evaluation, we repeated every task thrice and reported the average test (or validation) accuracy, which is the common measure of robustness to noisy labels [10, 14].

## 5.2 Robustness Comparison

**5.2.1 Results with Synthetic Noise.** Tables 2 and 3 show the test accuracy of the global model trained by eight methods for three Non-IID federated learning scenarios with four different noise settings, including symmetric, asymmetric, and mixed label noise. Overall, FedRN achieves the highest test accuracy in every case. Compared with FedAvg, FedRN’s accuracy improves relatively more as the noise rate increases from 0.0–0.4 to 0.0–0.8 in symmetric noise, because the performance discrepancy over clients drastically degrades with the increase in noise rate. In contrast, the other robust methods show considerably poor performance compared with FedRN, which is presumably attributed to the client-to-client performance discrepancy, despite the fact that they are incorporated with FedAvg. Meanwhile, the recent robust FL method, Robust FL, shows relatively poor performance in our setup since the method produces many false label corrections due to the clients with heavy label noise. Unlike, FedRN is robust even to heavy noise because all unreliable labels are excluded from training for high safety without correction. Our method can be extended with semi-supervised learning for exploiting even unreliable examples. We leave this as future work.

**5.2.2 Results with Real-world Noise.** Table 4 displays the validation accuracy of seven different methods on a real-world noisy mini-WebVision dataset. We report the top-1 and top-5 classification accuracy on the mini-WebVision validation set. FedRN maintains its performance dominance over multiple state-of-the-art robust methods for *real-world* label noise as well. FedRN ( $k = 2$ ) improves the top-1 accuracy by up to 9.2 compared with other robust methods.

## 5.3 In-depth Analysis on Selected Examples

All methods except FedAvg and Joint Optimization follow the pipeline of learning with sample selection. Hence, we evaluate the sample selection performance of them using the *two* metrics, namely label precision (LP) and recall (LR), which respectively represent the quality and quantity of examples selected as clean [14, 28],

$$\text{LP} = \frac{|\{(x, \tilde{y}) \in \mathcal{S}_c : \tilde{y} = y^*\}|}{|\mathcal{S}_c|}, \quad \text{LR} = \frac{|\{(x, \tilde{y}) \in \mathcal{S}_c : \tilde{y} = y^*\}|}{|\{(x, \tilde{y}) \in \mathcal{D}_c : \tilde{y} = y^*\}|},$$

where  $\mathcal{D}_c$  and  $\mathcal{S}_c$  denote the training data and its selected clean set of the client  $c$ , respectively.

Figure 3 shows the LP and LR curves averaged across all the clients per epoch on CIFAR-10 with symmetric and asymmetric label noise. The results of Robust FL are omitted due to its too low label precision and recall under the Non-IID scenario. After the warm-up period of 100 transmission rounds, FedRN shows the highest LP and LR with a large improvement of up to 0.12 and 0.11, respectively. In addition, its dominance over other robust learning methods is consistent regardless of the noise type. Therefore, these results indicate that the superior robustness of FedRN in Tables 2–4 is due to its high LP and LR.

$\alpha$	w.o. Fine-tuning		w. Fine-tuning	
	0.0–0.4	0.0–0.8	0.0–0.4	0.0–0.8
0.0	<b>65.03</b>	<b>58.20</b>	67.11	61.81
0.2	63.77	57.24	67.45	62.69
0.4	62.45	53.89	67.03	62.94
0.6	58.74	51.36	<b>67.62</b>	62.94
0.8	56.60	43.61	67.19	<b>62.97</b>
1.0	56.06	44.87	66.73	61.61

**Table 5: Effect of fine-tuning when trained on CIFAR-10 with symmetric noise and shard ( $S = 2$ ) settings.**

## 5.4 In-depth Analysis on Reliability Metric

We justify the use of two indicators for data expertise and similarity in Eqs. (4) and (5).

- **Data Expertise:** Figure 4 shows the training convergence rate (training accuracy) of all clients sorted in ascending order by their noise rate. Owing to the memorization effect of DNNs, the training accuracy of the local model indeed has a strong correlation of over 0.93 with respect to varying levels of label noise. These results confirm that the client’s training accuracy represents its data expertise.
- **Data Similarity:** Figure 5 shows similarity matrices between client’s softmax outputs for a Gaussian random noise when trained on data with or without label noise. As clients with similar Ids (indices) have similar data distributions in our Non-IID setting, the similarity around diagonal entries is high, as can be seen from Figure 5(a). Even with label noise shown in Figure 5(b), it turns out that the similarity trend remains. Therefore, it can be concluded that the similarity between the softmax output is a robust metric to find clients with similar data distributions.

