

---

# To Federate or Not To Federate: Incentivizing Client Participation in Federated Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Federated learning (FL) facilitates collaboration between a group of clients who  
2 seek to train a common machine learning model without directly sharing their  
3 local data. Although there is an abundance of research on improving the speed,  
4 efficiency, and accuracy of federated training, most works implicitly assume that  
5 all clients are willing to participate in the FL framework. Due to data heterogeneity,  
6 however, the global model may not work well for some clients, and they may  
7 instead choose to use their own local model. Such disincentivization of clients can  
8 be problematic from the server’s perspective because having more participating  
9 clients yields a better global model, and offers better privacy guarantees to the  
10 participating clients. In this paper, we propose an algorithm called INCFL that  
11 explicitly maximizes the fraction of clients who are incentivized to use the global  
12 model by dynamically adjusting the aggregation weights assigned to their updates.  
13 Our experiments show that INCFL increases the number of incentivized clients by  
14 30-55% compared to standard federated training algorithms, and can also improve  
15 the generalization performance of the global model on unseen clients.

## 16 1 Introduction

17 Federated learning (FL) is a distributed learning framework that enables the training of a machine  
18 learning model using a network of clients (e.g., mobile phones, hospitals) [1], without having to  
19 transfer the clients’ data to a central server. In the standard FL framework [1–6], clients perform a  
20 few updates using their local data, and a central server aggregates these updates into a single global  
21 model. As the global model is based on the union of all the local datasets, it is expected to generalize  
22 well for the entire client population. However, due to heterogeneity of the data between clients [7],  
23 not all clients stand to benefit from federation. While FL can produce a model that performs well  
24 on average, for some clients, it may perform even worse than a model trained in isolation on their  
25 limited local data. Our experiments (Section 4) demonstrate that local models trained in isolation on  
26 FL benchmarks can indeed outperform global models obtained by commonly used FL algorithms.

27 When a client participates in FL, it incurs the cost of contributing its local data and computational  
28 resources to the federation in return for receiving a global model. However, if a local model trained  
29 in isolation is better than the global model, the client may not be incentivized to participate in  
30 FL—causing it to opt out of contributing to and using the global model in the future. This lack of  
31 incentives for clients to participate in FL can be problematic from the server’s perspective. Having a  
32 large pool of clients willing to participate in training is beneficial, if not imperative, to ensure the  
33 performance of FL models [8, 9]. When a large number of clients participate, the global model is  
34 based on a larger pool of data, allowing better generalization to new clients that may join in the  
35 future. Having a larger number of participating clients can also improve privacy-utility trade-offs by  
36 mitigating the impact of each individual client on the global model [10–12].

37 In this work, we seek to answer the following pertinent question: *How can we actively incentivize*  
 38 *clients to use and contribute to a federated global model, rather than training local models in*  
 39 *isolation?* To address this question, we propose an algorithm called INCFL to train a global model  
 40 that dramatically improves the fraction of incentivized clients in comparison to standard FL algorithms.  
 41 Our key contributions are summarized as follows.

- 42 • In Section 2 we formalize the notion of client incentives by defining a metric called the incentivized  
 43 participation rate (IPR), which measures the fraction of clients willing to participate in the FL  
 44 framework. We propose to maximize a sigmoid relaxation of the IPR, which makes the objective  
 45 differentiable and enables the use of common gradient-based optimization algorithms.
- 46 • In Section 3 we propose a federated algorithm called INCFL to maximize incentivized client  
 47 participation. INCFL dynamically adjusts the weight assigned to each client’s local update when  
 48 aggregating the updates at the central server. The method allows for partial client availability for  
 49 training, it is applicable to general non-convex objectives (with convergence guarantees), and it is  
 50 stateless (does not require clients to maintain local parameters during training).
- 51 • In Section 4, we empirically validate the performance of INCFL by comparing it with standard FL  
 52 algorithms for multiple data sets. INCFL is able to increase the number of incentivized clients by  
 53 30-55%, and also ensures that the global model generalizes well to unseen clients.

54 As surveyed in [13], previous works investigating client incentives in FL have typically done so from  
 55 a game-theoretic perspective and for toy problems such as mean estimation. In contrast, our work  
 56 is generally applicable to non-convex objectives, and considers a server that seeks to train a single  
 57 global model that will be preferred by the maximum number of clients, thus incentivizing them to  
 58 participate in FL. We provide a more detailed review of prior work in Appendix A below.

## 59 2 Problem Formulation

60 We consider a FL setup where  $M$  clients are connected to a central server. For each client  $k \in$   
 61  $\{1, 2, \dots, M\}$ , its true local loss function is given by  $f_k(\mathbf{w}) = \mathbb{E}_{\xi \sim \mathcal{D}_k} [\ell(\mathbf{w}, \xi)]$  where  $\mathcal{D}_k$  is the true  
 62 data distribution of client  $k$ , and  $\ell(\mathbf{w}, \xi)$  is the composite loss function for the model  $\mathbf{w} \in \mathbb{R}^d$  for data  
 63 sample  $\xi$ . In practice, each client only has access to its local training dataset  $\mathcal{B}_k$  with  $|\mathcal{B}_k| = N_k$  data  
 64 samples sampled from  $\mathcal{D}_k$ . Client  $k$ ’s empirical local loss function is  $F_k(\mathbf{w}) = \frac{1}{|\mathcal{B}_k|} \sum_{\xi \in \mathcal{B}_k} \ell(\mathbf{w}, \xi)$ .  
 65 Our setup is applicable to both cross-device and cross-silo FL as we do not make any specific  
 66 assumptions about the nature of the clients or their constraints.

67 **Defining Client Incentives in FL.** The goal of each client is to find a model that minimizes its  
 68 true loss function, which we denote as  $\mathbf{w}_k^* := \operatorname{argmin}_{\mathbf{w}} f_k(\mathbf{w})$ . To do so, we consider that each  
 69 client does some solo training on its local dataset  $\mathcal{B}_k$  to obtain an approximate local model  $\hat{\mathbf{w}}_k$ . For  
 70 example,  $\hat{\mathbf{w}}_k$  can be found by running a few steps of SGD on the empirical loss  $F_k(\mathbf{w})$ . Since the  
 71 local dataset size is in general small,  $\hat{\mathbf{w}}_k$  may not generalize well to the true distribution  $\mathcal{D}_k$  of a  
 72 client. Therefore we say that a client is incentivized to participate in FL (i.e., use the federated model)  
 73 if the federated model gives better generalization performance than its local model.

74 **Definition 1 (Client Incentive).** *Given a global model  $\mathbf{w}$ , client  $k \in [M]$  is said to be incentivized to*  
 75 *participate in FL if  $f_k(\mathbf{w}) < f_k(\hat{\mathbf{w}}_k)$ , that is, the global model is better than its own local model.*

76 In practice, clients can have a separate validation dataset on which they compare the losses of the  
 77 global model and their local model to decide if they are incentivized to participate. In general,  $f_k(\hat{\mathbf{w}}_k)$   
 78 in Definition 1 acts as a performance benchmark for the global model and can also be replaced by a  
 79 different value depending on the specific need of a client.

80 **Standard FL Objective does not account for Client Incentives.** In standard FL, clients collab-  
 81 oratively minimize the objective  $F(\mathbf{w}) = \sum_{k=1}^M p_k F_k(\mathbf{w})$ , where the aggregation weights  $p_k$  are  
 82 usually set as  $p_k \propto |\mathcal{B}_k|$ . Observe that this objective function does not consider client incentives as  
 83 defined in Definition 1 and implicitly assumes that all clients will participate in training and use the  
 84 global model. However, due to clients’ data heterogeneity, this assumption may not hold in general.

85 **Incentivized Participation Rate (IPR).** Based on Definition 1, we formulate the following metric to  
 86 explicitly measure the fraction of clients that are incentivized for a given federated global model  $\mathbf{w}$ :

$$\text{Incentivized Participation Rate (IPR)} = \frac{1}{M} \sum_{k=1}^M \mathbb{I}\{f_k(\mathbf{w}) < f_k(\widehat{\mathbf{w}}_k)\}, \quad (1)$$

87 where  $\mathbb{I}$  is the indicator function. Note that IPR only looks at whether or not a client is incentivized  
 88 and not *how much* a client is incentivized (or disincentivized) since the decision to participate in  
 89 FL is binary. Another variation of (1) could be to measure the incentive margin of clients, e.g.  
 90  $\sum_k \max\{f_k(\widehat{\mathbf{w}}_k) - f_k(\mathbf{w}), 0\}$ , but this does not capture the motivation behind our work which is  
 91 to improve the number of the incentivized clients in FL. To the best of our knowledge, a similar  
 92 indicator based metric has not been explored previously in the FL literature.

## 93 2.1 Proposed INCFL Objective

94 A naïve approach to increase the number of incentivized clients with our definition of client incentives  
 95 in (1) is directly maximizing the IPR as follows:

$$\max_{\mathbf{w}} \left[ \frac{1}{M} \sum_{k=1}^M \mathbb{I}\{f_k(\mathbf{w}) < f_k(\widehat{\mathbf{w}}_k)\} \right] = \min_{\mathbf{w}} \left[ \frac{1}{M} \sum_{k=1}^M \text{sign}(f_k(\mathbf{w}) - f_k(\widehat{\mathbf{w}}_k)) \right]. \quad (2)$$

96 where  $\text{sign}(x) = 1$  if  $x \geq 0$  and 0 otherwise. There are two immediate difficulties in minimizing (2).  
 97 First, clients may not know their true data distribution  $\mathcal{D}_k$  to compute  $f_k(\mathbf{w}) - f_k(\widehat{\mathbf{w}}_k)$ . Secondly,  
 98 the sign function makes the objective nondifferentiable and limits the use of common gradient-based  
 99 methods. We resolve these issues by proposing a "proxy" for (2) with the following relaxations.

- 100 1. **Replacing the Sign function with the Sigmoid function  $\sigma(\cdot)$  [14]:** Replacing the non-  
 101 differentiable 0-1 loss with a smooth differentiable loss is a standard tool used in optimiza-  
 102 tion [15, 16]. Given the many candidates (e.g. hinge loss, ReLU, sigmoid), we find that using the  
 103 sigmoid function is essential for our objective to faithfully approximate the true objective in (2).  
 104 We discuss theoretical implications of using the sigmoid loss in more detail in Appendix B.1.
- 105 2. **Replacing  $\sigma(f_k(\mathbf{w}) - f_k(\widehat{\mathbf{w}}_k))$  with  $\sigma(F_k(\mathbf{w}) - F_k(\widehat{\mathbf{w}}_k))$ :** As clients do not have access to  
 106 their true distribution  $\mathcal{D}_k$  to compute  $f_k(\cdot)$  we propose to use an empirical estimate  $\sigma(F_k(\mathbf{w}) -$   
 107  $F_k(\widehat{\mathbf{w}}_k))$ . Since  $\widehat{\mathbf{w}}_k$  is locally trained, it is likely that  $F_k(\widehat{\mathbf{w}}_k) < f_k(\widehat{\mathbf{w}}_k)$ . On the other hand, the  
 108 global model  $\mathbf{w}$  is trained on the data of all clients, making it unlikely to overfit to the local data  
 109 of any particular client, leading to  $f_k(\mathbf{w}) \approx F_k(\mathbf{w})$  (see Appendix F.2). Hence, in most cases  
 110 we have  $f_k(\mathbf{w}) - f_k(\widehat{\mathbf{w}}_k) < F_k(\mathbf{w}) - F_k(\widehat{\mathbf{w}}_k)$  and since sigmoid is an increasing function, this  
 111 implies that  $\sigma(f_k(\mathbf{w}) - f_k(\widehat{\mathbf{w}}_k)) < \sigma(F_k(\mathbf{w}) - F_k(\widehat{\mathbf{w}}_k))$ . Therefore, with this relaxation we are  
 112 effectively trying to minimize an upper bound on our true objective.

113 With these relaxations, we present our proposed INCFL objective:

$$\text{INCFL Obj.} : \quad \min_{\mathbf{w}} \widetilde{F}(\mathbf{w}) = \min_{\mathbf{w}} \frac{1}{M} \sum_{i=1}^M \widetilde{F}_i(\mathbf{w}), \text{ where } \widetilde{F}_i(\mathbf{w}) := \sigma(F_i(\mathbf{w}) - F_i(\widehat{\mathbf{w}}_i)). \quad (3)$$

114 Our experimental results in Section 4 support our intuition of these relaxations and convincingly  
 115 demonstrate that minimizing our proposed objective leads to a much higher IPR than the standard FL  
 116 objective. Before discussing the details of how we minimize our objective, we take a closer look at  
 117 how our objective behaves for a mean estimation problem.

## 118 3 Proposed INCFL Algorithm

119 With the sigmoid approximation of the sign loss and for differentiable  $F_k(\mathbf{w})$ , our objective  $\widetilde{F}(\mathbf{w})$  in  
 120 (3) is differentiable and can be minimized with gradient descent and variants. Its gradient is given by:

$$\nabla \widetilde{F}(\mathbf{w}) = \frac{1}{M} \sum_{k=1}^M \underbrace{(1 - \widetilde{F}_k(\mathbf{w})) \widetilde{F}_k(\mathbf{w})}_{\text{aggregating weight} := q_k(\mathbf{w})} \nabla F_k(\mathbf{w}). \quad (4)$$

121 Observe that  $\nabla \tilde{F}(\mathbf{w})$  is a **weighted aggregate** of the gradi-  
 122 ents of the clients’ empirical losses, similar in spirit to the  
 123 gradient  $\nabla F(\mathbf{w})$  in standard FL. The key difference is that  
 124 in INCFL, the weights  $q_k(\mathbf{w}) := (1 - \tilde{F}_k(\mathbf{w}))\tilde{F}_k(\mathbf{w})$  are  
 125 *incentive-dependent*, and are dynamically updated based  
 126 on the current model  $\mathbf{w}$ , as we discuss below.

127 **Behavior of the Aggregation Weights  $q_k(\mathbf{w})$ .** For a  
 128 given  $\mathbf{w}$ , the behavior of the aggregation weights  $q_k(\mathbf{w})$   
 129 depend on the *empirical incentive gap*,  $F_k(\mathbf{w}) - F_k(\hat{\mathbf{w}}_k)$   
 130 (see Fig. 1) since  $\tilde{F}_k(\mathbf{w}) = \sigma(F_k(\mathbf{w}) - F_k(\hat{\mathbf{w}}_k))$ . When  
 131  $F_k(\mathbf{w}) \ll F_k(\hat{\mathbf{w}}_k)$ , it implies that the global model  $\mathbf{w}$   
 132 performs much better than the local model  $\hat{\mathbf{w}}_k$  at client  
 133  $k$ . Hence INCFL sets  $q_k(\mathbf{w}) \approx 0$  to focus on the updates  
 134 of other clients. Similarly if  $F_k(\mathbf{w}) \gg F_k(\hat{\mathbf{w}}_k)$ , INCFL  
 135 will set  $q_k(\mathbf{w}) \approx 0$ . This is because  $F_k(\mathbf{w}) \gg F_k(\hat{\mathbf{w}}_k)$   
 136 implies that the current model  $\mathbf{w}$  is incompatible with the  
 137 local model  $\hat{\mathbf{w}}_k$  at client  $k$  and hence it is better to avoid  
 138 optimizing for this client at the risk of disincentivizing  
 139 other clients. INCFL gives the highest weight to those  
 140 clients for which the global model performs similar to  
 141 their local models, i.e.  $F_k(\mathbf{w}) \approx F_k(\hat{\mathbf{w}}_k)$ , since this allows  
 142 it to increase IPR without hurting other clients’ performance.

143 **A Practical INCFL Solver.** Directly minimizing the INCFL objective using gradient descent can  
 144 be slow to converge and impractical since it requires all clients to be available for training. Instead,  
 145 we propose a practical INCFL algorithm, which uses multiple local updates at each client to speed  
 146 up convergence as done in standard FL [1] and allow for partial client availability. We replace the  
 147 gradient  $\nabla F_k(\mathbf{w})$  with the *local update*  $\Delta \mathbf{w}_k$  at a client, and aggregate these updates only from  
 148 clients that are available in that round.

149 Let us use the superscript  $(t, r)$  to denote the communication round  $t$  and local iteration index  $r$ . At  
 150 each round  $t$ , the server selects a new set of clients  $\mathcal{S}^{(t,0)}$  uniformly at random and sends the most  
 151 recent global model  $\mathbf{w}^{(t,0)}$  to clients in  $\mathcal{S}^{(t,0)}$ . The clients in  $\mathcal{S}^{(t,0)}$  then perform  $\tau$  local iterations  
 152 with local learning rate  $\eta$  to calculate their updates as follows:

$$\text{Perform Local SGD: } \mathbf{w}_k^{(t,r+1)} = \mathbf{w}_k^{(t,r)} - \eta \mathbf{g}(\mathbf{w}_k^{(t,r)}, \xi_k^{(t,r)}) \quad \text{for all } r \in \{0, \dots, \tau - 1\}, \quad (5)$$

$$\text{Compute Local Update: } \Delta \mathbf{w}_k^{(t,0)} = \mathbf{w}_k^{(t,\tau)} - \mathbf{w}_k^{(t,0)}, \quad (6)$$

153 where  $\mathbf{g}(\mathbf{w}_k^{(t,r)}, \xi_k^{(t,r)}) = \frac{1}{b} \sum_{\xi \in \xi_k^{(t,r)}} \nabla f(\mathbf{w}_k^{(t,r)}, \xi)$  is the stochastic gradient computed using a  
 154 mini-batch  $\xi_k^{(t,r)}$  of size  $b$  that is randomly sampled from client  $k$ ’s local dataset  $\mathcal{B}_k$ . The weight  
 155  $q_k(\mathbf{w}_k^{(t,0)})$  can be computed at each client by calculating the loss over its training data with  $\mathbf{w}_k^{(t,0)}$   
 156 which is a simple inference step. Clients in  $\mathcal{S}^{(t,0)}$  then send back their local updates  $\Delta \mathbf{w}_k^{(t,0)}$  and  
 157 weights  $q_k(\mathbf{w}_k^{(t,0)})$  to the server which updates the global model as follows:

$$\text{Global Update Rule: } \mathbf{w}^{(t+1,0)} = \mathbf{w}^{(t,0)} - \eta_g^{(t,0)} \sum_{k \in \mathcal{S}^{(t,0)}} q_k(\mathbf{w}^{(t,0)}) \Delta \mathbf{w}_k^{(t,0)}, \quad (7)$$

158 where  $\eta_g^{(t,0)} = \frac{\eta_g}{\sum_{k \in \mathcal{S}^{(t,0)}} q_k(\mathbf{w}^{(t,0)}) + \epsilon}$  is the adaptive server learning rate with a fixed global learning  
 159 rate  $\eta_g$  and constant  $\epsilon > 0$ . We discuss the reasoning for such an adaptive learning rate along with  
 160 the pseudo code and convergence bounds of our INCFL in Appendix C.

