

---

# Gradient-Driven Rewards to Guarantee Fairness in Collaborative Machine Learning

---

Xinyi Xu<sup>1,5\*</sup>, Lingjuan Lyu<sup>2\*</sup>, Xingjun Ma<sup>3</sup>, Chenglin Miao<sup>4</sup>,  
Chuan Sheng Foo<sup>5</sup>, and Bryan Kian Hsiang Low<sup>1</sup>

Department of Computer Science, National University of Singapore, Republic of Singapore<sup>1</sup>  
Sony AI<sup>2</sup>, School of Computer Science, Fudan University, People's Republic of China<sup>3</sup>

Department of Computer Science, University of Georgia, USA<sup>4</sup>

Institute for Infocomm Research, A\*STAR, Republic of Singapore<sup>5</sup>

{xuxinyi, lowkh}@comp.nus.edu.sg<sup>1</sup>, lingjuan.lv@sony.com<sup>2</sup>

danxjma@gmail.com<sup>3</sup>, cmiao@uga.edu<sup>4</sup>, foo\_chuan\_sheng@i2r.a-star.edu.sg<sup>5</sup>

## Abstract

In *collaborative machine learning* (CML), multiple agents pool their resources (e.g., data) together for a common learning task. In realistic CML settings where the agents are self-interested and not altruistic, they may be unwilling to share data or model information without adequate rewards. Furthermore, as the data/model information shared by the agents may differ in quality, designing rewards which are fair to them is important so that they would not feel exploited nor discouraged from sharing. In this paper, we adopt *federated learning* as the CML paradigm, propose a novel *cosine gradient Shapley value* (CGSV) to fairly evaluate the expected marginal contribution of each agent's uploaded model parameter update/gradient without needing an auxiliary validation dataset, and based on the CGSV, design a novel training-time gradient reward mechanism with a fairness guarantee by sparsifying the aggregated parameter update/gradient downloaded from the server as reward to each agent such that its resulting quality is commensurate to that of the agent's uploaded parameter update/gradient. We empirically demonstrate the effectiveness of our fair gradient reward mechanism on multiple benchmark datasets in terms of fairness, predictive performance, and time overhead.

## 1 Introduction

In *collaborative machine learning* (CML), multiple agents (e.g., researchers, organizations, companies) pool their resources (e.g., data) together for a common learning task. It spans a wide variety of real-world applications such as digital healthcare [49], clinical trial research [13, 23], wake word detection for smart voice assistants [27], and next word prediction on mobile devices [15].

*Federated learning* (FL) provides a natural paradigm of CML [18, 29, 41, 43, 57, 62]. In FL, the agents perform local model training (e.g., using stochastic gradient descent) and share their resulting model parameter updates/gradients via a *trusted server* [40, 56, 59]. An important distinction of our work here from the standard FL literature is that the agents are self-interested and hence not necessarily cooperative like the worker nodes in distributed learning. The implication is that to achieve competitive predictive performance for the learning task, it is imperative to incentivize/reward the agents for contributing/sharing high-quality information in the form of model parameter updates/gradients [47, 48, 52].

---

\*Equal contribution.

Our work here adopts FL as the CML paradigm for designing a fair reward mechanism such that the (self-interested) agents who contribute more would not feel exploited but be rewarded commensurately. This is often regarded as fairness in cooperative game theory [42], mechanism design [4], and computational social choice [11]. To design such a fair reward mechanism, we need to address three main questions:

Firstly, *what is a suitable notion of fairness?* The *Shapley value* (SV) [50] from cooperative game theory is an appealing choice and has been used in ML [14] and FL [54, 56]. However, existing SV-based works [19, 37, 54, 56] typically require the availability of (and all agents to agree on) an auxiliary validation dataset and significant time overhead from evaluating the agents’ contributions in the form of SVs and the resulting model training. To overcome these difficulties, we propose to instead exploit the alignment (specifically, cosine similarity) of an agent’s uploaded/contributed model parameter update/gradient vector (or that aggregated over some agents) to that aggregated over all agents (hence measuring its quality/value and circumventing the need for a validation dataset [12, 52]) for devising our proposed *cosine gradient Shapley value* (CGSV) (Sec. 3.2) which can be efficiently approximated with a bounded error (Sec. 3.3).