These empirical studies provide empirical evidence for the use of our two indicators,  $\text{Exp}(\cdot)$  and  $\text{Sim}(\cdot, \cdot)$ .

## 5.5 Reliable Neighbors with Fine-tuning

We analyze the effect of fine-tuning on data expertise and similarity used in neighbor search. Table 5 summarizes the performance of FedRN ( $k = 2$ ) with and without fine-tuning, where a small  $\alpha$  indicates that data similarity is considered more than data expertise. When fine-tuning is deactivated, the benefit of retrieving neighbors with high data similarity is remarkably dominant, because lower data similarity implies a larger distribution shift between the target and neighbor models. On the other hand, an opposite trend is observed when fine-tuning is activated. The test accuracy is high when data expertise is properly considered in conjunction with data similarity. Overall, fine-tuning does not only help mitigate the distribution shift problem but also significantly improves performance of FedRN.

## 5.6 Ablation Study on k-Reliable Neighbors

A larger number of reliable neighbors could increase the communication overhead in federated learning. We, therefore, investigated the change in performance of the global model according to the number of reliable neighbors. Tables 6 and 7 show the classification accuracy on CIFAR-10 and CIFAR-100 respectively with varying number of neighbors. The results show that FedRN with  $k = 1$

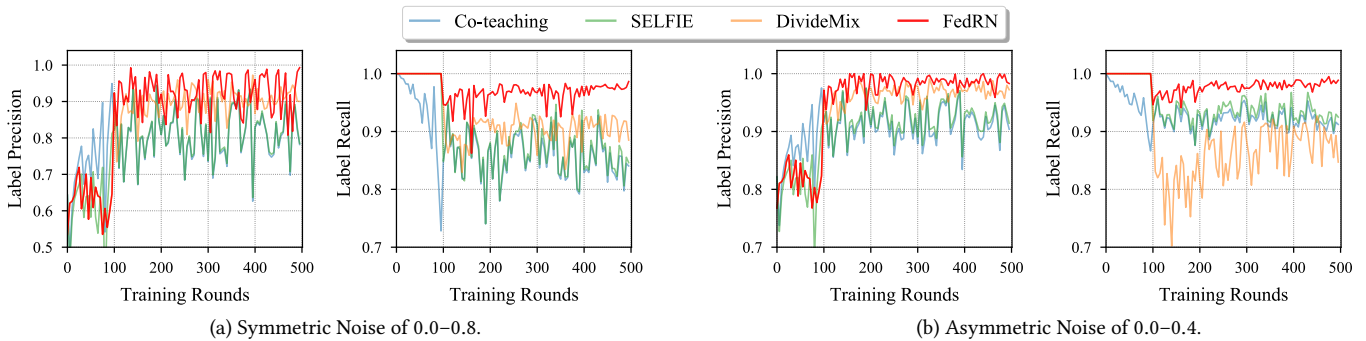


Figure 3: Label precision and label recall curves on CIFAR-10 with the shard ( $S = 2$ ) setting.

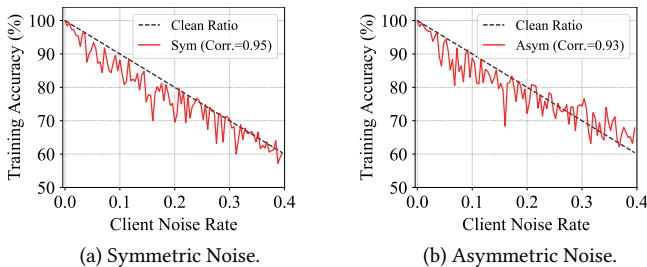


Figure 4: Correlation between the training accuracy and noise rate on CIFAR-10 with shard ( $S = 2$ ) of symmetric and asymmetric noises 0.0–0.4.

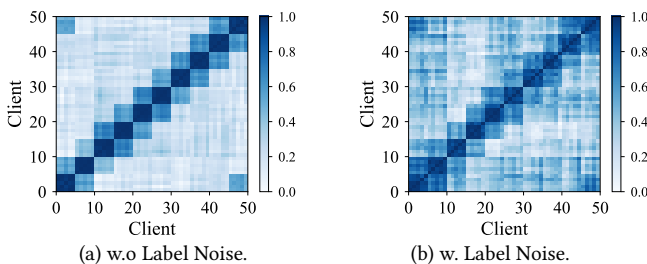


Figure 5: Similarity matrix between the client’s softmax output for a Gaussian random noise when trained on CIFAR-10 with shard ( $S = 2$ ) without noise (a) and with symmetric noise 0.0–0.4 (b).