## 161 4 Experimental Results

162 We evaluate INCFL in four different settings: (i) logistic regression on a synthetic federated dataset  
 163 (Synthetic(1,1) [2]), (ii) MLP trained on non-iid partitioned FMNIST [17], (iii) CNN trained on  
 164 non-iid partitioned CIFAR10 [18], and (iv) MLP for sentiment classification trained on Sent140 [19].  
 165 We compare INCFL with well-known stateless FL algorithms that train a single model such as

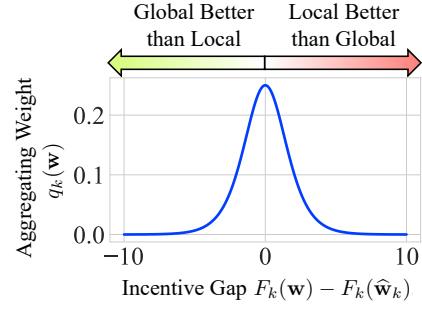


Figure 1: Aggregating weight  $q_k(\mathbf{w})$  for any client  $k$  versus the empirical incentive gap  $F_k(\mathbf{w}) - F_k(\hat{\mathbf{w}}_k)$ . The weight  $q_k(\mathbf{w})$  is small for clients that already have a very large incentive (global much better than local) or no incentive at all (local much better than global), and is highest for clients that are moderately incentivized (global similar to local).



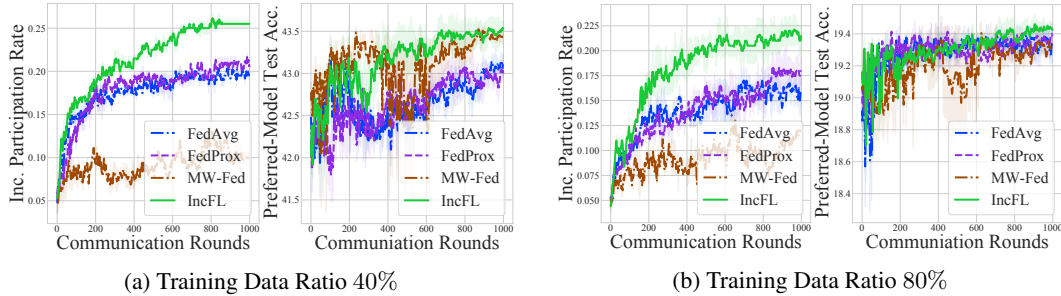


Figure 2: Incentivized participation rate (IPR), i.e., the fraction of incentivized clients, and preferred-model test accuracy evaluated on the test data for the synthetic data with the training clients. INCFL improves on both IPR and preferred-model test accuracy for both smaller (40%) and larger (80%) training data ratios where the IPR improvement of INCFL is larger for the smaller training data ratio.

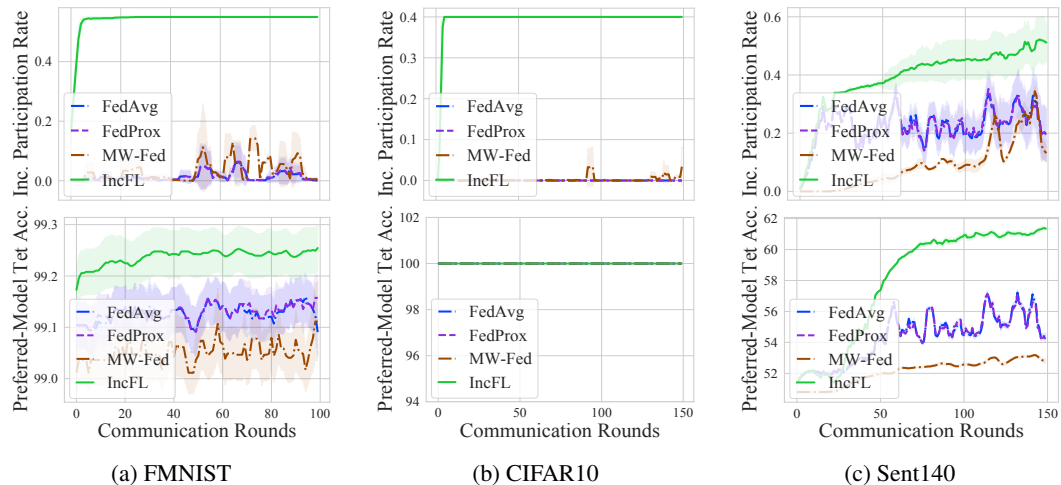


Figure 3: Incentivized participation rate (*upper*), i.e., the fraction of clients incentivized to participate in FL, and preferred-model test accuracy (*lower*) for the training clients' test data with different datasets. For all datasets, INCFL achieves at least 30% to up to 55% increase in the fraction of incentivized clients, while also achieving the maximum preferred-model test accuracy.

166 standard FedAvg [1], FedProx [2] which aims to tackle data heterogeneity, PerFedAvg [20] which  
 167 facilitates personalization, and MW-Fed [21] which incentivizes client participation. We provide  
 168 results on personalization jointly used with INCFL in Appendix F.2.

169 **Setup.** For Sent140, we consider 308 clients and for the other datasets we have 100 clients that  
 170 are used for training in FL. These clients are active at some point in training the global model, and  
 171 we name them as ‘seen clients’. In all experiments, 10 clients are sampled every communication  
 172 round. For FMNIST, data is partitioned into 5 clusters where 2 labels are assigned for each cluster  
 173 with no labels overlapping across clusters. Clients are randomly assigned to each cluster, and within  
 174 each cluster, clients are homogeneously distributed with the assigned labels. We similarly partition  
 175 CIFAR10, where clients are partitioned into 10 clusters with 1 label assigned to each cluster. For the  
 176 Sent140 dataset, clients are naturally partitioned with their twitter IDs. We also generate ‘unseen  
 177 clients’ the same way we generate the seen clients, with 619 clients for Sent140 and 100 clients for  
 178 all other datasets. These unseen clients represent new incoming clients that have not been seen before  
 179 during the training rounds of FL. The data of each client is partitioned to 60% : 40% for training and  
 180 test data ratio unless mentioned otherwise. Further details are deferred to Appendix F.1.

181 **Evaluation Metrics: IPR and Preferred-model Test Accuracy.** We evaluate INCFL and other  
 182 methods with these two key metrics: 1) Incentivized Participation Rate (IPR), defined in (1) and 2)  
 183 Preferred-Model Test Accuracy. Recall that IPR is the fraction of clients incentivized to participate in  
 184 FL and use the global server. The preferred-model test accuracy is the average of the clients’ test

Table 1: Incentivized participation rate (IPR) and preferred-model test accuracy of the final global models trained with different algorithms for the unseen clients’ test data.

	Incentivized Participation Rate (IPR)			Preferred-Model Test Acc.		
	FMNIST	CIFAR10	Sent140	FMNIST	CIFAR10	Sent140
FedAvg	0.08 ( $\pm 0.01$ )	0.00 ( $\pm 0.00$ )	0.37 ( $\pm 0.07$ )	98.53 ( $\pm 0.13$ )	100.00 ( $\pm 0.00$ )	57.05 ( $\pm 1.44$ )
FedProx	0.07 ( $\pm 0.01$ )	0.00 ( $\pm 0.00$ )	0.37 ( $\pm 0.07$ )	98.43 ( $\pm 0.21$ )	100.00 ( $\pm 0.00$ )	57.07 ( $\pm 1.42$ )
MW-Fed	0.05 ( $\pm 0.04$ )	0.02 ( $\pm 0.02$ )	0.17 ( $\pm 0.03$ )	98.32 ( $\pm 0.13$ )	100.00 ( $\pm 0.00$ )	55.57 ( $\pm 1.28$ )
INCFL	<b>0.55</b> ( $\pm 0.00$ )	<b>0.40</b> ( $\pm 0.00$ )	<b>0.43</b> ( $\pm 0.05$ )	<b>98.83</b> ( $\pm 0.06$ )	100.00 ( $\pm 0.00$ )	<b>57.16</b> ( $\pm 1.35$ )

185 accuracies computed on their preferred model, either the global model or their solo-trained local  
 186 model. Higher IPR is beneficial to the server, and higher preferred-model test accuracy is beneficial  
 187 to the clients. Thus, it is desirable for an algorithm to improve both these metrics.

188 **Incentivizing the Participation of the Seen Clients.** We first discuss the performance for seen  
 189 clients used during the training of the global model. In Fig. 2, we show that for the synthetic data,  
 190 INCFL incentivizes at least 5% more clients compared to the other baselines. The preferred-model  
 191 test accuracy achieved by INCFL is also highest amongst other baselines. Hence INCFL provides  
 192 a win-win for both the server and clients since clients have the highest accuracy from choosing the  
 193 better model and the server has the highest fraction of incentivized clients. In Fig. 3, for all DNN  
 194 experiments, INCFL significantly improves the IPR to at least 30% to at most 55%. Fig. 3 also shows  
 195 that the baselines can fail in incentivizing the clients with even 0% clients incentivized. INCFL also  
 196 improves on the preferred-model test accuracy than the other baselines for FMNIST and Sent140,  
 197 corroborating the win-win for both the server and clients. For CIFAR10, the preferred-model test  
 198 accuracy is 100% for all baselines while the IPR is significantly higher for INCFL. This shows  
 199 that while clients can always achieve 100% by either choosing the local or global model for best  
 200 performance, server can only gain a large fraction of incentivized clients when using INCFL.

201 **Incentivizing the Participation of the Unseen Clients.** We now show that INCFL is also better at  
 202 incentivizing the unseen clients that were not active during the training of the global model such as  
 203 new incoming clients. In Table 1, we show that INCFL consistently improves the IPR of such clients  
 204 by at least 40% for FMNIST and CIFAR 10, and 6% for Sent140. INCFL also achieves higher or at  
 205 least the same preferred-model test accuracy compared to that of all baselines for all datasets.

206 **Effect of Training Data Ratio.** In Fig. 2, we show the performance of INCFL with different ratios  
 207 of the training data to test data split for each client’s data. One can expect that if a client has a high  
 208 training data ratio, the solo-trained local model of a client sufficiently generalizes well to its test data,  
 209 and hence the client will be less incentivized to participate in FL. We show in that even if a client has  
 210 a high training data ratio (80% in Fig. 2(b)), INCFL is able to increase the fraction of incentivized  
 211 clients compared to other baselines but the improvement is smaller compared to when clients have  
 212 smaller training data ratio (40% in Fig. 2(a)). In general, clients are believed to have very few labeled  
 213 training data [22, 23], in which case INCFL can improve the fraction of incentivized clients greatly.

## 214 5 Concluding Remarks

215 In this work we carefully re-examine the fundamental assumption in FL that clients always stand to  
 216 benefit from federation. To do so, we formalize a intuitive notion of client incentives in FL based  
 217 on whether a global model has better generalization performance than a client’s local model. We  
 218 introduce a novel metric termed as Incentivized Participation Rate (IPR) to explicitly measure the  
 219 fraction of incentivized clients in FL and develop a corresponding framework INCFL to maximize  
 220 IPR. In contrast to existing work, INCFL allows the server to play an *active* role in incentivizing  
 221 clients by dynamically adjusting its aggregation procedure while training the global model. Moreover  
 222 INCFL is well-suited to both cross-device and cross-silo FL since it stateless and allows partial client  
 223 availability while training. We provide convergence guarantees for INCFL and show that in practice  
 224 it can dramatically improve IPR compared to standard FL. We believe our work will open up new  
 225 research directions in understanding the role played by the server in incentivizing clients for FL.  
 226 Future work includes jointly examining client incentives with privacy guarantees offered in FL.

227 **References**

- 228 [1] H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agøura y Arcas.  
229 Communication-Efficient Learning of Deep Networks from Decentralized Data. *International*  
230 *Conference on Artificial Intelligence and Statistics (AISTATS)*, April 2017.
- 231 [2] Anit Kumar Sahu, Tian Li, Maziar Sanjabi, Manzil Zaheer, Ameet Talwalkar, and Virginia  
232 Smith. Federated optimization for heterogeneous networks. In *Proceedings of the 3rd MLSys*  
233 *Conference*, January 2020.
- 234 [3] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J Reddi, Sebastian U Stich,  
235 and Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for on-device  
236 federated learning. *arXiv preprint arXiv:1910.06378*, 2019.
- 237 [4] Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H. Vincent Poor. Tackling the Objective  
238 Inconsistency Problem in Heterogeneous Federated Optimization. *preprint*, May 2020.
- 239 [5] A Khaled, K Mishchenko, and P Richtárik. Tighter theory for local SGD on identical and  
240 heterogeneous data. In *The 23rd International Conference on Artificial Intelligence and Statistics*  
241 *(AISTATS 2020)*, 2020.
- 242 [6] Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný,  
243 Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. In *International*  
244 *Conference on Learning Representations (ICLR)*, 2021.
- 245 [7] El Mahdi Chayti, Sai Praneeth Karimireddy, Sebastian U. Stich, Nicolas Flammarion, and Martin  
246 Jaggi. Linear speedup in personalized collaborative learning. *arXiv preprint arXiv:2111.05968*,  
247 2022.
- 248 [8] Zachary Charles, Zachary Garrett, Zhouyuan Huo, Sergei Shmulyian, and Virginia Smith. On  
249 large-cohort training for federated learning. In *Advances in neural information processing*  
250 *systems*, 2021.
- 251 [9] Jianyu Wang, Zachary Charles, Zheng Xu, Gauri Joshi, H Brendan McMahan, Maruan Al-  
252 Shedivat, Galen Andrew, Salman Avestimehr, Katharine Daly, Deepesh Data, et al. A field  
253 guide to federated optimization. *arXiv preprint arXiv:2107.06917*, 2021.
- 254 [10] K. Wei, J. Li, M. Ding, C. Ma, H. H. Yang, F. Farokhi, S. Jin, T. Q. S. Quek, and H. Vincent  
255 Poor. Federated learning with differential privacy: Algorithms and performance analysis. *IEEE*  
256 *Transactions on Information Forensics and Security*, 15:3454–3469, 2020.
- 257 [11] Robin C. Geyer, Tassilo Klein, and Moin Nabi. Differentially private federated learning: A  
258 client level perspective. *NIPS Workshop: Machine Learning on the Phone and other Consumer*  
259 *Devices*, 2017.
- 260 [12] H Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially  
261 private recurrent language models. *International Conference on Learning Representations*,  
262 2018.
- 263 [13] Xuezhen Tu, Kun Zhu, Nguyen Cong Luong, Dusit Niyato, Yang Zhang, and Juan Li. Incentive  
264 mechanisms for federated learning: From economic and game theoretic perspective. *arXiv*  
265 *preprint arXiv:2111.11850*, 2021.
- 266 [14] David H Von Seggern. *CRC standard curves and surfaces with mathematica*. CRC Press, 2017.
- 267 [15] Tan T. Nguyen and Scott Sanner. Algorithms for direct 0–1 loss optimization in binary  
268 classification. In *Proceedings of the 30th International Conference on Machine Learning*,  
269 volume 28, 2013.
- 270 [16] Hamed Masnadi-shirazi and Nuno Vasconcelos. On the design of loss functions for classification:  
271 theory, robustness to outliers, and savageboost. In *Advances in Neural Information Processing*  
272 *Systems(NIPS)*, 2008.
- 273 [17] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for  
274 benchmarking machine learning algorithms. <https://arxiv.org/abs/1708.07747>, aug 2017.