Secondly, *what is the choice of reward?* Various choices such as monetary rewards from a pre-allocated budget [65, 66] or the total revenue generated from the collaboration through FL [9, 10] have been proposed. Though it may seem natural to consider monetary rewards, it is not obvious how a common denomination between money and data/gradients [1, 46] can be readily established, which makes it challenging to apply these works in practice. Instead, we propose to consider the aggregated parameter updates/gradients downloaded from the server as rewards to the agents.

Finally, *how can the gradient reward mechanism ensure fairness?* Our proposed mechanism exploits a *sparsifying gradient* trick (Sec. 3.4) for controlling the quality of the aggregated parameter update/gradient downloaded from the server as reward to each agent (rather than post hoc [48, 52, 65]) such that its quality is commensurate to that of the agent’s uploaded/contributed parameter update/gradient [2, 7]. Consequently, an agent who uploads/contributes higher-quality parameter updates/gradients over the entire training process should eventually be rewarded with converged model parameters whose resulting training loss (and hence predictive performance) is closer to that of the server, as demonstrated in our fairness guarantee (Sec. 3.5) [52].

In summary, the contributions of our work here to CML and FL include the following:

- We propose a novel *cosine gradient Shapley value* (CGSV) (Sec. 3.2) to fairly evaluate the expected marginal contribution of each agent’s uploaded model parameter update/gradient without needing an auxiliary validation dataset and present an efficient approximation of CGSV with a bounded error (Sec. 3.3).
- Based on the approximate CGSV, we design a novel training-time gradient reward mechanism (Sec. 3.4) with a fairness guarantee (Sec. 3.5) by exploiting the trick of sparsifying the aggregated parameter update/gradient downloaded from the server as reward to each agent such that its resulting quality is commensurate to that of the agent’s uploaded/contributed parameter update/gradient.
- We empirically demonstrate the effectiveness of our fair gradient reward mechanism on multiple benchmark datasets in terms of fairness, predictive performance, and time overhead (Sec. 4).

## 2 Related Work

**Reward design and choice in CML.** In related topics such as FL [30, 36, 38, 47, 59, 63, 66], Bayesian CML [52], collaborative generative modeling [55], and data sharing [13, 23, 48], designing appropriate rewards to encourage collaboration (e.g., sharing real or synthetic data, gradients, or other information) is a non-trivial problem. A useful solution concept should provide a formal notion of fairness, a suitable form/denomination of reward, and a principled way to guarantee fairness via a carefully designed reward mechanism. Previous works have considered monetary rewards from a pre-allocated budget [65, 66] or the total revenue generated from the collaboration [9, 10], or simply an abstract yet quantifiable form of reward [47, 48]. Though it may seem natural to consider monetary rewards, it is not obvious how a common denomination between money and data/gradients [1, 46] can be readily established, which makes it challenging to apply these works in practice. The work of [66] has explored a different avenue of using a reverse auction to guarantee truthfulness in its mechanism instead of fairness.

**Fairness notions.** The *Shapley value* (SV) [50] from cooperative game theory is widely regarded as a principled notion of fairness [4, 11, 42] due to its several desirable properties such as symmetry and null player. Existing SV-based works have considered fairness in the sense of rewarding agents according to their contributions [19, 54, 56]. However, they typically require the availability of (and all agents to agree on) an auxiliary validation dataset [37, 52] and significant time overhead from evaluating the agents’ contributions in the form of SVs and the resulting model training [14, 19, 56]. In contrast, the work of [31] has adopted an egalitarian notion of fairness by aiming to equalize the final individual performance among agents, which is fundamentally different from SV.

Different from the fairness definition in [31], we adopt a fairness notion formalized by SV [14, 19, 52, 54, 56]. Our proposed work is novel in the application of SV: While previous works use the validation accuracy [14, 19, 54, 56], we exploit the cosine similarity between model parameter updates/gradient vectors [12] for devising our proposed *cosine gradient Shapley value* (CGSV) (Sec. 3.2) to fairly evaluate the expected marginal contribution of each agent’s uploaded model parameter update/gradient. Based on the CGSV, we design a novel training-time gradient reward mechanism (Sec. 3.4) with a fairness guarantee (Sec. 3.5) and empirically show that it outperforms several existing FL baselines in terms of predictive performance, fairness, and time overhead (Sec. 4.2).

