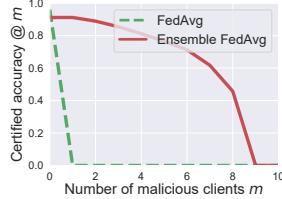
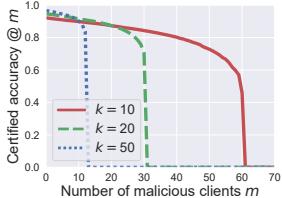


(a) MNIST

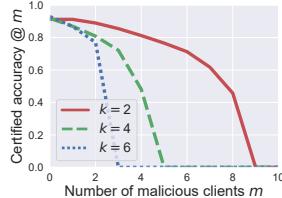


(b) HAR

Figure 2: FedAvg vs. ensemble FedAvg.



(a) MNIST



(b) HAR

Figure 3: Impact of  $k$  on our ensemble FedAvg.

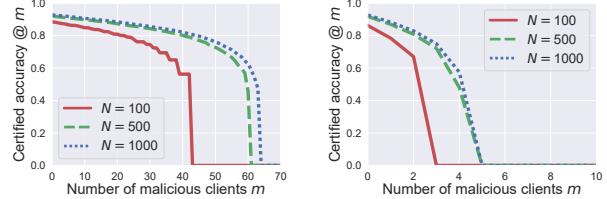
use the popular FedAvg (?) as the base federated learning algorithm. Recall that a base federated learning algorithm has hyperparameters (shown in Algorithm 1): *globalIter*, *localIter*, learning rate  $\eta$ , and batch size. Table 1 summarizes these hyperparameters for FedAvg in our experiments. In particular, we set the *globalIter* in Table 1 because FedAvg converges with such settings.

**Evaluation metric:** We use *certified accuracy* as our evaluation metric. Specifically, we define the *certified accuracy at  $m$  malicious clients* (denoted as CA@ $m$ ) for a federated learning method as the fraction of testing examples in the testing dataset  $\mathcal{D}$  whose labels are correctly predicted by the method and whose certified security levels are at least  $m$ . Formally, we define CA@ $m$  as follows:

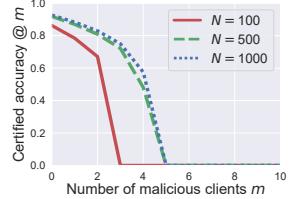
$$\text{CA}@m = \frac{\sum_{\mathbf{x}_t \in \mathcal{D}} \mathbb{I}(\hat{y}_t = y_t) \cdot \mathbb{I}(\hat{m}_t^* \geq m)}{|\mathcal{D}|}, \quad (9)$$

where  $\mathbb{I}$  is the indicator function,  $y_t$  is the true label for  $\mathbf{x}_t$ , and  $\hat{y}_t$  and  $\hat{m}_t^*$  are respectively the predicted label and certified security level for  $\mathbf{x}_t$ . Intuitively, CA@ $m$  means that when at most  $m$  clients are malicious, the accuracy of the federated learning method for  $\mathcal{D}$  is at least CA@ $m$  no matter what attacks the malicious clients use (i.e., no matter how the malicious clients tamper their local training data and model updates). Note that CA@0 reduces to the standard accuracy when there are no malicious clients.

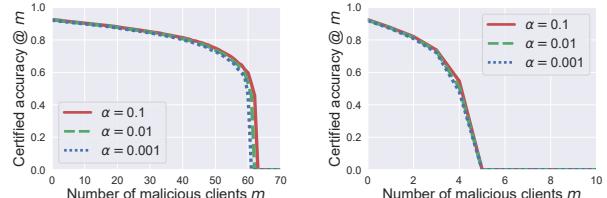
When we can compute the exact label probabilities via training  $\binom{n}{k}$  global models, the CA@ $m$  of our ensemble global model  $h$  computed using the certified security levels derived from Theorem 1 is deterministic. When  $\binom{n}{k}$  is large, we estimate predicted labels and certified security levels using Algorithm 2, and thus our CA@ $m$  has a confidence level  $1 - \alpha$  according to Theorem 3.



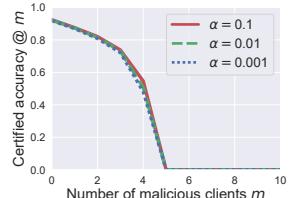
(a) MNIST



(b) HAR

Figure 4: Impact of  $N$  on our ensemble FedAvg.

(a) MNIST



(b) HAR

Figure 5: Impact of  $\alpha$  on our ensemble FedAvg.

**Parameter settings:** Our method has three parameters:  $N$ ,  $k$ , and  $\alpha$ . Unless otherwise mentioned, we adopt the following default settings for them:  $N = 500$ ,  $\alpha = 0.001$ ,  $k = 10$  for MNIST, and  $k = 2$  for HAR. Under such default setting for HAR, we have  $\binom{n}{k} = \binom{30}{2} = 435 < N = 500$  and we can compute the exact label probabilities via training 435 global models. Therefore, we have deterministic certified accuracy for HAR under the default setting. We will explore the impact of each parameter while using the default settings for the other two parameters. For HAR, we set  $k = 4$  when exploring the impact of  $N$  (i.e., Figure 4(b)) and  $\alpha$  (i.e., Figure 5(b)) since the default setting  $k = 2$  gives deterministic certified accuracy, making  $N$  and  $\alpha$  not relevant.