- 275 [18] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Learning multiple layers of features from  
276 tiny images. *CIFAR-10 (Canadian Institute for Advanced Research)*, 2009.
- 277 [19] Alec Go, Richa Bhayani, and Lei Huang. Twitter sentiment classification using distant supervi-  
278 sion. *CS224N Project Report, Stanford*, 2009.
- 279 [20] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning with  
280 theoretical guarantees: A model-agnostic meta-learning approach. In *Advances in Neural*  
281 *Information Processing Systems*, 2020.
- 282 [21] Avrim Blum, Nika Haghtalab, Richard Lanus Phillips, and Han Shao. One for one, or all for all:  
283 Equilibria and optimality of collaboration in federated learning. In *International Conference on*  
284 *Machine Learning*, 2021.
- 285 [22] Wonyong Jeong, Jaehong Yoon, Eunho Yang, and Sung Ju Hwang. Federated semi-supervised  
286 learning with inter-client consistency & disjoint learning. *International Conference on Learning*  
287 *Representations (ICLR)*, 2021.
- 288 [23] Nan Lu, Zhao Wang, Xiaoxiao Li, Gang Niu, Qi Dou, and Masashi Sugiyama. Federated learn-  
289 ing from only unlabeled data with class-conditional-sharing clients. *International Conference*  
290 *on Learning Representations (ICLR)*, 2022.
- 291 [24] Kate Donahue and Jon Kleinberg. Model-sharing games: Analyzing federated learning under  
292 voluntary participation. In *The Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI-*  
293 *21)*, 2021.
- 294 [25] Kate Donahue and Jon Kleinberg. Optimality and stability in federated learning: A game-  
295 theoretic approach. In *Advances in Neural Information Processing Systems*, 2021.
- 296 [26] Jingoo Han, Ahmad Faraz Khan, Syed Zawad, Ali Anwar, Nathalie Baracaldo Angel, Yi Zhou,  
297 Feng Yan, and Ali R. Butt. Tokenized incentive for federated learning. In *Proceedings of the*  
298 *Federated Learning Workshop at the Association for the Advancement of Artificial Intelligence*  
299 *(AAAI) Conference*, 2022.
- 300 [27] Jiawen Kang, Zehui Xiong, Dusit Niyato, Han Yu, Ying-Chang Liang, and Dong In Kim.  
301 Incentive design for efficient federated learning in mobile networks: A contract theory approach.  
302 In *2019 IEEE VTS Asia Pacific Wireless Communications Symposium (APWCS)*, pages 1–5,  
303 2019.
- 304 [28] Meng Zhang, Ermin Wei, and Randall Berry. Faithful edge federated learning: Scalability and  
305 privacy. *IEEE Journal on Selected Areas in Communications*, 39(12):3790–3804, 2021.
- 306 [29] M. Zhang, K. Sapra, S. Fidler, S. Yeung, and J. M. Alvarez. Personalized federated learning  
307 with first order model optimization. In *International Conference on Learning Representations*  
308 *(ICLR)*, 2021.
- 309 [30] Felix Grimberg, Mary-Anne Hartley, Sai P. Karimireddy, and Martin Jaggi. Optimal model  
310 averaging: Towards personalized collaborative learning. In *International Workshop on Federated*  
311 *Learning for User Privacy and Data Confidentiality in Conjunction with ICML 2021 (FL-*  
312 *ICML'21)*, 2021.
- 313 [31] Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated  
314 learning through personalization. In *Proceedings of the 38th International Conference on*  
315 *Machine Learning*, 2021.
- 316 [32] Yishay Mansour, Mehryar Mohri, Jae Ro, and Ananda Theertha Suresh. Three approaches  
317 for personalization with applications to federated learning. *arXiv preprint cs.LG/2002.10619*,  
318 2020.
- 319 [33] Tian Li, Maziari Sanjabi, and Virginia Smith. Fair resource allocation in federated learning.  
320 *arXiv preprint arXiv:1905.10497*, 2019.

- 321 [34] Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In  
322 Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International*  
323 *Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*,  
324 pages 4615–4625, Long Beach, California, USA, 09–15 Jun 2019. PMLR.
- 325 [35] Yûsaku Komatu. Elementary inequalities for mills’ ratio. *Reports of Statistical Application*  
326 *Research (Union of Japanese Scientific Engineers)*, 4:69–70, 1955.
- 327 [36] Sebastian U Stich. Local SGD converges fast and communicates little. In *International*  
328 *Conference on Learning Representations (ICLR)*, 2019.
- 329 [37] I. Bistriz, A. J. Mann, and N. Bambos. Distributed distillation for on-device learning. In  
330 *Advances in Neural Information Processing Systems*, 2020.
- 331 [38] Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Stochastic convex  
332 optimization. In *COLT*, 2009.
- 333 [39] Erlend S. Riis, Matthias J. Ehrhardt, G. R. W. Quispel, and Carola-Bibiane Schönlieb. A geomet-  
334 ric integration approach to nonsmooth, nonconvex optimisation. *Foundations of Computational*  
335 *Mathematics*, 2021.
- 336 [40] Divyansh Jhunjhunwala, Pranay Sharma, Aushim Nagarkatti, and Gauri Joshi. Fedvarp: Tack-  
337 ling the variance due to partial client participation in federated learning. *arXiv*, 2022.
- 338 [41] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen,  
339 David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern,  
340 Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fer-  
341 nández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler  
342 Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array  
343 programming with NumPy. *Nature*, 585(7825):357–362, September 2020.
- 344 [42] Pauli Virtanen, Ralf Gommers, Travis E Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau,  
345 Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, et al. Scipy 1.0:  
346 fundamental algorithms for scientific computing in python. *Nature methods*, 17(3):261–272,  
347 2020.
- 348 [43] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for  
349 word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages  
350 1532–1543, 2014.

351 **A Related Work**

352 **Game-Theoretic Study of Client Coalitions.** Similar to our work, [24, 25] consider the problem  
 353 of client incentives where each client can either train a local model in isolation or join a coalition with  
 354 other similar clients to train a common model. The authors study the ratio between the stable solution  
 355 (where no client is incentivized to shift to a different coalition) and the optimal solution (where  
 356 the sum of losses of all clients is minimized). Although this line of work establishes a number of  
 357 useful game-theoretic insights, it focuses on simple mean estimation and linear regression problems.  
 358 The server plays the role of a passive matchmaker, facilitating the formation of coalitions of clients.  
 359 In contrast, in our work, the server actively seeks to train a single global model that maximizes  
 360 client participation. This perspective alleviates some of the analysis complexities occurring in  
 361 game-theoretic formulations and allows us to consider general non-convex objective functions.

362 **Client Incentives for Contributing Data and Computation Resources.** Other recent works  
 363 such as [21, 26–28] develop mechanisms to incentivize clients to contribute data samples and local  
 364 computational resources to federated training, and compensate for these contributions. In [21], which  
 365 is closest to our work, the authors consider linear problems and analyze the existence and optimality  
 366 of equilibria in the problem of equitably distributing the burden of data contribution between clients.  
 367 They propose a heuristic algorithm called Multiplicative Weight (MW)-Fed, where the server instructs  
 368 clients for whom the global model is performing poorly to conduct more local updates. Unlike our  
 369 work, this approach does not explicitly maximize client participation, and it is not supported by  
 370 theoretical guarantees. However, we include it as a baseline in our experiments (Section 4).

371 **Personalized and Fair Federated Learning.** Finally, personalized or fair FL methods may offer  
 372 an auxiliary benefit of increasing incentivized client participation even though they are not formally  
 373 studied in the client incentives context. Instead of having all clients use a common global model,  
 374 personalized FL methods consider learning models unique to each client. In cross-device settings, a  
 375 common approach is to consider methods that fine-tune the global model to produce personalized  
 376 models [29, 30, 7, 31, 32]. This can naturally incentivize more clients to use the global model and  
 377 participate in FL. Our INCFL algorithm is orthogonal to and can be combined with such personalized  
 378 FL methods. In our experiments, we demonstrate the performance of INCFL combined with local  
 379 fine-tuning and compare it with [20], which uses meta-learning for personalization. Another related  
 380 area is fair FL, where a common goal is to train a global model whose accuracy has less variance  
 381 across the client population than standard FedAvg [33, 34]. A side benefit of these methods is that  
 382 they can incentivize the worst performing clients to participate. However, a downside is that the  
 383 performance of the global model may be degraded for the best performing clients, thus incentivizing  
 384 them to leave the federation. We show in additional experimental results in Appendix F.1 that  
 385 common fair FL methods are indeed not effective in improving the overall client participation rate.

386 **B Mean Estimation as a Toy Example for INCFL**

387 **B.1 Maximizing IPR in Two Client Mean Estimation**

388 We consider a setup with  $M = 2$  clients where each client aims to find the mean of its data distribution  
 389 by minimizing the true loss function  $f_k(w) = \mathbb{E}_{\xi_k} [(w - \xi_k)^2]$ ,  $\xi_k \sim \mathcal{N}(\theta_k, \nu^2) \forall k \in [2]$ . In  
 390 practice, clients only have  $N_k$  samples drawn from their distribution denoted by  $\mathcal{B}_k = \{e_{k,j}\}_{j=1}^{N_k}$  and  
 391 can only minimize their empirical loss function given by  $F_k(w) = \frac{1}{|\mathcal{B}_k|} \sum_{j=1}^{N_k} (w - e_{k,j})^2$ . Then the  
 392 solo trained models at each client will be their local empirical mean, i.e.  $\hat{w}_k = \hat{\theta}_k = \frac{1}{|\mathcal{B}_k|} \sum_{j=1}^{N_k} e_{k,j}$ .

393 **IPR for Standard FL Model Decreases Exponentially with Heterogeneity.** For simplicity let  
 394 us assume  $N_1 = N_2 = N$ . Let  $\gamma^2 = \nu^2/N$  be the variance of the local empirical means and  
 395  $\gamma_G^2 = ((\theta_1 - \theta_2)/2)^2 > 0$  be a measure of heterogeneity between the true means. The standard  
 396 FL objective will always set the FL model to be the average of the local empirical means (i.e.  
 397  $w = (\hat{\theta}_1 + \hat{\theta}_2)/2$ ) and does not take into account the heterogeneity among the clients. As a result,  
 398 the IPR of the global model decreases *exponentially* as  $\gamma_G^2$  increases.

399 **Lemma B.1.** *The expected IPR of the standard FL model is upper bounded by  $2 \exp\left(-\frac{\gamma_G^2}{5\gamma^2}\right)$ , where  
 400 the expectation is taken over the randomness in the local datasets  $\mathcal{B}_1, \mathcal{B}_2$ .*

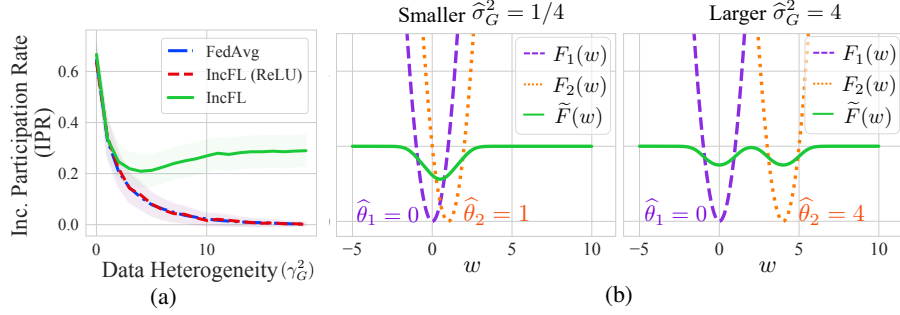


Figure 4: Results for the two client mean estimation in Appendix B.1; (a): IPR for FedAvg decays exponentially while IPR for INCFL is lower bounded by a constant. Replacing the sigmoid approximation with ReLU approximation in INCFL leads to the same solution as FedAvg; (b): IncFL adapts to the heterogeneity of the problem—for small heterogeneity it encourages collaboration by having a single global minima, for large heterogeneity it encourages separation by having far away local minimas.

401 **Maximizing IPR with Relaxed Objective.** We now explicitly maximize the IPR for this setting by  
 402 solving for a relaxed version of the objective in (2) as proposed earlier. We replace the true loss  $f_k(\cdot)$   
 403 by the empirical loss  $F_k(\cdot)$  and replace the 0-1 (sign) loss with a differentiable approximation  $h(\cdot)$ .

404 We first show that setting  $h(\cdot)$  to be a standard convex surrogate for the 0-1 loss (e.g. log loss,  
 405 exponential loss, ReLU) leads to our new objective behaving the same as the standard FL objective.

406 **Lemma B.2.** *Let  $h$  be any function that is convex, twice differentiable, and strictly increasing in*  
 407  *$[0, \infty)$ . Then our relaxed objective is strictly convex and has a unique minimizer at  $w^* = \left(\frac{\hat{\theta}_1 + \hat{\theta}_2}{2}\right)$ .*

408 **Maximizing the INCFL Objective Leads to Increased IPR.** Based on Lemma B.2, we see that  
 409 we need nonconvexity in  $h(\cdot)$  for the objective to behave differently than standard FL. We set  
 410  $h(x) = \sigma(x) = \frac{\exp(x)}{1 + \exp(x)}$ , as proposed in our INCFL objective in (3). We find that the INCFL

411 objective *adapts* to the empirical heterogeneity parameter  $\hat{\gamma}_G^2 = \left(\frac{\hat{\theta}_2 - \hat{\theta}_1}{2}\right)^2$ . If  $\hat{\gamma}_G^2 < 1$  (small data  
 412 heterogeneity), the objective encourages collaboration by setting the global model to be the average  
 413 of the local models. On the other hand, if  $\hat{\gamma}_G^2 > 2$  (large data heterogeneity), the objective encourages  
 414 *separation* by setting the global model close to either the local model of the first client or the local  
 415 model of the second client (see Fig. 4). Based on this observation, we have the following theorem.

416 **Theorem B.1.** *Let  $w$  be a local minima of the INCFL objective. The expected IPR using  $w$  is lower*  
 417 *bounded by  $\frac{1}{16} \exp\left(-\frac{1}{\hat{\gamma}^2}\right)$  where the expectation is over the randomness in the local dataset  $\mathcal{B}_1, \mathcal{B}_2$ .*

418 Note that our result above is independent of the heterogeneity parameter  $\hat{\gamma}_G^2$ . Therefore even with  
 419  $\hat{\gamma}_G^2 \gg 0$ , INCFL will keep incentivizing atleast one client by adapting its objective accordingly.  
 420 Additional discussion and proof details can be found in Appendix B.

421 We begin by recalling the setup discussed in Appendix B.1. We have a setup with  $M = 2$  clients  
 422 where each client aims to find the mean of its data distribution by minimizing the true loss function  
 423  $f_k(w) = \mathbb{E}_{\xi_k} [(w - \xi_k)^2]$ ,  $\xi_k \sim \mathcal{N}(\theta_k, \nu^2) \forall k \in [2]$ . Without loss of generality we assume  
 424 that  $\theta_2 \geq \theta_1$ . In practice, each client has  $N$  samples drawn from their distribution denoted by  
 425  $\mathcal{B}_k = \{e_{k,j}\}_{j=1}^N$  and can only minimize their empirical loss function given by

$$F_k(w) = \frac{1}{N} \sum_{j=1}^N (w - e_{k,j})^2 \quad (8)$$

$$= (w - \hat{\theta}_k)^2 + \frac{1}{N} \sum_{j=1}^N (\hat{\theta}_k - e_{k,j})^2 \quad (9)$$

426 where  $\hat{\theta}_k = \frac{1}{N} \sum_{j=1}^N e_{k,j}$  is the empirical mean at client  $k$ . We assume that clients set their solo  
 427 trained models as their empirical mean, i.e.  $\hat{w}_k = \hat{\theta}_k$ .

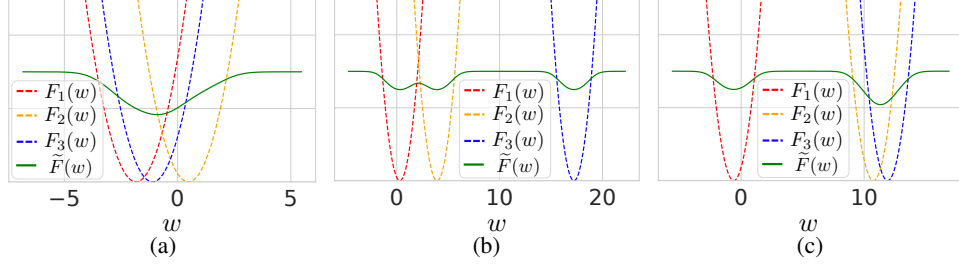


Figure 5: Results for the three client mean estimation; (a): case 1 when the true mean across clients are close to amongst each other where IncFL’s optimal solution is identical to that of FedAvg; (b): case 2 when the true mean across clients are all different from each other where IncFL’s optimal solution ensures that at least one of the clients will be incentivized participate with IncFL’s global model (unlike FedAvg); (c) case 3 when two clients’ true means are close to each other while the other client has a different mean. IncFL in this case, is able to ensure that the two clients participate while FedAvg is not able to make any client participate.

428 We define the following quantities

$$\gamma^2 := \frac{\nu^2}{N}; \quad \gamma_G^2 = \left( \frac{\theta_2 - \theta_1}{2} \right)^2; \quad (10)$$

429 Note that the distribution of the empirical means itself follows a normal distribution following the  
430 linear additivity of independent normal random variables.

$$\hat{\theta}_1 \sim \mathcal{N}(\theta_1, \gamma^2); \quad \hat{\theta}_2 \sim \mathcal{N}(\theta_2, \gamma^2) \quad (11)$$

### 431 B.2 Maximizing IPR in Three Client Mean Estimation.

432 We further examine the property of IncFL to incentivize clients with a 3 clients toy example which is  
433 an extension from what we have shown in Appendix B.1 for 2 clients. Reusing the notation from  
434 Appendix B.1, where  $\theta_i$  is the true mean at client  $i$  and  $\hat{\theta}_i \sim \mathcal{N}(\theta_i, 1)$  is the empirical mean of a  
435 client, our analysis can be divided into the following cases for the 3 client example (see Fig. 5):

- 436 • Case 1:  $\theta_1 \approx \theta_2 \approx \theta_3$ : This case captures the setting where the data at the clients is almost i.i.d. In  
437 this case, it makes sense for clients to collaborate together and therefore IncFL’s optimal solution  
438 will be the average of local empirical means (same as FedAvg).
- 439 • Case 2:  $\theta_1 \neq \theta_2 \neq \theta_3$ : This case captures the setting where the data at clients is completely  
440 disparate. In this case, none of the clients benefit from collaborating and therefore IncFL’s optimal  
441 solution will be the local model of one of the clients. This ensures atleast one of the clients will  
442 still be incentivized to use the IncFl global model unlike FedAvg.
- 443 • Case 3:  $\theta_1 \approx \theta_2 \neq \theta_3$ : The most interesting case happens when data at two of the clients is  
444 similar but the data at the third client is different. Without loss of generality we assume that data  
445 at clients 1 and 2 is similar and client 3 is different. In this case, although client 1 and 2 benefit  
446 from federating, FedAvg is unable to leverage that due to the heterogeneity at client 3. IncFL, on  
447 the other hand, will set the optimal solution to be the average of the local models of just client 1  
448 and client 2. This ensures that clients 1 and 2 will both continue to participate in the FL training  
449 process, thus maximizing the number of incentivized clients.

450 The behavior of IncFL in the three client setup clearly highlights the non-trivialness of our proposed  
451 IncFL’s formulation.