### 3 Fair Gradient Reward Mechanism

#### 3.1 Vanilla Federated Learning (FL) Problem Setting and Notations

The vanilla FL problem [56, 59] involves a set  $\mathcal{N} := \{i\}_{i=1,\dots,N}$  of  $N$  *honest* agents learning a  $D$ -dimensional vector  $\mathbf{w} \in \mathbb{R}^D$  of model parameters to minimize a loss function  $\mathbf{F}(\mathbf{w})$  that can be additively decomposed into  $N$  local differentiable loss functions  $\mathbf{F}_i(\mathbf{w})$  defined using the local dataset  $\mathcal{D}_i$  of agent  $i \in \mathcal{N}$  and weighted by its importance  $p_i \geq 0$  (e.g., proportional to  $|\mathcal{D}_i|$ ). That is,  $\mathbf{F}(\mathbf{w}) := \sum_{i \in \mathcal{N}} p_i \mathbf{F}_i(\mathbf{w})$  where  $\sum_{i \in \mathcal{N}} p_i = 1$ . We call  $\mathcal{N}$  the grand coalition; a coalition  $\mathcal{S} \subseteq \mathcal{N}$  is then a subset of the grand coalition  $\mathcal{N}$  of  $N$  agents. In iteration  $t = 0$ , every agent  $i \in \mathcal{N}$  starts with the same initialized parameter vector  $\mathbf{w}_{i,0} := \mathbf{w}_0$  as the server. In iteration  $t > 0$ , every agent  $i \in \mathcal{N}$  calculates a parameter update  $\Delta \mathbf{w}_{i,t} := -\eta_t \nabla \mathbf{F}_i(\mathbf{w}_{i,t-1})$  with step size  $\eta_t$  and gradient  $\nabla \mathbf{F}_i(\mathbf{w}_{i,t-1})$  w.r.t. parameter vector  $\mathbf{w}_{i,t-1}$  and uploads it to a *trusted* server who normalizes and aggregates all agents’ parameter updates as follows:

$$\mathbf{u}_{i,t} := \Gamma \Delta \mathbf{w}_{i,t} / \|\Delta \mathbf{w}_{i,t}\|, \quad \mathbf{u}_{\mathcal{N},t} := \sum_{i \in \mathcal{N}} r_{i,t-1} \mathbf{u}_{i,t} \quad (1)$$

where  $\Gamma$  is a normalization coefficient used to prevent gradient explosion [33, 45] and the importance coefficient  $r_{i,t-1}$  will be described later in Sec 3.4. So, we call (1) the *gradient aggregation step*. The *gradient download step* then follows where every agent  $i \in \mathcal{N}$  downloads the aggregated parameter update/gradient  $\mathbf{u}_{\mathcal{N},t}$  (1) from the server (as reward) for updating its model parameters  $\mathbf{w}_{i,t} := \mathbf{w}_{i,t-1} + \mathbf{u}_{\mathcal{N},t}$  to the same  $\mathbf{w}_t := \mathbf{w}_{t-1} + \mathbf{u}_{\mathcal{N},t}$  as the server. That is,  $\mathbf{w}_{i,t} = \mathbf{w}_t$  for all  $i \in \mathcal{N}$  and  $t \in \mathbb{Z}^+ \cup \{0\}$ . We define  $\mathbf{u}_{\mathcal{S},t}$  for any coalition  $\mathcal{S} \subseteq \mathcal{N}$  in a similar way as  $\mathbf{u}_{\mathcal{N},t}$  (1). For brevity, we omit  $t$  from our notations in Secs. 3.2 and 3.3 since we only refer to iteration  $t$ .

#### 3.2 Cosine Gradient Shapley Value (CGSV) for Fairness

In the gradient aggregation step (1), the quality/value of coalition  $\mathcal{S}$ ’s (normalized) aggregated parameter update/gradient  $\mathbf{u}_{\mathcal{S}}$  can be measured by its *cosine similarity*  $\cos(\mathbf{u}_{\mathcal{S}}, \mathbf{u}_{\mathcal{N}}) := \langle \mathbf{u}_{\mathcal{S}}, \mathbf{u}_{\mathcal{N}} \rangle / (\|\mathbf{u}_{\mathcal{S}}\| \|\mathbf{u}_{\mathcal{N}}\|)$  to the grand coalition  $\mathcal{N}$ ’s aggregated parameter update/gradient  $\mathbf{u}_{\mathcal{N}}$  [12, 28, 35]. We use this cosine similarity measure as our *gradient valuation function*  $\nu(\mathcal{S}) := \cos(\mathbf{u}_{\mathcal{S}}, \mathbf{u}_{\mathcal{N}})$ . Intuitively, if the direction of  $\mathbf{u}_{\mathcal{S}}$  aligns more closely with that of  $\mathbf{u}_{\mathcal{N}}$ , then its quality/value  $\nu(\mathcal{S})$  is higher. Using  $\nu$ , the contribution  $\phi_i$  of agent  $i \in \mathcal{N}$  is defined based on the notion of *Shapley value* (SV) [50] which measures its expected marginal contribution when joining the other agents preceding it in any permutation and satisfies certain desirable fairness properties [5], such as null player (i.e., an agent with no marginal contribution has zero SV), symmetry (i.e., agents with identical marginal contributions have equal SVs), among others, as formally discussed in Appendix A.1:

**Definition 1 (Cosine gradient Shapley value (CGSV)).** Let  $\Pi_{\mathcal{N}}$  be a set of all possible permutations of  $\mathcal{N}$  and  $\mathcal{S}_{\pi,i}$  be the coalition of agents preceding agent  $i$  in permutation  $\pi \in \Pi_{\mathcal{N}}$ . The CGSV of agent  $i \in \mathcal{N}$  is defined as

$$\phi_i := (1/N!) \sum_{\pi \in \Pi_{\mathcal{N}}} [\nu(\mathcal{S}_{\pi,i} \cup \{i\}) - \nu(\mathcal{S}_{\pi,i})]. \quad (2)$$

If  $\phi_i$  is negative, then it follows from the weighted sum of parameter updates/gradients in (1) that  $\mathbf{u}_i$  points in an opposite direction to some other parameter updates/gradients and hence has negative cosine similarities to them. In practice, due to the noisy training arising from the use of *stochastic gradient descent* (SGD) and/or a highly non-convex loss function,  $\phi_i$  may at times be negative even for an honest agent  $i$ . When the number of such cases is limited, training via SGD can still converge to yield a competitive predictive performance, as empirically validated in [12].

### 3.3 Efficient Approximation of CGSV

Since evaluating agent  $i$ 's CGSV  $\phi_i$  (2) exactly incurs  $\mathcal{O}(2^N D)$  time and is thus costly, we propose an efficient approximation by directly measuring the cosine similarity of its (normalized) parameter update/gradient  $\mathbf{u}_i$  to the grand coalition  $\mathcal{N}$ 's aggregated parameter update/gradient  $\mathbf{u}_{\mathcal{N}}$ , which reduces the incurred time by a factor of  $2^N$  and has a bounded error from  $\phi_i$  (Theorem 1):

$$\phi_i \approx \psi_i := \cos(\mathbf{u}_i, \mathbf{u}_{\mathcal{N}}). \quad (3)$$

**Theorem 1 (Approximation Error).** *Let  $I \in \mathbb{R}^+$ . Suppose that  $\|\mathbf{u}_i\| = \Gamma$  and  $|\langle \mathbf{u}_i, \mathbf{u}_{\mathcal{N}} \rangle| \geq 1/I$  for all  $i \in \mathcal{N}$ . Then,  $\phi_i - L_i \psi_i \leq I\Gamma^2$  where the multiplicative factor  $L_i$  can be normalized away.*

Its proof is in Appendix A.2. From Theorem 1, the approximation error is bounded and decreases quadratically with normalization coefficient  $\Gamma$ . However,  $\Gamma$  cannot be reduced to be arbitrarily small, which may cause  $|\langle \mathbf{u}_i, \mathbf{u}_{\mathcal{N}} \rangle| \geq 1/I$  not to hold. It also does not hold when  $\mathbf{u}_i$  is orthogonal to  $\mathbf{u}_{\mathcal{N}}$  or is close to the zero vector, hence implying the quality of that agent  $i$ 's parameter update/gradient is not high enough. So, every agent is encouraged to contribute a parameter update/gradient of sufficiently high quality in order to ensure the quality of the approximation  $\psi_i$  (Theorem 1).