$k$	Shard ( $S = 2$ )		Dirichlet ( $\beta = 0.5$ )	
	0.0–0.4	0.0–0.8	0.0–0.4	0.0–0.8
1	67.33	60.33	79.99	75.92
2	67.62	62.94	80.34	76.49
3	68.10	62.77	80.31	76.01
4	67.38	62.75	80.36	76.16
5	67.60	63.10	80.18	76.00

Table 6: Test accuracy (%) on CIFAR-10 using FedRN with different number of reliable neighbors.

and  $k = 2$  generally shows comparable performance than  $k > 2$ ; there is little improvement in performance with the increase in the number of reliable neighbors as long as  $k > 1$ . In other words, using only one or two reliable neighbors is sufficient to improve the robustness, because the proposed reliability score successfully retrieves the most helpful neighbor client among all the neighbors.

$k$	Shard ( $S = 20$ )			
	Symmetric		Asym	Mixed
	0.0–0.4	0.0–0.8	0.0–0.4	0.0–0.4
1	47.30	42.27	47.62	46.62
2	47.46	43.07	47.52	47.17
3	47.65	42.33	48.18	47.01
4	47.66	42.48	47.55	47.72
5	48.10	42.58	47.43	47.40

Table 7: Test accuracy (%) of on CIFAR-100 using FedRN with different number of reliable neighbors.

Therefore, in a scenario where the communication burden is significant, practitioners may reduce the communication overhead by only using a single reliable neighbor with FedRN to handle the problem of noisy labels.

### 5.7 Grid Search for Balancing Term

We adjusted  $\alpha$  to determine the contribution of data expertise and similarity in neighbor search. To ascertain the optimal value, we repeated the experiment on CIFAR-10 with various  $\alpha$  values, as summarized in Table 10. Because the fine-tuning technique is incorporated in FedRN, merging data expertise and similarity provides a synergistic enhancement in general. In particular, when  $\alpha = 0.6$ , the performance gain is the highest on average. Based on these results, we set  $\alpha$  to 0.6 for all experiments in Section 5.

### 5.8 Comparison with Random Neighbors

We analyzed the usefulness of  $k$ -reliable neighbors compared with  $k$ -random neighbors. For  $k$ -random neighbors, we replaced Lines 6 and 10 of Algorithm 1 with random client selection. Table 9 summarizes the test accuracy of FedRN with  $k$ -reliable and  $k$ -random neighbors. Since our fine-tuning technique was incorporated with  $k$ -random neighbors, using  $k$ -random neighbors improves the robustness against noisy labels, but FedRN with  $k$ -reliable neighbors maintains its dominance in every case. In detail, if the neighbors with low data expertise are selected, they misclassify mislabeled examples as clean examples in sample selection. Likewise, the neighbors with significantly different data distributions harm identifying clean examples from the target client’s noisy data due to the distribution shift problem. It can be confirmed experimentally that this issue becomes more prominent when the noise rate is high. In conclusion, we should leverage the reliability of neighbors to obtain a global model robust to noisy labels.



	Communication			Computation		
	Server→Client	Client→Server	Total Cost	Forward	Backward	Total Cost
FedAvg [24]	$M$	$M$	$2M$	$eF$	$eB$	$eF + eB$
Co-Teaching [14]	$2M$	$2M$	$4M$	$2eF$	$2eB$	$2eF + 2eB$
Joint-optm [30]	$M$	$M$	$2M$	$eF$	$eB$	$eF + eB$
SELFIE [27]	$M$	$M$	$2M$	$eF$	$eB$	$2eF + eB$
DivideMix [20]	$2M$	$2M$	$4M$	$(4m + 2)eF$	$2eB$	$(4m + 2)eF + 2eB$
Robust FL [36]	$M$	$M$	$2M$	$eF$	$eB$	$eF + eB$
FedRN	$(k + 1)M$	$M$	$(k + 2)M$	$(e + 2k + 1)F$	$(k + e)B$	$(e + 2k + 1)F + (k + e)B$

**Table 8: Analysis of the communication and computation costs in federated learning setting:  $M$  is the communication cost to send the model;  $F$  and  $B$  are the computational costs of forward and backward propagation, respectively;  $e$  is the number of local epochs for each communication round;  $k$  is the number of reliable neighbors.**