### 452 B.3 Proof of Lemma B.1

453 **Lemma B.1** *The expected IPR of the standard FL model is upper bounded by  $2 \exp\left(-\frac{\gamma_G^2}{5\gamma^2}\right)$ , where  
454 the expectation is taken over the randomness in the local datasets  $\mathcal{B}_1, \mathcal{B}_2$ .*



455 **Proof.**

456 The standard FL model is given by,

$$w = \frac{\hat{\theta}_1 + \hat{\theta}_2}{2} \quad (12)$$

457 Therefore the expected IPR is,

$$\mathbb{E} \left[ \frac{\mathbb{I}\{(w - \theta_1)^2 < (\hat{\theta}_1 - \theta_1)^2\} + \mathbb{I}\{(w - \theta_2)^2 < (\hat{\theta}_2 - \theta_2)^2\}}{2} \right] \quad (13)$$

$$= \frac{1}{2} \left[ \underbrace{\mathbb{P}\left((w - \theta_1)^2 < (\hat{\theta}_1 - \theta_1)^2\right)}_{T_1} + \underbrace{\mathbb{P}\left((w - \theta_2)^2 < (\hat{\theta}_2 - \theta_2)^2\right)}_{T_2} \right] \quad (14)$$

458 Next we bound  $T_1$  and  $T_2$ . Bounding  $T_1$  :

$$T_1 = \mathbb{P}\left((w - \theta_1)^2 < (\hat{\theta}_1 - \theta_1)^2\right) \quad (15)$$

$$= \mathbb{P}\left(\left(\frac{\hat{\theta}_1 + \hat{\theta}_2}{2} - \theta_1\right)^2 < (\hat{\theta}_1 - \theta_1)^2\right) \quad (16)$$

$$= \mathbb{P}\left(\left(\frac{\hat{\theta}_2 - \hat{\theta}_1}{2}\right)^2 + 2\left(\frac{\hat{\theta}_2 - \hat{\theta}_1}{2}\right)(\hat{\theta}_1 - \theta_1) < 0\right) \quad (17)$$

$$= \mathbb{P}\left(\left\{\left(\frac{\hat{\theta}_2 - \hat{\theta}_1}{2}\right)^2 + 2\left(\frac{\hat{\theta}_2 - \hat{\theta}_1}{2}\right)(\hat{\theta}_1 - \theta_1) < 0\right\} \cap \{\hat{\theta}_2 > \hat{\theta}_1\}\right) \\ + \mathbb{P}\left(\left\{\left(\frac{\hat{\theta}_2 - \hat{\theta}_1}{2}\right)^2 + 2\left(\frac{\hat{\theta}_2 - \hat{\theta}_1}{2}\right)(\hat{\theta}_1 - \theta_1) < 0\right\} \cap \{\hat{\theta}_2 \leq \hat{\theta}_1\}\right) \quad (18)$$

$$= \mathbb{P}\left(\left\{\left(\frac{\hat{\theta}_2 - \hat{\theta}_1}{2}\right) + 2(\hat{\theta}_1 - \theta_1) < 0\right\} \cap \{\hat{\theta}_2 > \hat{\theta}_1\}\right) \\ + \mathbb{P}\left(\left\{\left(\frac{\hat{\theta}_2 - \hat{\theta}_1}{2}\right) + 2(\hat{\theta}_1 - \theta_1) < 0\right\} \cap \{\hat{\theta}_2 \leq \hat{\theta}_1\}\right) \quad (19)$$

$$\leq \mathbb{P}\left(\left(\frac{\hat{\theta}_2 - \hat{\theta}_1}{2}\right) + 2(\hat{\theta}_1 - \theta_1) < 0\right) + \mathbb{P}\left(\hat{\theta}_2 - \hat{\theta}_1 \leq 0\right) \quad (20)$$

$$= \mathbb{P}(Z_1 < 0) + \mathbb{P}(Z_2 \leq 0) \quad \text{where } Z_1 \sim \mathcal{N}\left(\gamma_G, \frac{5}{2}\gamma^2\right), Z_2 \sim \mathcal{N}(2\gamma_G, 2\gamma^2) \quad (21)$$

$$\leq \exp\left(-\frac{\gamma_G^2}{5\gamma^2}\right) + \exp\left(-\frac{\gamma_G^2}{\gamma^2}\right) \quad (22)$$

$$\leq 2 \exp\left(-\frac{\gamma_G^2}{5\gamma^2}\right) \quad (23)$$

459 where (18) uses  $\mathbb{P}(A) = \mathbb{P}(A \cap B) + \mathbb{P}(A \cap B^c)$ , (20) uses  $\mathbb{P}(A \cap B) \leq \mathbb{P}(A)$ , (21) uses (11)  
460 and linear additivity of independent normal random variables, (22) uses a Chernoff bound.

461 We can similarly bound  $T_2$  to get  $T_2 \leq 2 \exp\left(-\frac{\gamma_G^2}{5\gamma^2}\right)$ . Thus the expected IPR of the standard FL  
462 model is upper bounded by  $2 \exp\left(-\frac{\gamma_G^2}{5\gamma^2}\right)$ .

463 **B.4 Proof of Lemma B.2**

464 **Lemma B.2** *Let  $h$  be any function that is convex, twice differentiable, and strictly increasing in*  
 465  *$[0, \infty)$ . Then our relaxed objective is strictly convex and has a unique minimizer at  $w^* = \left(\frac{\hat{\theta}_1 + \hat{\theta}_2}{2}\right)$ .*

466 **Proof.**

467 Let us denote our relaxed objective by  $v(w)$ . Then  $v(w)$  can be written as,

$$v(w) = \frac{1}{2} [h(F_1(w) - F(\hat{w}_1)) + h(F_2(w) - F(\hat{w}_2))] \quad (24)$$

$$= \underbrace{\frac{1}{2}h\left((w - \hat{\theta}_1)^2\right)}_{v_1(w)} + \underbrace{\frac{1}{2}h\left((w - \hat{\theta}_2)^2\right)}_{v_2(w)} \quad (25)$$

$$(26)$$

468 We first prove that  $v_1(w)$  is strictly convex. Let  $\lambda \in (0, 1)$  and  $(w_1, w_2)$  be any pair of points in  $\mathbb{R}^2$   
 469 such that  $w_1 \neq w_2$ . We have,

$$v_1(\lambda w_1 + (1 - \lambda)w_2) = \frac{1}{2}h\left((\lambda(w_1 - \hat{\theta}_1) + (1 - \lambda)(w_2 - \hat{\theta}_1))^2\right) \quad (27)$$

$$< \frac{1}{2}h\left(\lambda(w_1 - \hat{\theta}_1)^2 + (1 - \lambda)(w_2 - \hat{\theta}_1)^2\right) \quad (28)$$

$$\leq \frac{\lambda}{2}h\left((w_1 - \hat{\theta}_1)^2\right) + \frac{1 - \lambda}{2}h\left((w_2 - \hat{\theta}_1)^2\right) \quad (29)$$

$$= \lambda v_1(w_1) + (1 - \lambda)v_1(w_2) \quad (30)$$

470 where (28) follows from the strict convexity of  $f(w) = w^2$  and the fact that  $h(w)$  is strictly increasing  
 471 in the range  $[0, \infty)$ , (29) follows from the convexity of  $h(w)$ .

472 This completes the proof that  $v_1(w)$  is strictly convex. We can similarly prove that  $v_2(w)$  is strictly  
 473 convex and hence  $v(w)$  is strictly convex since summation of strictly convex functions is strictly  
 474 convex.

475 Also note that,

$$\nabla v(w) = \nabla h\left((w - \hat{\theta}_1)^2\right)(w - \hat{\theta}_1) + \nabla h\left((w - \hat{\theta}_2)^2\right)(w - \hat{\theta}_2) \quad (31)$$

476 It is easy to see that  $\nabla v(w) = 0$  at  $w = \left(\frac{\hat{\theta}_1 + \hat{\theta}_2}{2}\right)$ . Since  $v(w)$  is strictly convex this implies that  
 477  $w^* = \left(\frac{\hat{\theta}_1 + \hat{\theta}_2}{2}\right)$  will be a unique global minimizer. This completes the proof.

478 **B.5 Proof of Theorem B.1**

479 Before stating the proof of Theorem 3.1 we first state some intermediate results that will be used in  
 480 the proof.

481 The INCFLL objective can be written as,

$$v(w) = \frac{1}{2}\sigma\left((w - \hat{\theta}_1)^2\right) + \frac{1}{2}\sigma\left((w - \hat{\theta}_2)^2\right) \quad (32)$$

482 where  $\sigma(w) = 1/(1 + \exp(-w))$ .

483 We additionally define the following quantities,

$$i := \operatorname{argmin}\{\hat{\theta}_1, \hat{\theta}_2\}; \quad j := \operatorname{argmax}\{\hat{\theta}_1, \hat{\theta}_2\}; \quad \hat{\gamma}_G := \frac{\hat{\theta}_j - \hat{\theta}_i}{2} \quad (33)$$

484 Let  $q(w) = \sigma(w)(1 - \sigma(w))$ . The gradient of  $v(w)$  is given as,

$$\nabla v(w) = q\left((w - \hat{\theta}_1)^2\right)(w - \hat{\theta}_1) + q\left((w - \hat{\theta}_2)^2\right)(w - \hat{\theta}_2) \quad (34)$$

485 **Lemma B.3** For  $\widehat{\gamma}_G > 2$ ,  $w = \left(\frac{\widehat{\theta}_1 + \widehat{\theta}_2}{2}\right)$  will be a local maxima of the INCFL objective.

486 It is easy to see that  $w = \left(\frac{\widehat{\theta}_1 + \widehat{\theta}_2}{2}\right)$  will always be a stationary point of  $\nabla v(w)$ . Our goal is to  
 487 determine whether it will be a local minima or a local maxima. To do so, we calculate the hessian of  
 488  $v(w)$  as follows. Let  $f(w) = 2\sigma(w)(1 - \sigma(w))(1 - 2\sigma(w))$ . Then,

$$\nabla^2 v(w) = \underbrace{f\left(\left(w - \widehat{\theta}_1\right)^2\right)\left(w - \widehat{\theta}_1\right)^2 + q\left(\left(w - \widehat{\theta}_1\right)^2\right)}_{h_1(w)} + \underbrace{f\left(\left(w - \widehat{\theta}_2\right)^2\right)\left(w - \widehat{\theta}_2\right)^2 + q\left(\left(w - \widehat{\theta}_2\right)^2\right)}_{h_2(w)} \quad (35)$$

489 Note that  $h_1(w) = h_2(w)$  for  $w = \left(\frac{\widehat{\theta}_1 + \widehat{\theta}_2}{2}\right)$ . Hence it suffices to focus on the condition for which  
 490  $h_1(w) < 0$  at  $w = \left(\frac{\widehat{\theta}_1 + \widehat{\theta}_2}{2}\right)$ . We have,

$$h_1\left(\left(\widehat{\theta}_1 + \widehat{\theta}_2\right)/2\right) = f(\widehat{\gamma}_G^2)\widehat{\gamma}_G^2 + q(\widehat{\gamma}_G^2) \quad (36)$$

$$= q(\widehat{\gamma}_G^2)(2(1 - 2\sigma(\widehat{\gamma}_G^2))\widehat{\gamma}_G^2 + 1) \quad (37)$$

$$< 0 \quad \text{for } \widehat{\gamma}_G \geq 1.022 \quad (38)$$

491 where the last inequality follows from the fact that  $q(w) > 0$  for all  $w \in \mathbb{R}$  and  $2(1 - 2\sigma(w^2))w^2 + 1 <$   
 492  $0$  for  $w \geq 1.022$ . Thus for  $\widehat{\gamma}_G > 2$ ,  $w = \left(\frac{\widehat{\theta}_1 + \widehat{\theta}_2}{2}\right)$  will be a local maxima of the INCFL objective.

493 **Lemma B.4** For  $\widehat{\gamma}_G > 0$ , any local minima of  $v(w)$  lies in the range  $(\widehat{\theta}_i, \widehat{\theta}_i + 2] \cup [\widehat{\theta}_j - 2, \widehat{\theta}_j)$ .

494 Firstly note that since  $\widehat{\gamma}_G > 0$  we have  $\widehat{\theta}_j > \widehat{\theta}_i$ . Secondly note that since  $q(w) > 0$  for all  $w \in \mathbb{R}$ ,  
 495  $\nabla v(w) < 0$  for all  $w \leq \widehat{\theta}_i$  and  $\nabla v(w) > 0$  for all  $w \geq \widehat{\theta}_j$ . Therefore any root of the function  $\nabla v(w)$   
 496 must lie in the range  $(\widehat{\theta}_i, \widehat{\theta}_j)$ .

497 **Case 1:**  $0 < \widehat{\gamma}_G \leq 2$ .

498 In this case, the lemma is trivially satisfied since  $(\widehat{\theta}_i, \widehat{\theta}_j) \subset \left\{(\widehat{\theta}_i, \widehat{\theta}_i + 2] \cup [\widehat{\theta}_j - 2, \widehat{\theta}_j)\right\}$ .

499 **Case 2:**  $\widehat{\gamma}_G > 2$ .

500 Let  $x = w - \widehat{\theta}_i$  and  $g(x) = q(x^2)x$ . We can write  $\nabla v(w)$  as,

$$\nabla v(\widehat{\theta}_i + x) = g(x) - g(2\widehat{\gamma}_G - x) \quad (39)$$

501 It can be seen that for  $x > 2$ ,  $g(x)$  is a decreasing function. For  $x \in (2, \widehat{\gamma}_G)$  we have  $x >$   
 502  $2\widehat{\gamma}_G - x$  which implies  $g(x) > g(2\widehat{\gamma}_G - x)$ . Therefore  $\nabla v(\widehat{\theta}_i + x) > 0$  for  $x \in (2, \widehat{\gamma}_G)$ . Also  
 503  $\nabla v(\widehat{\theta}_i + 2\widehat{\gamma}_G - x) = -\nabla v(\widehat{\theta}_i + x)$  and therefore  $\nabla v(\widehat{\theta}_i + x) < 0$  for  $x \in (\widehat{\gamma}_G, 2\widehat{\gamma}_G - 2)$ .  
 504  $\nabla v(\widehat{\theta}_i + \widehat{\gamma}_G) = 0$  but this will be a local maxima for  $\widehat{\gamma}_G > 2$  as shown in Lemma B.3. Thus there  
 505 exists no local minima of  $v(w)$  for  $w \in (\widehat{\theta}_i + 2, \widehat{\theta}_j - 2)$

506 Combining both cases we see that any local minima of  $v(w)$  lies in the range  
 507  $\left\{(\widehat{\theta}_i, \widehat{\theta}_i + 2] \cup [\widehat{\theta}_j - 2, \widehat{\theta}_j)\right\}$ .

508 **Theorem B.1** Let  $w$  be a local minima of the INCFL objective. The expected IPR using  $w$  is lower  
 509 bounded by  $\frac{1}{16} \exp\left(-\frac{1}{\gamma^2}\right)$  where the expectation is over the randomness in the local dataset  $\mathcal{B}_1, \mathcal{B}_2$ .

510 **Proof.**

511 The IPR can be written as,

$$\frac{1}{2} \left[ \mathbb{P}\left((w - \theta_i)^2 < (\hat{\theta}_i - \theta_i)^2\right) + \mathbb{P}\left((w - \theta_j)^2 < (\hat{\theta}_j - \theta_j)^2\right) \right] \quad (40)$$

512 We focus on the case where  $\hat{\theta}_2 \neq \hat{\theta}_i$  implying  $\hat{\theta}_j > \hat{\theta}_i$  ( $\hat{\theta}_2 = \hat{\theta}_1$  is a zero-probability event and does  
 513 not affect our proof). Let  $w$  be any local minima of the INCFL objective. From Lemma B.4 we know  
 514 that  $w$  will lie in the range  $(\hat{\theta}_i, \hat{\theta}_i + 2] \cup [\hat{\theta}_j - 2, \hat{\theta}_j)$

515 **Case 1:**  $w \in (\hat{\theta}_i, \hat{\theta}_i + 2]$

$$\mathbb{P}\left((w - \theta_i)^2 < (\hat{\theta}_i - \theta_i)^2\right) = \mathbb{P}\left((w - \hat{\theta}_i)^2 + 2(w - \hat{\theta}_i)(\hat{\theta}_i - \theta_i) < 0\right) \quad (41)$$

$$= \mathbb{P}\left((w - \hat{\theta}_i) + 2(\hat{\theta}_i - \theta_i) < 0\right) \quad (42)$$

$$\geq \mathbb{P}\left(2 + 2(\hat{\theta}_i - \theta_i) < 0\right) \quad (43)$$

$$= \mathbb{P}\left(\hat{\theta}_i - \theta_i < -1\right) \quad (44)$$

$$\geq \mathbb{P}\left(\{\hat{\theta}_1 < \hat{\theta}_2\} \cap \{\hat{\theta}_1 - \theta_1 < -1\}\right) \quad (45)$$

$$= \mathbb{P}\left(\hat{\theta}_1 < \hat{\theta}_2\right) \mathbb{P}\left(\hat{\theta}_1 - \theta_1 < -1 | \hat{\theta}_1 < \hat{\theta}_2\right) \quad (46)$$

$$\geq \mathbb{P}\left(\hat{\theta}_1 < \hat{\theta}_2\right) \mathbb{P}\left(\hat{\theta}_1 - \theta_1 < -1\right) \quad (47)$$

$$= \mathbb{P}\left(\hat{\theta}_1 < \hat{\theta}_2\right) \mathbb{P}\left(Z > 1/\gamma\right) \quad \text{where } Z \sim \mathcal{N}(0, 1) \quad (48)$$

$$\geq \frac{1}{8} \exp\left(-\frac{1}{\gamma^2}\right) \quad (49)$$

516 (42) uses the fact that  $(w - \hat{\theta}_i) > 0$ , (43) uses  $(w - \hat{\theta}_i) \leq 2$ , (45) uses  $\mathbb{P}(A) \geq \mathbb{P}(A \cap B)$  and  
 517 definition of  $i$ . (47) uses the following argument. If  $\theta_1 - 1 \geq \hat{\theta}_2$  then  $\mathbb{P}\left(\hat{\theta}_1 - \theta_1 < -1 | \hat{\theta}_1 < \hat{\theta}_2\right) =$   
 518 1. If  $\theta_1 - 1 < \hat{\theta}_2$  then  $\mathbb{P}\left(\hat{\theta}_1 - \theta_1 < -1 | \hat{\theta}_1 < \hat{\theta}_2\right) = \mathbb{P}\left(\hat{\theta}_1 - \theta_1 < -1\right) / \mathbb{P}\left(\hat{\theta}_1 < \hat{\theta}_2\right) \geq$   
 519  $\mathbb{P}\left(\hat{\theta}_1 - \theta_1 < -1\right)$ . (48) uses  $\hat{\theta}_1 - \theta_1 \sim \mathcal{N}(0, \gamma^2)$ , (49) uses  $\mathbb{P}\left(\hat{\theta}_1 < \hat{\theta}_2\right) \geq \frac{1}{2}$  and  $\mathbb{P}(Z \geq x) \geq$   
 520  $\frac{2 \exp(-x^2/2)}{\sqrt{2\pi}(\sqrt{4+x^2}+x)} \geq \frac{1}{4} \exp(-x^2)$  where  $Z \sim \mathcal{N}(0, 1)$  [35].