We have performed a simple experiment to compare the quality of our approximation  $\psi_i$  with that of a sampling-based  $(\epsilon, \delta)$ -approximation  $\bar{\phi}_i$  [39], the latter of which is widely used by existing works in data valuation and CML/FL [14, 19, 52, 56, 60]. In this experiment, we have drawn  $N$  random  $D$ -dimensional vectors from a standard multivariate normal distribution to simulate  $\mathbf{u}_1, \dots, \mathbf{u}_N$  and calculated the resulting exact CGSVs  $\phi := (\phi_i)_{i=1, \dots, N}$ , our approximation  $\psi := (\psi_i)_{i=1, \dots, N}$ , and the sampling-based  $(0.1, 0.1)$ -approximation  $\bar{\phi} := (\bar{\phi}_i)_{i=1, \dots, N}$ . Fig. 1 shows the results for  $\ell_1$  error,  $\ell_2$  error, and the incurred time averaged over 10 runs: Our approximation  $\psi$  performs better in all three metrics with varying  $D$  (right figure) and the performance gap widens with an increasing number  $N$  of agents (left figure).

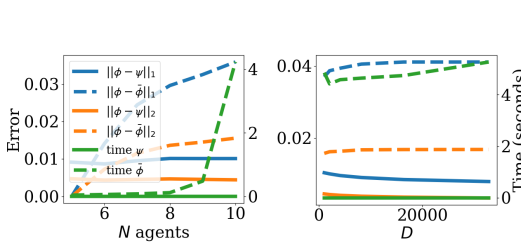


Figure 1: Comparison of  $\ell_1$  error (blue),  $\ell_2$  error (orange), and incurred time (green) (i.e., averaged over 10 runs) between our approximation  $\psi$  (solid lines) vs. a sampling-based approximation  $\bar{\phi}$  (dashed lines) [39] of the exact CGSVs  $\phi$  with (left) varying number  $N$  of agents and  $D = 1024$ , and (right) varying vector dimension  $D$  and  $N = 10$ . For all metrics, lower is better.

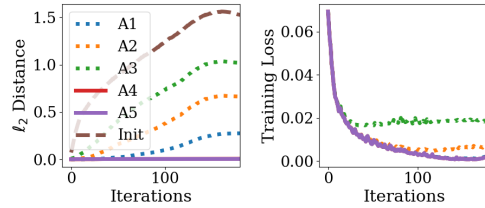


Figure 2: (Left)  $\ell_2$  distance between model parameters of agent  $i = 1, \dots, 5$  (abbreviated to  $A_i$ ) vs. that of the server, and (right) corresponding training loss for an FL problem with  $N = 5$  agents using local MNIST datasets of 600 images each to collaboratively learn 2-layer CNN parameters where the datasets of A1 (blue), A2 (orange), and A3 (green) have 20%, 40%, and 60% randomly corrupted labels, respectively. The brown line denotes  $\ell_2$  distance between  $\mathbf{w}_0$  (initialization) vs. server's model parameters.

### 3.4 Server-Side Training-Time Gradient Reward Mechanism

We will now describe the exact details of the gradient aggregation and download steps performed by the server to implement our proposed fair gradient reward mechanism:

**Gradient Aggregation Step.** With a specified normalization coefficient  $\Gamma$  and an initialized coefficient  $r_{i,0}$ , the server performs normalization and aggregation of all agents' parameter updates into  $\mathbf{u}_{\mathcal{N},t}$  using (1), as previously discussed in the FL problem setting (Sec. 3.1). Then, the server computes our approximation  $\psi_{i,t}$  (3) of the CGSV  $\phi_{i,t}$  (2) and updates (and normalizes) the importance coefficient  $r_{i,t}$  in iteration  $t$  via a moving average of  $\psi_{i,t}$  given the relative weight  $\alpha$  on  $r_{i,t-1}$  from previous iteration  $t - 1$ :

$$r_{i,t} := \alpha r_{i,t-1} + (1 - \alpha) \psi_{i,t}, \quad r_{i,t} \leftarrow r_{i,t} / \sum_{i' \in \mathcal{N}} r_{i',t} \quad (4)$$

where  $r_{i,0} := 0$ . Note that  $r_{i,t}$  (4) is used for deriving the sparsified gradient (5) in the gradient download step as well as the aggregation of all agents' parameter updates into  $\mathbf{u}_{\mathcal{N},t+1}$  (1) in iteration  $t + 1$ . The use of a moving average of  $\psi_{i,t}$  to compute  $r_{i,t}$  (4) provides a smoothed estimate without abrupt fluctuations and reduces the effect of noisy training due to the use of SGD in practice [31, 56]. It also allows a flexible weighting over the iterations of the entire training process: In particular, setting  $\alpha < 1$  can effectively mitigate the noise from random initialization of model parameters  $\mathbf{w}_0$  because the weight on  $\psi_{i,t'}$  in earlier iteration  $t' < t$  decays exponentially with  $t$  [54].