Type	Symmetric Noise					
	Shard ( $S = 2$ )		Shard ( $S = 5$ )		Dirichlet	
	0.0–0.4	0.0–0.8	0.0–0.4	0.0–0.8	0.0–0.4	0.0–0.8
$k$ -random	66.82	60.11	76.37	71.56	79.79	75.11
$k$ -reliable	<b>67.62</b>	<b>62.94</b>	<b>76.81</b>	<b>72.85</b>	<b>80.34</b>	<b>76.49</b>

**Table 9: Performance comparison with  $k$ -random neighbors on CIFAR-10 when  $k = 2$ .**

$\alpha$	Shard ( $S = 2$ )		Dirichlet ( $\beta = 0.5$ )		Mean
	0.0–0.4	0.0–0.8	0.0–0.4	0.0–0.8	
0.0	67.11	61.81	80.18	75.94	$71.26 \pm 8.33$
0.2	67.45	62.69	80.35	75.95	$71.61 \pm 8.00$
0.4	67.03	62.94	79.95	76.27	$71.55 \pm 7.90$
0.6	67.62	62.94	80.34	76.49	<b><math>71.85 \pm 7.97</math></b>
0.8	67.19	62.97	80.18	76.28	$71.76 \pm 7.80$
1.0	66.73	61.61	80.00	76.07	$71.10 \pm 8.43$

**Table 10: Test accuracy (%) with different  $\alpha$  values.**

## 5.9 Detailed Cost Analysis

We analyze the communication and computation cost of FedRN compared with other six methods, as summarized in Table 8. In this comparison, the communication cost is split into two parts, “server→client” and “client→server,” while the computation cost is split into two other parts, “forward” and “backward,” and the “total cost” aggregates all the costs for each perspective.

- **Communication Cost:** FedRN requires  $3M$  or  $4M$  communication cost in total when using one or two reliable neighbors, which is a sufficient number for general use cases. The remaining communication cost for calculating data similarity and expertise (i.e.,  $Acc(c)$  and  $p(\tilde{x}; \Theta_c)$ ) is negligible. By contrast, Co-Teaching and DivideMix both require  $4M$  communication cost in total because they maintain two models for co-training. In summary, FedRN leads to greater performance with comparable or lower communication cost compared with other baseline methods.
- **Computation Cost:** We exclude the cost of fitting GMMs and searching  $k$ -reliable neighbors because their costs are negligible compared to those of the forward and backward steps. Along the pipeline of our local update phase, FedRN first requires  $kF$  and  $kB$  computation cost to perform fine-tuning with  $k$ -reliable neighbors for one epoch. Next, to construct the clean set,  $(k + 1)F$  communication cost is needed due to the inference with  $k + 1$

available models. As the constructed clean set is deterministic over all the local epochs, only  $eF + eB$  computation cost is required to train a single model for the given local epoch  $e$  in each communication round. Consequently, FedRN needs  $(e + 2k + 1)F + (k + e)B$  computation cost in total. In a typical deep learning pipeline, the cost of the backward step is relatively smaller than that of the forward step. Hence, the computation cost for the forward step is the major issue. In light of this fact, the forward cost of FedRN is comparable to those of FedAvg, Joint-optm, and SELFIE, and less than those of Co-teaching and DivideMix, because the number of local epochs is larger than the number of reliable neighbors,  $e > k$ . Among the compared methods, DivideMix is the slowest method in that it performs  $m$  data augmentations with two models for each local epoch, so its forward cost is  $(4m + 2)eF$ .

Therefore, with one or two reliable neighbors, FedRN greatly improves the robustness of the global model without adding much communication and computation cost in the scenario of federated learning with noisy labels.

## 6 CONCLUSION

The scenario of considering both data heterogeneity and label noise is a new important research direction of practical federated learning. In this paper, we propose a novel federated learning method called FedRN. In the local update phase,  $k$ -reliable neighbors with high data expertise or similarity are retrieved, and they identify clean examples collaboratively based on their reliability scores. The models are then trained on the clean set, and the weights are averaged in the model aggregation phase. As a result, FedRN achieves high label precision and recall for sample selection, leading to a robust global model even in the presence of label noise. We conducted extensive experiments on three real-world or synthetic noisy datasets. The results verified that FedRN improves the robustness against label noise consistently and significantly. Overall, the use of reliable neighbors will inspire future studies for robustness.