521 In the case where  $w \in (\hat{\theta}_j - 2, \hat{\theta}_j]$  a similar technique can be used to lower bound  
 522  $\mathbb{P}\left((w - \theta_j)^2 < (\hat{\theta}_j - \theta_j)^2\right)$ . Thus the IPR of any local minima of the INCFL objective is lower  
 523 bounded by  $\frac{1}{16} \exp\left(-\frac{1}{\gamma^2}\right)$ .

524 **C Additional Discussion on Our Proposed INCFL Algorithm**

525 We first present our pseudo-code for INCFL below in Algorithm 1.

---

**Algorithm 1** Proposed Client Incentivizing FL Framework: INCFL

---

- 1: **Input:** mini-batch size  $b$ , local iteration steps  $\tau$ , training loss  $F_i(\widehat{\mathbf{w}}_i)$  for each client  $i \in [M]$
  - 2: **Output:** Global model  $\mathbf{w}^{(T,0)}$ , **Initialize:** Global model  $\mathbf{w}^{(0,0)}$
  - 3: **For**  $t = 0, \dots, T - 1$  **communication rounds do:**
  - 4:   **Global server do:**
  - 5:     Select  $m$  clients for  $\mathcal{S}^{(t,0)}$  uniformly at random and send  $\mathbf{w}^{(t,0)}$  to clients in  $\mathcal{S}^{(t,0)}$
  - 6:   **Clients**  $k \in \mathcal{S}^{(t,0)}$  **in parallel do:**
  - 7:     Set  $\mathbf{w}_k^{(t,0)} = \mathbf{w}^{(t,0)}$ , and calculate  $q_k(\mathbf{w}_k^{(t,0)}) = \sigma(F_k(\mathbf{w}_k^{(t,0)}) - F_k(\widehat{\mathbf{w}}_k))$
  - 8:     **For**  $r = 0, \dots, \tau - 1$  **local iterations do:**
  - 9:       Update  $\mathbf{w}_k^{(t,r+1)} \leftarrow \mathbf{w}_k^{(t,r)} - \eta_l \mathbf{g}(\mathbf{w}_k^{(t,r)}, \xi_k^{(t,r)})$
  - 10:      Send  $\Delta \mathbf{w}_k^{(t,0)} = \mathbf{w}_k^{(t,0)} - \mathbf{w}_k^{(t,\tau)}$  and aggregation weight  $q_k(\mathbf{w}_k^{(t,0)})$  to the server
  - 11:   **Global server do:**
  - 12:     Update global model with  $\mathbf{w}^{(t+1,0)} = \mathbf{w}^{(t,0)} - \eta_g^{(t,0)} \sum_{k \in \mathcal{S}^{(t,0)}} q_k(\mathbf{w}_k^{(t,0)}) \Delta \mathbf{w}_k^{(t,0)}$
- 

526 **Adaptive Server Learning Rate for INCFL.** With  $L_c$  continuous and  $L_s$  smooth  $F_k(\mathbf{w})$ ,  $\forall k \in [M]$   
 527 (see Assumption C.1), the objective  $\tilde{F}(\mathbf{w})$  is  $\tilde{L}_s$  smooth where  $\tilde{L}_s = \frac{L_s}{M} \sum_{k=1}^M q_k(\mathbf{w}) + \frac{L_c}{4}$  (see  
 528 Appendix D.2). Hence, the optimal learning rate  $\tilde{\eta}$  for the INCFL is given by,  $\tilde{\eta} = 1/\tilde{L}_s =$   
 529  $M\eta / (\sum_{k=1}^M q_k(\mathbf{w}) + \epsilon)$ , where  $\eta = \frac{1}{L_s}$  is the optimal learning rate for standard FL and  $\epsilon = \frac{ML_c}{4L_s}$   
 530  $> 0$  is a constant. The denominator of the optimal  $\tilde{\eta}$  is proportional to the sum of the aggregation  
 531 weights  $q_k(\mathbf{w})$  and acts as a dynamic normalizing factor. Therefore, we propose using adaptive global  
 532 learning rate  $\eta_g^{(t,0)} = \eta_g / (\sum_{k \in \mathcal{S}^{(t,0)}} q_k(\mathbf{w}_k^{(t,0)}) + \epsilon)$  with hyperparameters  $\eta_g$  and  $\epsilon$ .

533 **INCFL’s Theoretical Learning Rate Behavior for Fig. 4 (b).**  
 534 Here, we provide a plot of INCFL’s theoretical learning rate for the mean estimation example in Fig. 4(b) in Fig. 6 to show how  
 535 the learning rate changes for different regions of the model. We  
 536 show this plot as a proof of concept on the adaptive learning rate  
 537 we discuss above. For the sigmoid function which is used for our  
 538 INCFL objective, using a global notion of smoothness can cause  
 539 gradient descent to be too slow since global smoothness is deter-  
 540 mined by behavior at  $w = 0$  where  $w$  is the model. In this case, it is  
 541 better to use a local estimate of smoothness in the flat regions where  
 542  $|w| \gg 0$ . Recall that  $\nabla^2 \sigma(w) = \sigma(w)(1 - \sigma(w))(1 - 2\sigma(w)) <$   
 543  $\sigma(x)(1 - \sigma(w))$  and therefore setting the learning rate proportional  
 544 to  $\frac{1}{\sigma(w)(1 - \sigma(w))}$  can increase the learning rate in flat regions where  
 545  $\sigma(w)$  is close to 1 or 0. Following a similar argument, we can  
 546 show that the learning rate in our objective should be proportional  
 547 to  $1 / (\sum_{i=1}^M \sigma(F_i(w) - F_i(\hat{w}^*)) (1 - \sigma(F_i(w) - F_i(\hat{w}^*)))$ .  
 548

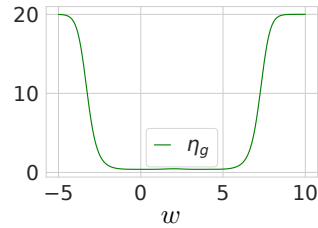


Figure 6: Behavior of Theoretical Learning Rate of INCFL for the mean estimation example in Fig. 4(b). As expected from the theoretical learning rate formula, we see a higher learning rate in regions where the function is flat.

549 **Ease of Implementing INCFL.** INCFL enjoys the following prop-  
 550 erties: i) it does not modify the local SGD procedure clients perform  
 551 in standard FL, ii) it allows for partial client participation, and iii) it is stateless. By stateless, we  
 552 mean that clients do not carry varying local parameters throughout training rounds preventing any  
 553 problems from stale parameters which can be exacerbated with partial client participation. Note that  
 554  $F_k(\widehat{\mathbf{w}}_k)$  needed for calculating  $q_k(\mathbf{w})$ ,  $k \in [M]$  can be considered as an input parameter for INCFL,  
 555 that is computed once and saved as a constant beforehand at each client by training  $\widehat{\mathbf{w}}_k$  on its local  
 556 dataset for a few SGD steps.

## 557 C.1 Convergence Properties of INCFL

558 In this section we show the convergence properties of the global model trained with INCFL. Our  
 559 convergence analysis shows that the gradient norm of our global model goes to zero and therefore we  
 560 converge to a stationary point of our objective  $\tilde{F}(\mathbf{w})$ .

561 First we introduce the assumptions and definitions utilized for our convergence analysis below.

562 **Assumption C.1** (Continuity & Smoothness of  $F_k(\mathbf{w})$ ,  $\forall k$ ). *The local objective functions for all*  
 563 *clients,  $F_1(\mathbf{w})$ , ...,  $F_M(\mathbf{w})$ , are all  $L_c$ -continuous and  $L_s$ -smooth for all  $\mathbf{w}$ .*

564 **Assumption C.2** (Unbiased Stochastic Gradient with Bounded Variance for  $F_k(\mathbf{w})$ ,  $\forall k$ ). *For the*  
 565 *mini-batch  $\xi_k$  uniformly sampled at random from  $\mathcal{B}_k$  from user  $k$ , the resulting stochastic gradient is*  
 566 *unbiased, i.e.,  $\mathbb{E}[\mathbf{g}_k(\mathbf{w}_k, \xi_k)] = \nabla F_k(\mathbf{w}_k)$ . Also, the variance of stochastic gradients is bounded:*  
 567  *$\mathbb{E}[\|\mathbf{g}_k(\mathbf{w}_k, \xi_k) - \nabla F_k(\mathbf{w}_k)\|^2] \leq \sigma_g^2$  for  $k \in [M]$ .*

568 **Assumption C.3** (Bounded Dissimilarity of  $F(\mathbf{w})$ ). *There exists constants  $\beta^2 \geq 1$ ,  $\kappa^2 \geq 0$  such*  
 569 *that  $\frac{1}{M} \sum_{i=1}^M \|\nabla F_i(\mathbf{w})\|^2 \leq \beta^2 \|\frac{1}{M} \sum_{i=1}^M \nabla F_i(\mathbf{w})\|^2 + \kappa^2$  for any  $\mathbf{w}$ .*

570 Assumption C.1-C.3 are standard assumptions frequently used in the optimization literature [36, 3, 37,  
 571 4], including the  $L_c$ -continuity assumption [38, 39]. Note that we do not have any assumptions over  
 572 our proposed objective function  $\tilde{F}(\mathbf{w})$  and only use the conventional assumptions used in FL for the  
 573 standard objective function  $F(\mathbf{w})$  to prove the convergence of INCFL over  $\tilde{F}(\mathbf{w})$  in Theorem C.1.

574 **Theorem C.1** (Convergence to the INCFL Objective  $\tilde{F}(\mathbf{w})$ ). *Under Assumption C.1-C.3, suppose the*  
 575 *server uniformly selects  $m$  out of  $M$  clients without replacement in each global round of Algorithm 1.*  
 576 *With  $\eta_l = \frac{1}{\sqrt{T}\tau}$ ,  $\eta_g = \sqrt{\tau m}$ , for a sufficiently large  $T$  our optimization error is bounded as follows:*

$$\min_{t \in [T]} \mathbb{E} \left[ \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 \right] \leq \mathcal{O} \left( \frac{\sigma_g^2}{\sqrt{m\tau T}} \right) + \mathcal{O} \left( \frac{\sigma_g^2}{T\tau} \right) + \mathcal{O} \left( \frac{\sqrt{\tau}}{\sqrt{Tm}} \right) + \mathcal{O} \left( \frac{\kappa^2 + \beta^2}{T} \right) \quad (50)$$

577 where  $\mathcal{O}$  subsumes all constants (including  $L_s$  and  $L_c$ ).

578 Theorem C.1 shows that with a sufficiently large number of communication rounds  $T$  we reach a  
 579 stationary point of our objective function  $\tilde{F}(\mathbf{w})$ . The proof is deferred to Appendix D.2 where we  
 580 also show a version of this theorem that contains the learning rates  $\eta_g$  and  $\eta_l$  with the constants.

## 581 D Convergence Proof

### 582 D.1 Preliminaries

583 First, we introduce the key lemmas used for the convergence analysis.

584 **Lemma D.1** (Bounded Dissimilarity for  $\tilde{F}(\mathbf{w})$ ). *With Assumption C.1 and Assumption C.3 we have*  
 585 *the bounded dissimilarity with respect to  $\tilde{F}(\mathbf{w})$  as:*

$$\frac{1}{M} \sum_{i=1}^M \|\nabla \tilde{F}_i(\mathbf{w})\|^2 \leq \beta'^2 \|\nabla \tilde{F}(\mathbf{w})\|^2 + \kappa'^2 \quad (51)$$

586 where  $\beta'^2 = 2\beta^2$ ,  $\kappa'^2 = 4\beta^2 L_c^2 + \kappa^2$

587 *Proof.* One can easily show that

$$\frac{1}{M} \sum_{i=1}^M \|\nabla \tilde{F}_i(\mathbf{w})\|^2 = \frac{1}{M} \sum_{i=1}^M q_i(\mathbf{w})^2 \|\nabla F_i(\mathbf{w})\|^2 \leq \frac{1}{M} \sum_{i=1}^M \|\nabla F_i(\mathbf{w})\|^2 \quad (52)$$

588 due to  $q_i(\mathbf{w}) \leq 1$ . Hence we have from Assumption C.3 and Cauchy-Schwarz inequality that

$$\frac{1}{M} \sum_{i=1}^M \|\nabla \tilde{F}_i(\mathbf{w})\|^2 \leq \frac{1}{M} \sum_{i=1}^M \|\nabla F_i(\mathbf{w})\|^2 \leq \beta^2 \|\nabla F(\mathbf{w}) - \nabla \tilde{F}(\mathbf{w}) + \nabla \tilde{F}(\mathbf{w})\|^2 + \kappa^2 \quad (53)$$

$$\leq 2\beta^2 \|\nabla F(\mathbf{w}) - \nabla \tilde{F}(\mathbf{w})\|^2 + 2\beta^2 \|\nabla \tilde{F}(\mathbf{w})\|^2 + \kappa^2 \quad (54)$$

589 We bound the first term in (54) as

$$\|\nabla F(\mathbf{w}) - \nabla \tilde{F}(\mathbf{w})\|^2 = \left\| \sum_{i=1}^M \frac{(1 - q_i(\mathbf{w}))}{M} \nabla F_i(\mathbf{w}) \right\|^2 \leq \frac{1}{M} \sum_{i=1}^M \|(1 - q_i(\mathbf{w})) \nabla F_i(\mathbf{w})\|^2 \quad (55)$$

$$\leq \frac{2}{M} \sum_{i=1}^M \|\nabla F_i(\mathbf{w})\|^2 \leq 2L_c^2 \quad (56)$$

590 where in (56) we use  $q_i(\mathbf{w}) \leq 1, \forall i \in [M]$  and Assumption C.1. Then from (54) we have

$$\frac{1}{M} \sum_{i=1}^M \|\nabla \tilde{F}_i(\mathbf{w})\|^2 \leq 2\beta^2 \|\nabla \tilde{F}(\mathbf{w})\|^2 + \kappa^2 + 4\beta^2 L_c^2 \quad (57)$$

591 completing the proof.  $\square$

592 **Lemma D.2** (Smoothness of  $\tilde{F}(\mathbf{w})$ ). *If Assumption C.1 is satisfied we have that the incentive local*  
 593 *objectives,  $\tilde{F}_1(\mathbf{w}), \dots, \tilde{F}_M(\mathbf{w})$ , are also  $\tilde{L}_s$ -smooth for any  $\mathbf{w}$  where  $\tilde{L}_s = L_c^2/4 + q_i(\mathbf{w})L_s$ .*

594 *Proof.* Recall the definitions of  $\tilde{F}(\mathbf{w})$  below:

$$\tilde{F}(\mathbf{w}) = \frac{1}{M} \sum_{i=1}^M \tilde{F}_i(\mathbf{w}), \quad \tilde{F}_i(\mathbf{w}) := \sigma(F_i(\mathbf{w}) - F_i(\hat{\mathbf{w}}_i^*)) \quad (58)$$

595 Let  $\|\cdot\|_{op}$  denote the spectral norm of a matrix. Accordingly, with the model parameter vector  $\mathbf{w} \in \mathbb{R}^d$ ,  
 596 we have the spectral norm of the Hessian of  $\tilde{F}_i(\mathbf{w})$ ,  $\forall i \in [M]$  as:

$$\|\nabla^2 \tilde{F}_i(\mathbf{w})\|_{op} = \|q_i(\mathbf{w})[(\nabla F_i(\mathbf{w}) \nabla F_i(\mathbf{w})^T)(1 - q_i(\mathbf{w})) + \nabla^2 F_i(\mathbf{w})]\|_{op} \quad (59)$$

597 where  $q_i(\mathbf{w}) = \text{Sigmoid}(F_i(\mathbf{w}) - F_i(\hat{\mathbf{w}}_i^*))$  and  $\nabla F_i(\mathbf{w}) \in \mathbb{R}^{d \times 1}$  is the gradient vector for the local  
 598 objective  $F_i(\mathbf{w})$  and  $\nabla^2 F_i(\mathbf{w}) \in \mathbb{R}^{d \times d}$  is the Hessian of  $F_i(\mathbf{w})$ . We can bound the RHS of (59) as  
 599 follows

$$\|\nabla^2 \tilde{F}_i(\mathbf{w})\|_{op} = \|q_i(\mathbf{w})(1 - q_i(\mathbf{w}))(\nabla F_i(\mathbf{w}) \nabla F_i(\mathbf{w})^T) + q_i(\mathbf{w}) \nabla^2 F_i(\mathbf{w})\|_{op} \quad (60)$$

$$\leq \|q_i(\mathbf{w})(1 - q_i(\mathbf{w}))(\nabla F_i(\mathbf{w}) \nabla F_i(\mathbf{w})^T)\|_{op} + \|q_i(\mathbf{w}) \nabla^2 F_i(\mathbf{w})\|_{op} \quad (61)$$

$$= q_i(\mathbf{w})(1 - q_i(\mathbf{w})) \|(\nabla F_i(\mathbf{w}) \nabla F_i(\mathbf{w})^T)\|_{op} + q_i(\mathbf{w}) \|\nabla^2 F_i(\mathbf{w})\|_{op} \quad (62)$$

$$= q_i(\mathbf{w})(1 - q_i(\mathbf{w})) \|\nabla F_i(\mathbf{w})\|^2 + q_i(\mathbf{w}) \|\nabla^2 F_i(\mathbf{w})\|_{op} \quad (63)$$

$$\leq \frac{L_c^2}{4} + q_i(\mathbf{w})L_s \quad (64)$$

600 where we use triangle inequality in (61), and use  $\|\mathbf{x}\mathbf{y}^T\|_{op} = \|\mathbf{x}\| \|\mathbf{y}\|$  in (63), and use  $q_i(\mathbf{w}) \leq$   
 601  $1$  along with Assumption C.1 in (64). Since the norm of the Hessian of  $\tilde{F}_i(\mathbf{w})$  is bounded by  
 602  $\frac{L_c^2}{4} + q_i(\mathbf{w})L_s$  we complete the proof.  $\square$