**Gradient Download Step.** Recall from the *vanilla* FL problem setting (Sec. 3.1) that in each iteration  $t$ , this step involves all agents downloading an identical aggregated parameter update/gradient  $\mathbf{u}_{\mathcal{N},t}$  (1) from the server (as reward) for updating their model parameters to the same  $\mathbf{w}_t$  (as the server), which is expected to converge to yield a competitive predictive performance [8, 32]. However, such *equal* rewards to all agents is unfair and will discourage any agent from uploading/contributing a parameter update/gradient of higher quality [37, 63] when it can afford to. To ensure fairness, each agent should download some form of aggregated parameter update/gradient as reward that is commensurate to the quality/value of its uploaded/contributed parameter update/gradient. Consequently, an agent who uploads/contributes higher-quality parameter updates/gradients over the entire training process should eventually be rewarded with converged model parameters whose resulting training loss (and hence predictive performance) is closer to that of the server (Theorem 2).

To achieve this, we adopt the trick of *sparsifying*<sup>2</sup> the aggregated parameter update/gradient  $\mathbf{u}_{\mathcal{N},t}$  downloaded from the server as reward to agent  $i$  in each iteration  $t$ . Specifically, we zero out fewer of its smallest components (hence higher-quality gradient reward) when the importance coefficient  $r_{i,t}$  (4) (i.e., moving average of the approximate CGSV  $\psi_{i,t}$ ) is larger:

$$\mathbf{v}_{i,t} := \text{mask}(\mathbf{u}_{\mathcal{N},t}, q_{i,t}), \quad q_{i,t} := \lfloor D \tanh(\beta r_{i,t}) / \max_{i' \in \mathcal{N}} \tanh(\beta r_{i',t}) \rfloor \quad (5)$$

where  $\text{mask}(\mathbf{u}, q)$  retains the largest  $\max(0, q)$  components (in magnitude) of  $\mathbf{u}$  and zeros out all of its other components [2, 61], and  $\beta \geq 1$  specifies the degree of altruism: Greater altruism  $\beta$  gives any agent with a smaller  $r_{i,t}$  a larger improvement in the quality of its gradient reward, i.e., a larger reduction in the sparsity of its downloaded  $\mathbf{v}_{i,t}$  as reward. In the extreme case of  $\beta = \infty$ , we recover the *vanilla* FL problem setting (Sec. 3.1) where all agents are rewarded equally with  $\mathbf{u}_{\mathcal{N},t}$  (i.e., best-quality gradient reward  $\mathbf{v}_{i,t} = \mathbf{u}_{\mathcal{N},t}$  for all  $i \in \mathcal{N}$  with no sparsification), albeit with importance coefficients  $r_{i,t}$  possibly differing across agents  $i \in \mathcal{N}$  and dynamically updated over iteration  $t \in \mathbb{Z}^+$ . Hence, increasing  $\beta$  from 1 to  $\infty$  trades off fairness for equality in gradient rewards by being more altruistic to any agent with a smaller  $r_{i,t}$ ; we empirically show the effect of varying  $\beta$  on training loss in Fig. 7 of Sec. 4.2. Note the agent  $i^* := \arg\max_{i' \in \mathcal{N}} \tanh(\beta r_{i',t})$  with the largest possible  $r_{i^*,t}$  does not benefit from such altruism since it already downloads the best-quality gradient reward (i.e.,  $\mathbf{u}_{\mathcal{N},t}$ ) according to (5).