## 7 ACKNOWLEDGEMENT

This work was supported by Institute of Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No.2019-0-00075, Artificial Intelligence Graduate School Program(KAIST) and No. 2022-0-00871, Development of AI Autonomy and Knowledge Enhancement for AI Agent Collaboration)

## REFERENCES

- [1] Durmus Alp Emre Acar, Yue Zhao, Ramon Matas, Matthew Mattina, Paul Whatmough, and Venkatesh Saligrama. 2020. Federated learning based on dynamic regularization. In *ICLR*.
- [2] Eric Arazo, Diego Ortego, Paul Albert, Noel O'Connor, and Kevin McGuinness. 2019. Unsupervised label noise modeling and loss correction. In *ICML*.
- [3] Devansh Arpit, Stanislaw Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. 2017. A closer look at memorization in deep networks. In *ICML*.
- [4] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin Raffel. 2019. MixMatch: A holistic approach to semi-supervised learning. In *NeurIPS*.
- [5] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir Ivanov, Chloe Kiddon, Jakub Konečný, Stefano Mazzocchi, H Brendan McMahan, et al. 2019. Towards federated learning at scale: System design. In *SysML*.
- [6] Christopher Briggs, Zhong Fan, and Peter Andras. 2020. Federated learning with hierarchical clustering of local updates to improve training on non-IID data. In *IJCNN*.
- [7] Huili Chen, Jie Ding, Eric Tramel, Shuang Wu, Anit Kumar Sahu, Salman Avestimehr, and Tao Zhang. 2022. Self-Aware Personalized Federated Learning. [arXiv preprint arXiv:2204.08069 \(2022\)](#).
- [8] Hong-You Chen and Wei-Lun Chao. 2021. FEDBE: Making bayesian model ensemble applicable to federated learning. In *ICLR*.
- [9] Jiangui Chen, Ruqing Zhang, Jiafeng Guo, Yixing Fan, and Xueqi Cheng. 2021. FedMatch: Federated Learning Over Heterogeneous Question Answering Data. In *CIKM*.
- [10] Pengfei Chen, Ben Ben Liao, Guangyong Chen, and Shengyu Zhang. 2019. Understanding and utilizing deep neural networks trained with noisy labels. In *ICML*.
- [11] Yiqiang Chen, Xiaodong Yang, Xin Qin, Han Yu, Biao Chen, and Zhiqi Shen. 2020. Focus: Dealing with label quality disparity in federated learning. In *IJCAIW*.
- [12] Liam Collins, Hamed Hassani, Aryan Mokhtari, and Sanjay Shakkottai. 2021. Exploiting Shared Representations for Personalized Federated Learning. In *ICML*.
- [13] Moming Duan, Duo Liu, Xianzhang Chen, Yujuan Tan, Jinting Ren, Lei Qiao, and Liang Liang. 2019. Astraea: Self-balancing federated learning for improving classification accuracy of mobile deep learning applications. In *ICCD*.
- [14] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi Sugiyama. 2018. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In *NeurIPS*.
- [15] Lang Huang, Chao Zhang, and Hongyang Zhang. 2020. Self-adaptive training: beyond empirical risk minimization. In *NeurIPS*.
- [16] Wonyong Jeong, Jaehong Yoon, Eunho Yang, and Sung Ju Hwang. 2021. Federated semi-supervised learning with inter-client consistency. In *ICLR*.
- [17] Di Jiang, Yuanfeng Song, Yongxin Tong, Xueyang Wu, Weiwei Zhao, Qian Xu, and Qiang Yang. [n. d.]. Federated topic modeling. In *CIKM*.
- [18] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. 2018. MentorNet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *ICML*.
- [19] Jae-Gil Lee, Yuji Roh, Hwanjun Song, and Steven Euijong Whang. 2021. Machine Learning Robustness, Fairness, and their Convergence. In *KDD*.
- [20] Junnan Li, Richard Socher, and Steven CH Hoi. 2020. DivideMix: Learning with noisy labels as semi-supervised learning. In *ICLR*.
- [21] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. 2018. Federated optimization in heterogeneous networks. [arXiv preprint arXiv:1812.06127 \(2018\)](#).
- [22] Junyu Luo, Jianlei Yang, Xucheng Ye, Xin Guo, and Weisheng Zhao. [n. d.]. FedSkel: Efficient Federated Learning on Heterogeneous Systems with Skeleton Gradients Update. In *CIKM*.