## 603 D.2 Proof of Theorem C.1 – Full Client Participation

604 For ease of writing, we define the following auxiliary variables for any client  $i \in [M]$ :

$$\text{Weighted Stochastic Gradient: } \mathbf{h}_i^{(t,0)} := q_i(\mathbf{w}^{(t,0)}) \sum_{r=0}^{\tau-1} \mathbf{g}(\mathbf{w}_i^{(t,r)}, \xi_i^{(t,r)}), \quad (65)$$

$$\text{Weighted Gradient: } \bar{\mathbf{h}}_i^{(t,0)} := q_i(\mathbf{w}^{(t,0)}) \sum_{r=0}^{\tau-1} \nabla F_i(\mathbf{w}_i^{(t,r)}), \quad (66)$$

$$\text{Normalized Global Learning Rate: } \eta_g^{(t,0)} := \eta_g / \left( \sum_{i=1}^M q_i(\mathbf{w}^{(t,0)}) + \epsilon \right) \quad (67)$$

605 where  $\epsilon$  is a constant added to the denominator to prevent the denominator from being 0. From  
 606 Algorithm 1 with full client participation, our proposed algorithm has the following effective update  
 607 rule for the global model at the server:

$$\mathbf{w}^{(t+1,0)} = \mathbf{w}^{(t,0)} - \eta_g^{(t,0)} \eta_l \sum_{k=1}^M \mathbf{h}_k^{(t,0)} \quad (68)$$

608 With the update rule in (68), defining  $\tilde{\eta}^{(t,0)} := \eta_g^{(t,0)} \eta_l \tau M$  and using Lemma D.2 we have

$$\begin{aligned} \mathbb{E} \left[ \tilde{F}(\mathbf{w}^{(t+1,0)}) \right] - \tilde{F}(\mathbf{w}^{(t,0)}) &\leq -\tilde{\eta}^{(t,0)} \mathbb{E} \left[ \left\langle \nabla \tilde{F}(\mathbf{w}^{(t,0)}), \frac{1}{M\tau} \sum_{i=1}^M \mathbf{h}_i^{(t,0)} \right\rangle \right] \\ &\quad + \frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2}{2} \mathbb{E} \left[ \left\| \frac{1}{M\tau} \sum_{i=1}^M \mathbf{h}_i^{(t,0)} \right\|^2 \right] \end{aligned} \quad (69)$$

$$\begin{aligned} &= -\tilde{\eta}^{(t,0)} \mathbb{E} \left[ \left\langle \nabla \tilde{F}(\mathbf{w}^{(t,0)}), \frac{1}{M\tau} \sum_{i=1}^M (\mathbf{h}_i^{(t,0)} - \bar{\mathbf{h}}_i^{(t,0)}) \right\rangle \right] - \tilde{\eta}^{(t,0)} \mathbb{E} \left[ \left\langle \nabla \tilde{F}(\mathbf{w}^{(t,0)}), \frac{1}{M\tau} \sum_{i=1}^M \bar{\mathbf{h}}_i^{(t,0)} \right\rangle \right] \\ &\quad + \frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2}{2} \mathbb{E} \left[ \left\| \frac{1}{M\tau} \sum_{i=1}^M \mathbf{h}_i^{(t,0)} \right\|^2 \right] \end{aligned} \quad (70)$$

$$\begin{aligned} &= -\frac{\tilde{\eta}^{(t,0)}}{2} \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 - \frac{\tilde{\eta}^{(t,0)}}{2} \mathbb{E} \left[ \left\| \frac{1}{M\tau} \sum_{i=1}^M \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] + \frac{\tilde{\eta}^{(t,0)}}{2} \mathbb{E} \left[ \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) - \frac{1}{M\tau} \sum_{i=1}^M \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \\ &\quad + \frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2}{2M^2\tau^2} \mathbb{E} \left[ \left\| \sum_{i=1}^M \mathbf{h}_i^{(t,0)} \right\|^2 \right] \end{aligned} \quad (71)$$

609 For the last term in (71), we can bound it as

$$\frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2}{2M^2\tau^2} \mathbb{E} \left[ \left\| \sum_{i=1}^M \mathbf{h}_i^{(t,0)} \right\|^2 \right] \leq \frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2}{M^2\tau^2} \sum_{i=1}^M \mathbb{E} \left[ \left\| \mathbf{h}_i^{(t,0)} - \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] + \frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2}{M^2\tau^2} \mathbb{E} \left[ \left\| \sum_{i=1}^M \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \quad (72)$$

$$= \frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2}{M^2\tau^2} \sum_{i=1}^M \mathbb{E} \left[ \left\| q_i(\mathbf{w}^{(t,0)}) \sum_{r=0}^{\tau-1} (\mathbf{g}(\mathbf{w}_i^{(t,r)}, \xi_i^{(t,r)}) - \nabla F_i(\mathbf{w}_i^{(t,r)})) \right\|^2 \right] + \frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2}{M^2\tau^2} \mathbb{E} \left[ \left\| \sum_{i=1}^M \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \quad (73)$$

$$= \frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2}{M^2\tau^2} \sum_{i=1}^M q_i(\mathbf{w}^{(t,0)})^2 \sum_{r=0}^{\tau-1} \mathbb{E} \left[ \left\| \mathbf{g}(\mathbf{w}_i^{(t,r)}, \xi_i^{(t,r)}) - \nabla F_i(\mathbf{w}_i^{(t,r)}) \right\|^2 \right] + \frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2}{M^2\tau^2} \mathbb{E} \left[ \left\| \sum_{i=1}^M \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \quad (74)$$

$$= \frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2}{M^2\tau^2} \sum_{i=1}^M q_i(\mathbf{w}^{(t,0)})^2 \tau \sigma_g^2 + \frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2}{M^2\tau^2} \mathbb{E} \left[ \left\| \sum_{i=1}^M \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \quad (75)$$

$$\leq \frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2 \sigma_g^2}{M\tau} + \tilde{L}_s(\tilde{\eta}^{(t,0)})^2 \mathbb{E} \left[ \left\| \frac{1}{M\tau} \sum_{i=1}^M \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \quad (76)$$



610 where (72) is due to the Cauchy-Schwartz inequality and (75) is due to Assumption C.2 and (76) is  
 611 due to  $q_i(\mathbf{w}) \leq 1, \forall i \in [M]$ . Merging (76) into (71) we have

$$\begin{aligned} \mathbb{E} \left[ \tilde{F}(\mathbf{w}^{(t+1,0)}) \right] - \tilde{F}(\mathbf{w}^{(t,0)}) &\leq -\frac{\tilde{\eta}^{(t,0)}}{2} \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 + \frac{\tilde{\eta}^{(t,0)}}{2} \mathbb{E} \left[ \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) - \frac{1}{M\tau} \sum_{i=1}^M \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \\ &\quad + \frac{\tilde{L}_s (\tilde{\eta}^{(t,0)})^2 \sigma_g^2}{M\tau} + \left( (\tilde{\eta}^{(t,0)})^2 \tilde{L}_s - \frac{\tilde{\eta}^{(t,0)}}{2} \right) \mathbb{E} \left[ \left\| \frac{1}{M\tau} \sum_{i=1}^M \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \end{aligned} \quad (77)$$

612 Now we aim at bounding the second term in the RHS of (77) as follows:

$$\frac{\tilde{\eta}^{(t,0)}}{2} \mathbb{E} \left[ \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) - \frac{1}{M\tau} \sum_{i=1}^M \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \quad (78)$$

$$= \frac{\tilde{\eta}^{(t,0)}}{2} \mathbb{E} \left[ \left\| \frac{1}{M} \sum_{i=1}^M q_i(\mathbf{w}^{(t,0)}) \nabla F_i(\mathbf{w}^{(t,0)}) - \frac{1}{M\tau} \sum_{i=1}^M q_i(\mathbf{w}^{(t,0)}) \sum_{r=0}^{\tau-1} \nabla F_i(\mathbf{w}_i^{(t,r)}) \right\|^2 \right] \quad (79)$$

$$= \frac{\tilde{\eta}^{(t,0)}}{2} \mathbb{E} \left[ \left\| \frac{1}{M\tau} \sum_{i=1}^M q_i(\mathbf{w}^{(t,0)}) \sum_{r=0}^{\tau-1} \left( \nabla F_i(\mathbf{w}^{(t,0)}) - \nabla F_i(\mathbf{w}_i^{(t,r)}) \right) \right\|^2 \right] \quad (80)$$

$$\leq \frac{\tilde{\eta}^{(t,0)}}{2M\tau} \sum_{i=1}^M q_i(\mathbf{w}^{(t,0)})^2 \sum_{r=0}^{\tau-1} \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}^{(t,0)}) - \nabla F_i(\mathbf{w}_i^{(t,r)}) \right\|^2 \right] \quad (81)$$

$$= \frac{L_s^2 \tilde{\eta}^{(t,0)}}{2M\tau} \sum_{i=1}^M q_i(\mathbf{w}^{(t,0)})^2 \sum_{r=0}^{\tau-1} \mathbb{E} \left[ \left\| \mathbf{w}^{(t,0)} - \mathbf{w}_i^{(t,r)} \right\|^2 \right] \quad (82)$$

613 where (81) is due to Jensen's inequality and (82) is due to Lemma D.2. We can bound the difference  
 614 of the global model and local model for any client  $i \in [M]$  as follows:

$$\mathbb{E} \left[ \left\| \mathbf{w}^{(t,0)} - \mathbf{w}_i^{(t,r)} \right\|^2 \right] = \eta_l^2 \mathbb{E} \left[ \left\| \sum_{l=0}^{r-1} \mathbf{g}(\mathbf{w}_i^{(t,l)}, \xi_i^{(t,l)}) \right\|^2 \right] \quad (83)$$

$$\leq 2\eta_l^2 \mathbb{E} \left[ \left\| \sum_{l=0}^{r-1} \mathbf{g}(\mathbf{w}_i^{(t,l)}, \xi_i^{(t,l)}) - \nabla F_i(\mathbf{w}_i^{(t,l)}) \right\|^2 \right] + 2\eta_l^2 \mathbb{E} \left[ \left\| \sum_{l=0}^{r-1} \nabla F_i(\mathbf{w}_i^{(t,l)}) \right\|^2 \right] \quad (84)$$

$$\leq 2\eta_l^2 \sigma_g^2 r + 2\eta_l^2 \mathbb{E} \left[ \left\| \sum_{l=0}^{r-1} \nabla F_i(\mathbf{w}_i^{(t,l)}) \right\|^2 \right] \quad (85)$$

615 where (84) is due to Cauchy-Schwarz inequality and (85) is due to Assumption C.2. We bound the  
 616 last term in (85) as follows:

$$\mathbb{E} \left[ \left\| \sum_{l=0}^{r-1} \nabla F_i(\mathbf{w}_i^{(t,l)}) \right\|^2 \right] \leq r \sum_{l=0}^{r-1} \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}_i^{(t,l)}) \right\|^2 \right] \leq \tau \sum_{l=0}^{\tau-1} \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}_i^{(t,l)}) \right\|^2 \right] \quad (86)$$

$$\leq 2\tau \sum_{l=0}^{\tau-1} \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}_i^{(t,l)}) - \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \right] + 2\tau^2 \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \right] \quad (87)$$

$$\leq 2L_s^2 \tau \sum_{l=0}^{\tau-1} \mathbb{E} \left[ \left\| \mathbf{w}_i^{(t,l)} - \mathbf{w}^{(t,0)} \right\|^2 \right] + 2\tau^2 \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \right] \quad (88)$$

617 where (86) is due to Jensen's inequality, and (87) is due to Cauchy-Schwarz inequality, and (88) is  
618 due to Lemma D.2. Combining (88) with (85) we have that

$$\mathbb{E} \left[ \left\| \mathbf{w}^{(t,0)} - \mathbf{w}_i^{(t,r)} \right\|^2 \right] \leq 2\eta_l^2 \sigma_g^2 r + 4L_s^2 \eta_l^2 \tau \sum_{l=0}^{\tau-1} \mathbb{E} \left[ \left\| \mathbf{w}^{(t,0)} - \mathbf{w}_i^{(t,l)} \right\|^2 \right] + 4\eta_l^2 \tau^2 \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \right] \quad (89)$$

619 Reorganizing (89) and taking the summation  $r \in [\tau]$  on both sides we have,

$$(1 - 4L_s^2 \eta_l^2 \tau^2) \sum_{r=0}^{\tau-1} \mathbb{E} \left[ \left\| \mathbf{w}^{(t,0)} - \mathbf{w}_i^{(t,r)} \right\|^2 \right] \leq 2\eta_l^2 \sigma_g^2 \sum_{r=0}^{\tau-1} r + 4\eta_l^2 \tau^3 \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \right] \quad (90)$$

$$\leq \eta_l^2 \sigma_g^2 \tau^2 + 4\eta_l^2 \tau^3 \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \right] \quad (91)$$

620 With  $\eta_l \leq 1/(2\sqrt{2}\tau L_s)$ , we have that  $1/(1 - 4L_s^2 \eta_l^2 \tau^2) \leq 2$  and hence can further bound (91) as

$$\sum_{r=0}^{\tau-1} \mathbb{E} \left[ \left\| \mathbf{w}^{(t,0)} - \mathbf{w}_i^{(t,r)} \right\|^2 \right] \leq 2\eta_l^2 \sigma_g^2 \tau^2 + 8\eta_l^2 \tau^3 \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \right] \quad (92)$$

621 Finally, plugging in (92) to (82) we have

$$\frac{\tilde{\eta}^{(t,0)}}{2} \mathbb{E} \left[ \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) - \frac{1}{M\tau} \sum_{i=1}^M \tilde{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \quad (93)$$

$$\leq \frac{L_s^2 \tilde{\eta}^{(t,0)}}{2M\tau} \sum_{i=1}^M q_i(\mathbf{w}^{(t,0)})^2 \left( 2\eta_l^2 \sigma_g^2 \tau^2 + 8\eta_l^2 \tau^3 \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \right] \right)$$

$$\leq L_s^2 \tilde{\eta}^{(t,0)} \eta_l^2 \sigma_g^2 \tau + 4\eta_l^2 \tau^2 L_s^2 \tilde{\eta}^{(t,0)} \frac{1}{M} \sum_{i=1}^M \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \right] \quad (94)$$

$$\leq L_s^2 \tilde{\eta}^{(t,0)} \eta_l^2 \sigma_g^2 \tau + 4\eta_l^2 \tau^2 L_s^2 \tilde{\eta}^{(t,0)} (\beta'^2 \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 + \kappa'^2) \quad (95)$$

622 where (94) uses  $q_i(\mathbf{w}) \leq 1, \forall i \in [M]$  and (95) uses Lemma D.1. Merging (95) to (77) we have

$$\mathbb{E} \left[ \tilde{F}(\mathbf{w}^{(t+1,0)}) \right] - \tilde{F}(\mathbf{w}^{(t,0)})$$

$$\leq -\frac{\tilde{\eta}^{(t,0)}}{2} \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 + \tilde{\eta}^{(t,0)} \left( \tilde{\eta}^{(t,0)} \tilde{L}_s - \frac{1}{2} \right) \mathbb{E} \left[ \left\| \frac{1}{M\tau} \sum_{i=1}^M \tilde{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \quad (96)$$

$$+ \frac{\tilde{L}_s (\tilde{\eta}^{(t,0)})^2 \sigma_g^2}{M\tau} + \tilde{\eta}^{(t,0)} L_s^2 \eta_l^2 \sigma_g^2 \tau + 4\tilde{\eta}^{(t,0)} \eta_l^2 \tau^2 L_s^2 \beta'^2 \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 + 4\tilde{\eta}^{(t,0)} \eta_l^2 \tau^2 L_s^2 \kappa'^2$$

623 With  $\eta_l \eta_g \leq 1/(4\tau L_s)$  we have that  $\tilde{\eta}^{(t,0)} \tilde{L}_s - \frac{1}{2} \leq -1/4$  and thus can further simplify (96) to

$$\mathbb{E} \left[ \tilde{F}(\mathbf{w}^{(t+1,0)}) \right] - \tilde{F}(\mathbf{w}^{(t,0)}) \leq -\frac{\tilde{\eta}^{(t,0)}}{2} \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 + 4\tilde{\eta}^{(t,0)} \eta_l^2 \tau^2 L_s^2 \beta'^2 \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2$$

$$+ \frac{\tilde{L}_s (\tilde{\eta}^{(t,0)})^2 \sigma_g^2}{M\tau} + \tilde{\eta}^{(t,0)} L_s^2 \eta_l^2 \sigma_g^2 \tau + 4\tilde{\eta}^{(t,0)} \eta_l^2 \tau^2 L_s^2 \kappa'^2 \quad (97)$$

$$= \tilde{\eta}^{(t,0)} \left( 4\eta_l^2 \tau^2 L_s^2 \beta' - \frac{1}{2} \right) \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 + \frac{\tilde{L}_s (\tilde{\eta}^{(t,0)})^2 \sigma_g^2}{M\tau} + \tilde{\eta}^{(t,0)} L_s^2 \eta_l^2 \sigma_g^2 \tau + 4\tilde{\eta}^{(t,0)} \eta_l^2 \tau^2 L_s^2 \kappa'^2 \quad (98)$$