## Experimental Results

**Single-global-model FedAvg vs. ensemble FedAvg:** Figure 2 compares single-global-model FedAvg and ensemble FedAvg with respect to certified accuracy on the two datasets. When there are no malicious clients (i.e.,  $m = 0$ ), single-global-model FedAvg is more accurate than ensemble FedAvg. This is because ensemble FedAvg uses a subsample of clients to train each global model. However, single-global-model FedAvg has 0 certified accuracy when just one client is malicious. This is because a single malicious client can arbitrarily manipulate the global model learnt by FedAvg (?). However, the certified accuracy of ensemble FedAvg reduces to 0 when up to 61 and 9 clients (6.1% and 30%) are malicious on MNIST and HAR, respectively. Note that it is unknown whether existing Byzantine-robust federated learning methods have non-zero certified accuracy when  $m > 0$ , and thus we cannot compare ensemble FedAvg with them.

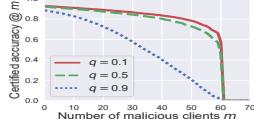


Figure 6: Impact of the degree of non-IID  $q$  on MNIST.

**Impact of  $k$ ,  $N$ , and  $\alpha$ :** Figure 3, 4, and 5 show the impact of  $k$ ,  $N$ , and  $\alpha$ , respectively.  $k$  achieves a trade-off between accuracy under no malicious clients and security under malicious clients. Specifically, when  $k$  is larger, the ensemble global model is more accurate at  $m = 0$ , but the certified accuracy drops more quickly to 0 as  $m$  increases. This is because when  $k$  is larger, it is more likely for the sampled  $k$  clients to include malicious ones. The certified accuracy increases as  $N$  or  $\alpha$  increases. This is because training more global models or a larger  $\alpha$  allows Algorithm 2 to estimate tighter probability bounds and larger certified security levels. When  $N$  increases from 100 to 500, the certified accuracy increases significantly. However, when  $N$  further grows to 1,000, the increase of certified accuracy is marginal. Our results show that we don’t need to train too many global models in practice, as the certified accuracy saturates when  $N$  is larger than some threshold.

**Impact of degree of non-IID  $q$ :** Figure 6 shows the certified accuracy of our ensemble FedAvg on MNIST when the clients’ local training data have different degrees of non-IID. We observe that the certified accuracy drops when  $q$  increases from 0.5 to 0.9, which represents a high degree of non-IID. However, the certified accuracy is still high when  $m$  is small for  $q = 0.9$ , e.g., the certified accuracy is still 83% when  $m = 10$ . This is because although each global model trained using a subsample of clients is less accurate when the local training data are highly non-IID, the ensemble of multiple global models is still accurate.

## Related Work

In federated learning, the first category of studies (?????) aim to design federated learning methods that can learn more accurate global models and/or analyze their convergence properties. For instance, FedMA (?) constructs the global model via matching and averaging the hidden elements in a neural network with similar feature extraction signatures. The second category of studies (?????????????????) aim to improve the communication efficiency between the clients and server via sparsification, quantization, and/or encoding of the model updates sent from the clients to the server. The third category of studies (?????????) aim to explore the privacy/fairness issues of federated learning and their defenses. These studies often assume a single global model is shared among the clients. Smith et al. (?) proposed to learn a customized model for each client via multi-task learning.

Our work is on security of federated learning, which is orthogonal to the studies above. Multiple studies (????) showed that the global model’s accuracy can be significantly downgraded by malicious clients. Existing defenses against

malicious clients leverage Byzantine-robust aggregation rules such as Krum (?), trimmed mean (?), coordinate-wise median (?), and Bulyan (?). However, they cannot provably guarantee that the global model’s predicted label for a testing example is not affected by malicious clients. As a result, they may be broken by strong attacks that carefully craft the model updates sent from the malicious clients to the server, e.g., (?). We propose ensemble federated learning whose predicted label for a testing example is provably not affected by a bounded number of malicious clients.

We note that ensemble methods were also proposed as provably secure defenses (e.g., (?)) against data poisoning attacks. However, they are insufficient to defend against malicious clients that can manipulate both the local training data and the model updates. In particular, a provably secure defense against data poisoning attacks guarantees that the label predicted for a testing example is unaffected by a bounded number of poisoned training examples. However, a single malicious client can poison an arbitrary number of its local training examples, breaking the assumption of provably secure defenses against data poisoning attacks.

## Conclusion

In this work, we propose ensemble federated learning and derive its tight provable security guarantee against malicious clients. Moreover, we propose an algorithm to compute the certified security levels. Our empirical results on two datasets show that our ensemble federated learning can effectively defend against malicious clients with provable security guarantees. Interesting future work includes estimating the probability bounds deterministically and considering the internal structure of a base federated learning algorithm to further improve our provable security guarantees.

## Acknowledgement

We thank the anonymous reviewers for insightful reviews. This work was supported by NSF grant No.1937786.

Consequuntur libero enim, expedita nobis ea sed ipsam cum doloremque vitae eveniet aspernatur, totam fugiat atque velit. Quisquam temporibus ipsam iusto rerum sunt velit excepturi iste, provident unde error totam ab quisquam deserunt nam consequatur numquam suscipit, nemo maxime incident sed ullam dolorem? Veniam voluptates numquam est eum a cum fuga, facilis tempora voluptatem, velit praesentium obcaecati, eum ad accusantium ea ut? Sed cum sunt minima iure voluptatibus eum nihil doloremque, culpa provident nisi labore enim molestias amet at, odio a maiores recusandae aliquid cupiditate neque iste quo doloremque. Laborum eos optio ipsam earum, quae quibusdam adipisci expedita id dignissimos alias porro, nisi dolor molestias placeat officia deleniti quaerat dolorem, minima explicabo voluptatibus illum nulla libero deleniti quis dignissimos atque, in ut cumque natus at autem consecetur? Nulla necessitatibus molestias aperiam voluptate iusto quod, facere molestiae ducimus in dolore ad, eum