- [23] Amirhossein Malekijoo, Mohammad Javad Fadaeieslam, Hanieh Malekijoo, Morteza Homayounfar, Farshid Alizadeh-Shabdiz, and Reza Rawassizadeh. 2021. FEDZIP: A Compression Framework for Communication-Efficient Federated Learning. [arXiv preprint arXiv:2102.01593 \(2021\)](#).
- [24] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In *AISTATS*.
- [25] Khalil Muhammad, Qin Qin Wang, Diarmuid O'Reilly-Morgan, Elias Tragos, Barry Smyth, Neil Hurley, James Geraci, and Aonghus Lawlor. 2020. Fedfast: Going beyond average for faster training of federated recommender systems. In *KDD*.
- [26] Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. 2017. Making deep neural networks robust to label noise: A loss correction approach. In *CVPR*.
- [27] Hwanjun Song, Minseok Kim, and Jae-Gil Lee. 2019. SELFIE: Refurbishing unclean samples for robust deep learning. In *ICML*.
- [28] Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. 2021. Robust Learning by Self-Transition for Handling Noisy Labels. In *KDD*.
- [29] Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. 2022. Learning from noisy labels with deep neural networks: A survey. *IEEE Transactions on Neural Networks and Learning Systems (2022)*.
- [30] Daiki Tanaka, Daiki Ikami, Toshihiko Yamasaki, and Kiyoharu Aizawa. 2018. Joint optimization framework for learning with noisy labels. In *CVPR*.
- [31] Tiffany Tuor, Shiqiang Wang, Bong Jun Ko, Changchang Liu, and Kin K Leung. 2020. Overcoming noisy and irrelevant data in federated learning. In *ICPR*.
- [32] Xiaobo Xia, Tongliang Liu, Nannan Wang, Bo Han, Chen Gong, Gang Niu, and Masashi Sugiyama. 2019. Are anchor points really indispensable in label-noise learning?. In *NeurIPS*.
- [33] Jingyi Xu, Zihan Chen, Tony QS Quek, and Kai Fong Ernest Chong. 2022. FedCorr: Multi-Stage Federated Learning for Label Noise Correction. [arXiv preprint arXiv:2204.04677 \(2022\)](#).
- [34] Chengxu Yang, Qipeng Wang, Mengwei Xu, Zhenpeng Chen, Kaigui Bian, Yunxin Liu, and Xuanzhe Liu. 2021. Characterizing Impacts of Heterogeneity in Federated Learning upon Large-Scale Smartphone Data. In *TheWebConf*.
- [35] Qian Yang, Jianyi Zhang, Weituo Hao, Gregory P Spell, and Lawrence Carin. 2021. Flop: Federated learning on medical datasets using partial networks. In *KDD*.
- [36] Seunghan Yang, Hyoungseob Park, Junyoung Byun, and Changick Kim. 2022. Robust federated learning with noisy labels. *IEEE Intelligent Systems (2022)*.
- [37] Timothy Yang, Galen Andrew, Hubert Eichner, Haicheng Sun, Wei Li, Nicholas Kong, Daniel Ramage, and Françoise Beaufays. 2018. Applied federated learning: Improving google keyboard query suggestions. [arXiv preprint arXiv:1812.02903 \(2018\)](#).
- [38] Xingrui Yu, Bo Han, Jiangchao Yao, Gang Niu, Ivor Tsang, and Masashi Sugiyama. 2019. How does disagreement help generalization against label corruption?. In *ICML*.
- [39] Bixiao Zeng, Xiaodong Yang, Yiqiang Chen, Hanchao Yu, and Yingwei Zhang. 2022. CLC: A Consensus-based Label Correction Approach in Federated Learning. *ACM Transactions on Intelligent Systems and Technology (2022)*.
- [40] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. 2017. Understanding deep learning requires rethinking generalization. In *ICLR*.
- [41] Qingsong Zhang, Bin Gu, Zhiyuan Dang, Cheng Deng, and Heng Huang. [n. d.]. Desirable Companion for Vertical Federated Learning: New Zeroth-Order Gradient Based Algorithm. In *CIKM*.
- [42] Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandrara. 2018. Federated learning with non-iid data. [arXiv preprint arXiv:1806.00582 \(2018\)](#).
- [43] Guoqing Zheng, Ahmed Hassan Awadallah, and Susan Dumais. 2021. Meta label correction for noisy label learning. In *AAAI*.
- [44] Hangyu Zhu, Jinjin Xu, Shiqing Liu, and Yaochu Jin. 2021. Federated Learning on Non-IID Data: A Survey. [arXiv preprint arXiv:2106.06843 \(2021\)](#).