624 With local learning rate  $\eta_l \leq \min\{1/(4\tau L_s), 1/(4\beta'\tau L_s)\}$  we have that

$$\mathbb{E} \left[ \tilde{F}(\mathbf{w}^{(t+1,0)}) \right] - \tilde{F}(\mathbf{w}^{(t,0)}) \leq -\frac{\tilde{\eta}^{(t,0)}}{4} \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 + \frac{\tilde{L}_s (\tilde{\eta}^{(t,0)})^2 \sigma_g^2}{M\tau} + \tilde{\eta}^{(t,0)} L_s^2 \eta_l^2 \sigma_g^2 \tau$$

$$+ 4\tilde{\eta}^{(t,0)} \eta_l^2 \tau^2 L_s^2 \kappa'^2 \quad (99)$$

625 and we use the property of  $\tilde{\eta}^{(t,0)}$  that  $\frac{M\tau\eta_l\eta_g}{M+\epsilon} \leq \tilde{\eta}^{(t,0)} \leq \frac{M\tau\eta_l\eta_g}{\epsilon}$  to get

$$\begin{aligned} \mathbb{E} \left[ \tilde{F}(\mathbf{w}^{(t+1,0)}) \right] - \tilde{F}(\mathbf{w}^{(t,0)}) &\leq -\frac{M\tau\eta_l\eta_g}{4(M+\epsilon)} \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 + \frac{\tilde{L}_s M \tau \eta_l^2 \eta_g^2 \sigma_g^2}{\epsilon^2} \\ &\quad + \frac{M\tau^2 L_s^2 \eta_l^3 \eta_g \sigma_g^2}{\epsilon} + \frac{4M\eta_l^3 \eta_g \tau^3 L_s^2 \kappa'^2}{\epsilon} \end{aligned} \quad (100)$$

626 Taking the average across all rounds on both sides of (100) we get

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 \right] &\leq \frac{4(M+\epsilon) \left( \tilde{F}(\mathbf{w}^{(0,0)}) - \tilde{F}_{\text{inf}} \right)}{M\tau\eta_l\eta_g T} + \frac{16\eta_l^2 \tau^2 L_s^2 \kappa'^2 (M+\epsilon)}{\epsilon} \\ &\quad + \frac{4L_s^2 \eta_l^2 \tau \sigma_g^2 (M+\epsilon)}{\epsilon} + \frac{4\eta_g \eta_l \tilde{L}_s \sigma_g^2 (M+\epsilon)}{\epsilon^2} \end{aligned} \quad (101)$$

627 and prove

$$\begin{aligned} \min_{t \in [T]} \mathbb{E} \left[ \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 \right] &\leq \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 \right] \leq \frac{4(M+\epsilon) \left( \tilde{F}(\mathbf{w}^{(0,0)}) - \tilde{F}_{\text{inf}} \right)}{M\tau\eta_l\eta_g T} \\ &\quad + \frac{16\eta_l^2 \tau^2 L_s^2 \kappa'^2 (M+\epsilon)}{\epsilon} + \frac{4L_s^2 \eta_l^2 \tau \sigma_g^2 (M+\epsilon)}{\epsilon} + \frac{4\eta_g \eta_l \tilde{L}_s \sigma_g^2 (M+\epsilon)}{\epsilon^2} \end{aligned} \quad (102)$$

628 Further, using  $\tilde{L}_s = \frac{L_s}{M} \sum_{k=1}^M q_k(\mathbf{w}) + \frac{L_c}{4}$  and  $\epsilon = \frac{ML_c}{4L_s} > 0$  from the optimal learning rate we  
629 have the bound in (102) to be

$$\begin{aligned} \min_{t \in [T]} \mathbb{E} \left[ \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 \right] &\leq \frac{(4L_s + L_c) \left( \tilde{F}(\mathbf{w}^{(0,0)}) - \tilde{F}_{\text{inf}} \right)}{L_s \tau \eta_l \eta_g T} + \frac{64\eta_l^2 \tau^2 L_s^2 \kappa'^2 (4L_s + L_c)}{L_c} \\ &\quad + \frac{4L_s^2 \eta_l^2 \tau \sigma_g^2 (4L_s + L_c)}{L_c} + \frac{64L_s \eta_g \eta_l \sigma_g^2 (L_s + L_c/4)^2}{ML_c^2} \end{aligned} \quad (103)$$

630 By setting the global and local learning rate as  $\eta_g = \sqrt{\tau M}$  and  $\eta_l = \frac{1}{\sqrt{T\tau}}$  we can further optimize  
631 the bound as

$$\begin{aligned} \min_{t \in [T]} \mathbb{E} \left[ \left\| \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 \right] &\leq \frac{(4L_s + L_c) \left( \tilde{F}(\mathbf{w}^{(0,0)}) - \tilde{F}_{\text{inf}} \right)}{L_s \sqrt{T M \tau}} + \frac{64L_s^2 \kappa'^2 (4L_s + L_c)}{L_c T} \\ &\quad + \frac{4L_s^2 \sigma_g^2 (4L_s + L_c)}{T \tau L_c} + \frac{64L_s \sigma_g^2 (L_s + L_c/4)^2}{\sqrt{T M \tau}} \end{aligned} \quad (104)$$

632 completing the full client participation proof of Theorem C.1.

### 633 D.3 Proof of Theorem C.1 – Partial Client Participation

634 We present the convergence guarantees of INCFL for partial client participation in this section. With  
635 partial client participation, we have the update rule in (68) changed to

$$\mathbf{w}^{(t+1,0)} = \mathbf{w}^{(t,0)} - \eta_g^{(t,0)} \eta_l \sum_{k \in \mathcal{S}^{(t,0)}} \mathbf{h}_k^{(t,0)} \quad (105)$$

636 where the  $m$  clients are sampled uniformly at random without replacement for  $\mathcal{S}^{(t,0)}$  at each commu-  
637 nication round  $t$  by the server and  $\eta_g^{(t,0)} = m\eta_g / (\sum_{k \in \mathcal{S}^{(t,0)}} q_k(\mathbf{w}^{(t,0)}) + \epsilon)$  for positive constant  $\epsilon$ .

638 Then with the update rule in (105) and Lemma D.2, defining  $\tilde{\eta}^{(t,0)} = \eta_g^{(t,0)} \eta_l \tau m$  we have

$$\begin{aligned} \mathbb{E} \left[ \tilde{F}(\mathbf{w}^{(t+1,0)}) \right] - \tilde{F}(\mathbf{w}^{(t,0)}) &\leq \mathbb{E} \left[ -\tilde{\eta}^{(t,0)} \left\langle \nabla \tilde{F}(\mathbf{w}^{(t,0)}), \frac{1}{m\tau} \sum_{i \in \mathcal{S}^{(t,0)}} \mathbf{h}_i^{(t,0)} \right\rangle \right] \\ &\quad + \mathbb{E} \left[ \frac{\tilde{L}_s (\tilde{\eta}^{(t,0)})^2}{2} \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}^{(t,0)}} \mathbf{h}_i^{(t,0)} \right\|^2 \right] \end{aligned} \quad (106)$$

639 For the first term in the RHS of (106) we have that due to the uniform sampling of clients (see Lemma  
640 4 in [40]), it becomes analogous to the derivation for full client participation. Hence, with the property  
641 of  $\frac{m\tau\eta_l\eta_g}{m+\epsilon} \leq \tilde{\eta}^{(t,0)} \leq \frac{m\tau\eta_l\eta_g}{\epsilon}$  and using the previous bounds in (95), we result in the final bound for  
642 the first term in the RHS of (106) as below:

$$\mathbb{E} \left[ -\tilde{\eta}^{(t,0)} \left\langle \nabla \tilde{F}(\mathbf{w}^{(t,0)}), \frac{1}{m\tau} \sum_{i \in \mathcal{S}^{(t,0)}} \mathbf{h}_i^{(t,0)} \right\rangle \right] \leq \left( -\frac{m\tau\eta_l\eta_g}{m+\epsilon} + \frac{4\eta_l^3\tau^3 L_s^2 \beta'^2 \eta_g m}{\epsilon} \right) \|\nabla \tilde{F}(\mathbf{w}^{(t,0)})\|^2 + \frac{4L_s^2\tau^3\eta_l^3 m\eta_g \kappa'^2}{\epsilon} + \frac{L_s^2\tau^2\eta_l^2 m\eta_g \sigma_g^2}{\epsilon} \quad (107)$$

643 For the second term in the RHS of (106), with  $C = \tilde{L}_s(m\tau\eta_l\eta_g/\epsilon)^2$  we have the following:

$$\mathbb{E} \left[ \frac{\tilde{L}_s(\tilde{\eta}^{(t,0)})^2}{2} \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}^{(t,0)}} \mathbf{h}_i^{(t,0)} \right\|^2 \right] \leq C \mathbb{E} \left[ \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}^{(t,0)}} (\mathbf{h}_i^{(t,0)} - \bar{\mathbf{h}}_i^{(t,0)}) \right\|^2 \right] + C \mathbb{E} \left[ \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}^{(t,0)}} \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \quad (108)$$

$$= \frac{C}{m^2\tau^2} \mathbb{E} \left[ \sum_{i \in \mathcal{S}^{(t,0)}} \|\mathbf{h}_i^{(t,0)} - \bar{\mathbf{h}}_i^{(t,0)}\|^2 \right] + C \mathbb{E} \left[ \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}^{(t,0)}} \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \quad (109)$$

$$= \frac{C}{mM\tau^2} \sum_{i=1}^M \mathbb{E} \left[ \|\mathbf{h}_i^{(t,0)} - \bar{\mathbf{h}}_i^{(t,0)}\|^2 \right] + C \mathbb{E} \left[ \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}^{(t,0)}} \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \quad (110)$$

$$\leq \frac{C\sigma_g^2}{m\tau} + C \mathbb{E} \left[ \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}^{(t,0)}} \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \quad (111)$$

644 where (110) follows due to, again, the uniform sampling of clients and the rest follows identical steps  
645 for full client participation in the derivation for (72). Note that

$$C = \left( \frac{L_s}{M} \sum_{k=1}^M q_k(\mathbf{w}) + \frac{L_c}{4} \right) (m\tau\eta_l\eta_g/\epsilon)^2 \leq (L_s + \frac{L_c}{4})(m\tau\eta_l\eta_g/\epsilon)^2 \quad (112)$$

646 For the second term in (111) we have that

$$\mathbb{E} \left[ \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}^{(t,0)}} \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] = \mathbb{E} \left[ \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}^{(t,0)}} \left( \bar{\mathbf{h}}_i^{(t,0)} - \nabla \tilde{F}_i(\mathbf{w}^{(t,0)}) + \nabla \tilde{F}_i(\mathbf{w}^{(t,0)}) \right) \right\|^2 \right] \quad (113)$$

$$\leq \underbrace{3 \mathbb{E} \left[ \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}^{(t,0)}} \left( \bar{\mathbf{h}}_i^{(t,0)} - \nabla \tilde{F}_i(\mathbf{w}^{(t,0)}) \right) \right\|^2 \right]}_{A_1} + \underbrace{\frac{3}{\tau^2} \mathbb{E} \left[ \left\| \frac{1}{m} \sum_{i \in \mathcal{S}^{(t,0)}} \nabla \tilde{F}_i(\mathbf{w}^{(t,0)}) - \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 \right]}_{A_2} + 3 \mathbb{E} \left[ \left\| \frac{1}{\tau} \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 \right] \quad (114)$$

647 First we bound  $A_1$  in (114) as follows:

$$3\mathbb{E} \left[ \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}(t,0)} \left( \bar{\mathbf{h}}_i^{(t,0)} - \nabla \tilde{F}_i(\mathbf{w}^{(t,0)}) \right) \right\|^2 \right] \quad (115)$$

$$= 3\mathbb{E} \left[ \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}(t,0)} q_i(\mathbf{w}^{(t,0)}) \sum_{r=0}^{\tau-1} \left( \nabla F_i(\mathbf{w}_i^{(t,r)}) - \nabla F_i(\mathbf{w}^{(t,0)}) \right) \right\|^2 \right] \quad (116)$$

$$\leq \frac{3}{m\tau} \mathbb{E} \left[ \sum_{i \in \mathcal{S}(t,0)} \sum_{r=0}^{\tau-1} \left\| \nabla F_i(\mathbf{w}_i^{(t,r)}) - \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \right] \quad (117)$$

$$= \frac{3}{M\tau} \sum_{i=1}^M \sum_{r=0}^{\tau-1} \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}_i^{(t,r)}) - \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \right] \quad (118)$$

$$\leq \frac{3L_s^2}{M\tau} \sum_{i=1}^M \sum_{r=0}^{\tau-1} \mathbb{E} \left[ \left\| \mathbf{w}^{(t,0)} - \mathbf{w}_i^{(t,r)} \right\|^2 \right] \quad (118)$$

648 where (116) is due to Jensen's inequality and  $q_i(\mathbf{w}) \leq 1$  and (117) is due to the uniform sampling of  
 649 clients, and (118) is due to Assumption C.1. Using (77) we have already derived, bound (118) further  
 650 to:

$$3\mathbb{E} \left[ \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}(t,0)} \left( \bar{\mathbf{h}}_i^{(t,0)} - \nabla \tilde{F}_i(\mathbf{w}^{(t,0)}) \right) \right\|^2 \right] \leq 6L_s^2\eta_l^2\sigma_g^2\tau + \frac{24L_s^2\eta_l^2\tau^2}{M} \sum_{i=1}^M \mathbb{E} \left[ \left\| \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \right] \quad (119)$$

$$\leq 6L_s^2\eta_l^2\sigma_g^2\tau + 24L_s^2\eta_l^2\tau^2(\beta'^2\|\nabla \tilde{F}(\mathbf{w}^{(t,0)})\|^2 + \kappa'^2) \quad (120)$$

651 where (120) is due to Lemma D.1.

652 Next we bound  $A_2$  as follows:

$$\frac{3}{\tau^2} \mathbb{E} \left[ \left\| \frac{1}{m} \sum_{i \in \mathcal{S}(t,0)} \nabla \tilde{F}_i(\mathbf{w}^{(t,0)}) - \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 \right] \quad (121)$$

$$= \frac{3(M-m)}{\tau^2 m M (M-1)} \sum_{i=1}^M \mathbb{E} \left[ \left\| \nabla \tilde{F}_i(\mathbf{w}^{(t,0)}) - \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 \right]$$

$$= \frac{3(M-m)}{\tau^2 m M (M-1)} \sum_{i=1}^M \left\| \nabla q_i(\mathbf{w}^{(t,0)}) F_i(\mathbf{w}^{(t,0)}) - \frac{1}{M} \sum_{i=1}^M q_i(\mathbf{w}^{(t,0)}) \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \quad (122)$$

$$\leq \frac{6(M-m)}{\tau^2 m M (M-1)} \sum_{i=1}^M \left( \left\| \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 + \left\| \frac{1}{M} \sum_{i=1}^M q_i(\mathbf{w}^{(t,0)}) \nabla F_i(\mathbf{w}^{(t,0)}) \right\|^2 \right) \quad (123)$$

$$\leq \frac{12(M-m)L_c^2}{\tau^2 m (M-1)} \quad (124)$$

653 where (121) is due to the variance under uniform sampling without replacement (see Lemma 4 in  
 654 [40]) and (123) is due to the Cauchy-Schwarz inequality and (124) is due to Assumption C.1.

655 Mering the bounds for  $A_1$  and  $A_2$  to (114) we have that

$$\mathbb{E} \left[ \left\| \frac{1}{m\tau} \sum_{i \in \mathcal{S}(t,0)} \bar{\mathbf{h}}_i^{(t,0)} \right\|^2 \right] \leq 6L_s^2 \eta_l^2 \sigma_g^2 \tau + 24L_s^2 \eta_l^2 \tau^2 \beta'^2 \|\nabla \tilde{F}(\mathbf{w}^{(t,0)})\|^2 \quad (125)$$

$$+ 24L_s^2 \eta_l^2 \tau^2 \kappa'^2 + \frac{12(M-m)L_c^2}{\tau^2 m(M-1)} + 3\mathbb{E} \left[ \left\| \frac{1}{\tau} \nabla \tilde{F}(\mathbf{w}^{(t,0)}) \right\|^2 \right] \\ = \left( 24L_s^2 \eta_l^2 \tau^2 \beta'^2 + \frac{3}{\tau^2} \right) \|\nabla \tilde{F}(\mathbf{w}^{(t,0)})\|^2 + 6L_s^2 \eta_l^2 \tau (\sigma_g^2 + 4\tau \kappa'^2) + \frac{12(M-m)L_c^2}{\tau^2 m(M-1)} \quad (126)$$

656 Then we can plug in (126) back to (111) and plugging in (107) to (106), we can derive the bound in  
657 (106) as

$$\mathbb{E} \left[ \tilde{F}(\mathbf{w}^{(t+1,0)}) \right] - \tilde{F}(\mathbf{w}^{(t,0)}) \\ \leq \left( -\frac{m\tau\eta_l\eta_g}{m+\epsilon} + \frac{4\eta_l^3\eta_g\tau^3 L_s^2 \beta'^2 m}{\epsilon} + \nu (\tau\eta_l\eta_g)^2 (24L_s^2 \eta_l^2 \tau^2 \beta'^2 + 3) \right) \|\nabla \tilde{F}(\mathbf{w}^{(t,0)})\|^2 + (\tau\eta_l\eta_g)^2 \nu \frac{\sigma_g^2}{m\tau} \\ + (\tau\eta_l\eta_g)^2 \nu \left( 6L_s^2 \eta_l^2 \tau (\sigma_g^2 + 4\tau \kappa'^2) + \frac{12(M-m)L_c^2}{\tau^2 m(M-1)} \right) + \frac{4L_s^2 \tau^3 \eta_l^3 m \eta_g \kappa'^2}{\epsilon} + \frac{L_s^2 \tau^2 \eta_l^2 m \eta_g \sigma_g^2}{\epsilon} \quad (127)$$

658 where  $\nu = L_s + L_c/4$ . With  $\eta_l \leq 1/4\beta'\tau L_s$ ,  $\epsilon = m$ , and  $\eta_g\eta_l \leq \frac{1}{9\tau\nu}$ , we can further bound above  
659 as

$$\mathbb{E} \left[ \tilde{F}(\mathbf{w}^{(t+1,0)}) \right] - \tilde{F}(\mathbf{w}^{(t,0)}) \leq -\frac{\eta_l\eta_g\tau}{4} \|\nabla \tilde{F}(\mathbf{w}^{(t,0)})\|^2 + (\tau\eta_l\eta_g)^2 \nu \frac{\sigma_g^2}{m\tau} \\ + (\tau\eta_l\eta_g)^2 \nu \left( 6L_s^2 \eta_l^2 \tau (\sigma_g^2 + 4\tau \kappa'^2) + \frac{12(M-m)L_c^2}{\tau^2 m(M-1)} \right) + 4L_s^2 \tau^3 \eta_l^3 \eta_g \kappa'^2 + L_s^2 \tau^2 \eta_l^2 \eta_g \sigma_g^2 \quad (128)$$

660 Taking the average across all rounds on both sides of (128) and rearranging the terms we get

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \|\nabla \tilde{F}(\mathbf{w}^{(t,0)})\|^2 \right] \leq \frac{4 \left( \tilde{F}(\mathbf{w}^{(0,0)}) - \tilde{F}_{\text{inf}} \right)}{T\eta_l\eta_g\tau} + 4\sigma_g^2 \eta_l \left( \frac{\eta_g\nu}{m} + \frac{2L_s^2 \eta_l \tau}{3} + L_s^2 \tau \right) \\ + \frac{80L_s^2 \eta_l^2 \tau^2 \kappa'^2}{3} + \frac{48\eta_l\eta_g\nu(M-m)L_c^2}{\tau m(M-1)} \quad (129)$$

661 With the small enough learning rate  $\eta_l = 1/(\sqrt{T}\tau)$  and  $\eta_g = \sqrt{\tau m}$  one can prove that

$$\min_{t \in [T]} \mathbb{E} \left[ \|\nabla \tilde{F}(\mathbf{w}^{(t,0)})\|^2 \right] \leq \frac{4 \left( \tilde{F}(\mathbf{w}^{(0,0)}) - \tilde{F}_{\text{inf}} \right) + 4\sigma_g^2 \nu}{\sqrt{T}\tau m} + \frac{4\sigma_g^2 L_s^2}{\sqrt{T}} + \frac{8\sigma_g^2 L_s^2}{3\tau T} \\ + \frac{80L_s^2 \kappa'^2}{T} + \frac{48\nu(M-m)L_c^2 \sqrt{\tau}}{\sqrt{T}m} \\ = \mathcal{O} \left( \frac{\sigma_g^2}{\sqrt{T}\tau m} \right) + \mathcal{O} \left( \frac{\sigma_g^2}{\tau T} \right) + \mathcal{O} \left( \frac{\kappa'^2}{T} \right) + \mathcal{O} \left( \frac{\sqrt{\tau}}{\sqrt{T}m} \right) \quad (131)$$

662 completing the proof for Theorem C.1 for partial client participation.

## 663 E Simulation Details for Fig. 4a

664 For the mean estimation simulation for Fig. 4(a), we set the true means for the two clients as  
665  $\theta_1 = 0$ ,  $\theta_2 = 2\gamma_G$  where  $\gamma_G \in [0, \sqrt{20}]$ . The simulation was performed using NumPy [41] and  
666 SciPy [42]. The empirical means  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are sampled from the distribution  $\mathcal{N}(\theta_1, 1)$  and  $\mathcal{N}(\theta_2, 1)$

667 respectively where the number of samples are assumed to be identical for simplicity. For local  
668 training we assume clients set their local models as their local empirical means which is analogous to  
669 clients performing a large number of local SGD steps to obtain the local minima of their empirical  
670 loss. For the global objective (standard FL, INCFL (ReLU), INCFL) a local minima is found using  
671 the `scipy.optimize` function in the SciPy package. For each  $\gamma_G^2 \in [0, \sqrt{20}]$ , the average IPR is  
672 calculated over 10000 runs for each global objective.

## 673 F Experiment Details and Additional Results

674 All experiments are conducted on clusters equipped with one NVIDIA TitanX GPU. The algorithms  
675 are implemented in PyTorch 1.11.0. All experiments are run with 3 different random seeds and  
676 the average performance with the standard deviation is shown. The code used for all experiments is  
677 included in the supplementary material.

### 678 F.1 Experiment Details

679 **Obtaining  $\hat{\mathbf{w}}_i$ ,  $i \in [M]$  for INCFL Results in Section 4.** In INCFL, we use  $\hat{\mathbf{w}}_i$ ,  $i \in [M]$  to  
680 calculate the aggregating weights (see Algorithm 1). For all experiments with INCFL, we obtain  
681  $\hat{\mathbf{w}}_i$ ,  $i \in [M]$  at each client by each client taking 100 local SGD steps on its local dataset with its own  
682 separate local model before starting federated training. We use the same batch-size and learning rate  
683 used for the local training at clients done after we start the federated training (line 8-9 in Algorithm 1).  
684 The specific values are mentioned in the next paragraph.

685 **Local Training and Hyperparameters.** For all experiments, we do a grid search over the required  
686 hyperparameters to find the best performing ones. Specifically, we do a grid search over the learning  
687 rate:  $\eta_l \eta_g \in \{0.1, 0.05, 0.01, 0.005, 0.001\}$ , batchsize:  $b \in \{32, 64, 128\}$ , and local iterations:  
688  $\tau \in \{10, 30, 50\}$  to find the hyper-parameters with the highest test accuracy for each benchmark.  
689 For all benchmarks we use the best hyper-parameter for each benchmark after doing a grid search  
690 over feasible parameters referring to their source codes that are open-sourced. For a fair comparison  
691 across all benchmarks we do not use any learning rate decay or momentum.

692 **Logistic Regression on the Synthetic Dataset.** We conduct simulations on synthetic data which  
693 allows precise manipulation of heterogeneity. Using the methodology constructed in [2], we use the  
694 dataset with large data heterogeneity, `Synthetic(1,1)`. We have in total 100 devices where the local  
695 dataset sizes for each device follows the power law. The dimension used for logistic regression is  
696  $\mathbb{R}^{61 \times 10}$  where 10 is the output dimension.

697 **DNN Experiments.** For FMNIST, we train a deep multi-layer perceptron network with 2 hidden  
698 layers of units [64, 30] with dropout after the first hidden layer where the input is the normalized  
699 flattened image and the output is consisted of 10 units each of one of the 0-9 labels. For CIFAR10,  
700 we train a deep convolutional neural network with 2 convolutional layers with max pooling and 4  
701 hidden fully connected linear layers of units [120, 100, 84, 50]. The input is the normalized flattened  
702 convolution output and the output is consisted of 10 units each of one of the 0-9 labels. For Sent140,  
703 we train a deep multi-layer perceptron network with 3 hidden layers of units [128, 86, 30] with  
704 pre-trained 200D average-pooled GloVe embedding [43]. The input is the embedded 200D vector  
705 and the output is a binary classifier determining whether the tweet sentiment is positive or negative  
706 with labels 0 and 1 respectively. All clients have at least 50 data samples.

### 707 F.2 Additional Experimental Results

708 **Local Tuning for Personalization.** Personalized FL methods can be used to fine-tune the global  
709 model at each client before comparing it with that client’s locally trained model. INCFL can be  
710 combined with these methods by simply allowing clients to perform some fine-tuning iterations  
711 before computing the aggregation weights in Step 7 of Algorithm 1. Both for clients that are active  
712 during training and unseen test clients, we show in Table 2 that INCFL increases the fraction of  
713 incentivized clients by at least 10% as compared to all baselines. For FMNIST, CIFAR10, and  
714 Sent140, the improvement in IPR over other methods is up to 27%, 39%, and 28% respectively for  
715 active clients and 17%, 35%, and 4% respectively for the unseen incoming clients.

Table 2: Incentivized participation rate (IPR) of locally-tuned models with 5 local steps from the final global models trained with different algorithms for seen clients and unseen clients (the corresponding preferred-model test accuracy is in Appendix F.2).

	Seen Clients			Unseen Clients		
	FMNIST	CIFAR10	Sent140	FMNIST	CIFAR10	Sent140
FedAvg	0.38 ( $\pm 0.06$ )	0.19 ( $\pm 0.07$ )	0.25 ( $\pm 0.09$ )	0.39 ( $\pm 0.06$ )	0.20 ( $\pm 0.07$ )	0.42 ( $\pm 0.06$ )
FedProx	0.40 ( $\pm 0.07$ )	0.17 ( $\pm 0.07$ )	0.26 ( $\pm 0.09$ )	0.41 ( $\pm 0.07$ )	0.19 ( $\pm 0.07$ )	0.43 ( $\pm 0.12$ )
PerFedAvg	0.45 ( $\pm 0.05$ )	0.26 ( $\pm 0.02$ )	0.24 ( $\pm 0.10$ )	0.46 ( $\pm 0.06$ )	0.28 ( $\pm 0.04$ )	0.47 ( $\pm 0.06$ )
MW-Fed	0.28 ( $\pm 0.07$ )	0.01 ( $\pm 0.01$ )	0.08 ( $\pm 0.01$ )	0.39 ( $\pm 0.04$ )	0.06 ( $\pm 0.03$ )	0.20 ( $\pm 0.01$ )
INCFL	<b>0.55</b> ( $\pm 0.01$ )	<b>0.40</b> ( $\pm 0.00$ )	<b>0.36</b> ( $\pm 0.05$ )	<b>0.56</b> ( $\pm 0.01$ )	<b>0.41</b> ( $\pm 0.01$ )	<b>0.55</b> ( $\pm 0.01$ )

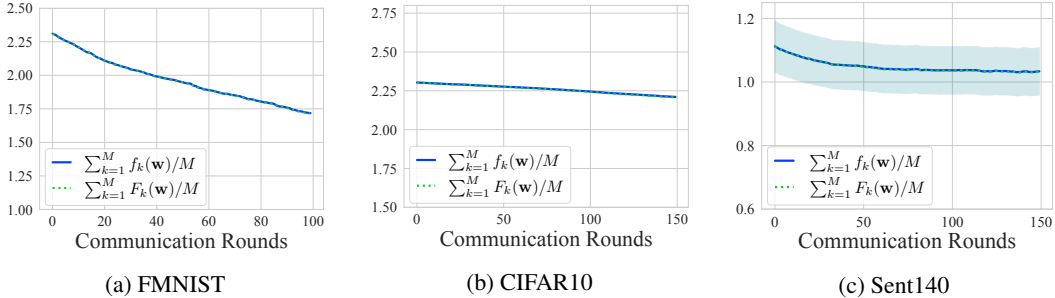


Figure 7: Comparison of the average of the true local losses across all clients ( $\sum_{k=1}^M f_k(\mathbf{w})/M$ ) and the empirical local losses across all clients ( $\sum_{k=1}^M F_k(\mathbf{w})/M$ ) where the former is calculated on the test dataset and the latter is calculated on the training dataset for the global model  $\mathbf{w}$ . We show that the average of the true local losses is nearly identical to the average empirical local loss across all clients empirically validating our relaxation of replacing  $f_k(\mathbf{w})$  with  $F_k(\mathbf{w})$ .

716 **Ablation Study on  $f_k(\mathbf{w}) \approx F_k(\mathbf{w})$ .** One of the two key relaxations we use for INCFL (see  
717 Section 2.1) is that we replace  $f_k(\mathbf{w}) - f_k(\hat{\mathbf{w}}_k)$  with  $F_k(\mathbf{w}) - F_k(\hat{\mathbf{w}}_k)$ . In other words, we replace  
718 the true loss  $f_k(\mathbf{w}) = \mathbb{E}_{\xi \sim \mathcal{D}_k}[\ell(\mathbf{w}, \xi)]$  with the empirical loss  $F_k(\mathbf{w}) = \frac{1}{|\mathcal{B}_k|} \sum_{\xi \in \mathcal{B}_k} \ell(\mathbf{w}, \xi)$  for  
719 all clients  $k \in [M]$ . We have used the likely conjecture that the global model  $\mathbf{w}$  is trained on the  
720 data of all clients, making it unlikely to overfit to the local data of any particular client, leading to  
721  $f_k(\mathbf{w}) \approx F_k(\mathbf{w})$ . We show in Fig. 7 that this is indeed the case. For all DNN experiments, we  
722 show that the average true local loss across all clients, i.e.,  $\sum_{k=1}^M f_k(\mathbf{w})/M$  is nearly identical to the  
723 average empirical local loss across all clients, i.e.,  $\sum_{k=1}^M F_k(\mathbf{w})/M$  given the training of the global  
724 model  $\mathbf{w}$  throughout the communication rounds. This empirically validates our relaxation of the true  
725 local losses to the empirical local losses.

726 **Preferred-model Test Accuracy for the Local-Tuning Results in Table 2.** In Table 2, we have  
727 shown how INCFL can largely increase the fraction of incentivized clients compared to the other base-  
728 lines even when jointly used with local-tuning. In Table 3, we show the corresponding preferred-model  
729 test accuracies. We show that for the seen clients that were active during training, INCFL achieves at  
730 least the same or higher preferred-model test accuracy than the other methods for all the different  
731 datasets. Hence, the clients are able to also gain from INCFL by achieving the highest accuracy in  
732 average with their preferred models (either global model or solo-trained local model). For the unseen  
733 clients with FMNIST, FedProx achieves a slightly higher preferred-model test accuracy (+0.05) than  
734 INCFL but with a much lower IPR of 0.46 (see Table 2) as INCFL’s IPR is 0.56. For the other datasets  
735 with unseen clients, INCFL achieves at least the same or higher preferred-model test accuracy than  
736 the other methods. This demonstrates that INCFL consistently largely improves the IPR compared to  
737 the other methods while losing very little, if any, in terms of the preferred-model test accuracy.

738 **Comparison with Algorithms for Fairness** Fair FL methods [33, 34] aim in training a global  
739 model that yields small variance across the clients’ test accuracies. These methods may incentivize  
740 the worst performing clients to participate, but potentially at the cost of disincentivizing the best  
741 performing clients. We show in Table 4 that the common fair FL methods are indeed not effective in



Table 3: Preferred-model test accuracy with the locally-tuned models with 5 local steps from the final global models trained with different algorithms for seen clients’ and unseen clients’ test data (the corresponding IPR is in Table 2).

	Seen Clients			Unseen Clients		
	FMNIST	CIFAR10	Sent140	FMNIST	CIFAR10	Sent140
FedAvg	99.37 ( $\pm 0.24$ )	100.00 ( $\pm 0.00$ )	55.71 ( $\pm 0.46$ )	99.50 ( $\pm 0.02$ )	100.00 ( $\pm 0.00$ )	58.79 ( $\pm 0.67$ )
FedProx	99.35 ( $\pm 0.23$ )	100.00 ( $\pm 0.00$ )	55.75 ( $\pm 0.80$ )	<b>99.55</b> ( $\pm 0.09$ )	100.00 ( $\pm 0.00$ )	58.82 ( $\pm 0.72$ )
PerFedAvg	99.20 ( $\pm 0.25$ )	100.00 ( $\pm 0.00$ )	55.74 ( $\pm 0.80$ )	98.98 ( $\pm 0.55$ )	100.00 ( $\pm 0.00$ )	58.82 ( $\pm 0.72$ )
MW-Fed	99.27 ( $\pm 0.39$ )	100.00 ( $\pm 0.00$ )	55.06 ( $\pm 0.38$ )	99.47 ( $\pm 0.08$ )	100.00 ( $\pm 0.00$ )	57.36 ( $\pm 0.71$ )
INCFL	<b>99.40</b> ( $\pm 0.30$ )	100.00 ( $\pm 0.00$ )	<b>55.82</b> ( $\pm 0.82$ )	99.50 ( $\pm 0.02$ )	100.00 ( $\pm 0.00$ )	<b>58.88</b> ( $\pm 0.77$ )

Table 4: Incentivized participation rate (IPR) and preferred-model test accuracy for the seen clients’ test data with the final global models trained via INCFL and q-FFL [33] which aims in improving fairness. The baseline q-FFL with large  $q$ , e.g.  $q = 10$ , emulates the behavior of another well-known algorithm for improving fairness named AFL [34].

	Incentivized Participation Rate (IPR)			Preferred-Model Test Acc.		
	FMNIST	CIFAR10	Sent140	FMNIST	CIFAR10	Sent140
q-FFL ( $q = 1$ )	0.03 ( $\pm 0.01$ )	0.00 ( $\pm 0.00$ )	0.09 ( $\pm 0.06$ )	99.24 ( $\pm 0.05$ )	100.00 ( $\pm 0.00$ )	53.10 ( $\pm 2.63$ )
q-FFL ( $q = 10$ )	0.00 ( $\pm 0.00$ )	0.00 ( $\pm 0.00$ )	0.09 ( $\pm 0.00$ )	98.90 ( $\pm 0.01$ )	100.00 ( $\pm 0.00$ )	52.71 ( $\pm 1.40$ )
INCFL	<b>0.55</b> ( $\pm 0.00$ )	<b>0.40</b> ( $\pm 0.00$ )	<b>0.41</b> ( $\pm 0.07$ )	<b>99.29</b> ( $\pm 0.03$ )	100.00 ( $\pm 0.00$ )	<b>53.93</b> ( $\pm 1.87$ )

742 improving the overall clients’ incentivized participation rate. We see that the fair FL methods achieve  
743 an incentivized participation rate lower than 0.01 for all datasets while INCFL achieves at least 0.40  
744 for all datasets. Moreover, the preferred-model test accuracy is also higher for INCFL compared to  
745 the fair FL methods. This underwhelming performance of fair FL methods in incentivizing clients can  
746 be due to the fact that fair FL methods try to find the global model that performs well, in overall, over  
747 *all* clients which results in failing to incentivize *any* client.