Suppose that there exists a known threshold  $\underline{r} > 0$  s.t.  $r_{i,t} \geq \underline{r}$  for all  $i \in \mathcal{N}$  and  $t \in \mathbb{Z}^+$  and we want to limit the sparsity of any downloaded  $\mathbf{v}_{i,t}$  or, equivalently, ensure the minimum quality of any gradient reward: Specifically, given a predefined threshold  $c \in (0, 1]$ , we want to guarantee  $q_{i,t} \geq \lfloor D \times c \rfloor$  holds for all  $i \in \mathcal{N}$  and  $t \in \mathbb{Z}^+$ . By setting  $\beta$  s.t.  $\tanh(\beta \underline{r}) \geq c$ , it follows from (5) and  $\max_{i' \in \mathcal{N}} \tanh(\beta r_{i',t}) \leq 1$  that  $\tanh(\beta r_{i,t}) / \max_{i' \in \mathcal{N}} \tanh(\beta r_{i',t}) \geq \tanh(\beta r_{i,t}) \geq \tanh(\beta \underline{r}) \geq c$  and hence  $q_{i,t} \geq \lfloor D \times c \rfloor$  ensues. By using the property that  $\tanh(\beta \underline{r}) = (\exp(2\beta \underline{r}) - 1) / (\exp(2\beta \underline{r}) + 1)$ ,  $\beta \geq \ln((1 + c)/(1 - c)) / (2\underline{r})$  can be derived and used for setting  $\beta$ . It further informs us that reducing the sparsity of any downloaded  $\mathbf{v}_{i,t}$  or, equivalently, improving the minimum quality of any gradient reward (i.e., by increasing  $c$ ) requires greater altruism  $\beta$  to be introduced, while improving the minimum quality of uploaded/contributed parameter updates/gradients by any agent over the entire training process (hence larger  $\underline{r}$ ) eases the need of introducing greater altruism  $\beta$ .

<sup>2</sup>Sparsifying a parameter update/gradient vector means zeroing out some of its components and leaving the others unchanged [7, 33].

To see why the sparsifying gradient trick (5) can ensure fairness, we illustrate its effect in an FL problem with  $N = 5$  agents using local MNIST datasets of 600 images each to collaboratively learn the parameters of a 2-layer *convolutional neural network* (CNN) where the datasets of agents 1, 2, and 3 have 20%, 40%, and 60% randomly corrupted labels, respectively. The uploaded/contributed parameter updates/gradients thus decrease in quality from agents 1 to 3 (i.e.,  $\psi_{1,t} = 0.194$ ,  $\psi_{2,t} = 0.088$ , and  $\psi_{3,t} = 0.043$  on average) due to increasingly noisy labels in their datasets, while agents 4 and 5 upload/contribute parameter updates/gradients of high quality (i.e.,  $\psi_{4,t} = 0.331$  and  $\psi_{5,t} = 0.342$  on average) due to uncorrupted labels in their datasets. Consequently, agents 1 to 3 have increasing sparsity (resp., 34.9%, 67.6%, and 83.0% on average) while agents 4 and 5 have little/no sparsity (resp., 3.5% and 1.1% on average) in their downloaded  $\mathbf{v}_{i,t}$  as rewards ( $\beta = 1$ ). Fig. 2 shows that the converged model parameters of agents 1 to 3 grow in  $\ell_2$  distance from that of the server (hence increasing training loss) while agents 4 and 5 have the closest converged model parameters (hence lowest training loss).

We provide the pseudocodes performed by the server and agent  $i \in \mathcal{N}$  in each iteration  $t$  below. We will discuss in Sec 4.2 how the hyperparameters  $\Gamma$  in (1),  $\alpha$  in (4), and  $\beta$  in (5) are set in our experiments.

Server ( $t$ )	Agent ( $i, t$ )
1: <b>for all</b> $i \in \mathcal{N}$ <b>do</b> 2:   Download $\Delta \mathbf{w}_{i,t}$ from agent $i$ 3: $\triangleright$ <b>Gradient Aggregation Step</b> 4: Compute $\mathbf{u}_{i,t}$ and $\mathbf{u}_{\mathcal{N},t}$ (1) 5: <b>for all</b> $i \in \mathcal{N}$ <b>do</b> 6:   Compute $\psi_{i,t}$ (3) and $r_{i,t}$ (4) 7: $\triangleright$ <b>Gradient Download Step</b> 8: <b>for all</b> $i \in \mathcal{N}$ <b>do</b> 9:   Compute $\mathbf{v}_{i,t}$ (5) for download by agent $i$	1: Upload $\Delta \mathbf{w}_{i,t} = -\eta_t \nabla \mathbf{F}_i(\mathbf{w}_{i,t-1})$ to server 2: Download $\mathbf{v}_{i,t}$ from server 3: Update $\mathbf{w}_{i,t} = \mathbf{w}_{i,t-1} + \mathbf{v}_{i,t}$

### 3.5 Fairness Guarantee

We have previously discussed the intuition underlying our notion of fairness in Sec. 3.4 that an agent who uploads/contributes higher-quality parameter updates/gradients over the entire training process should eventually be rewarded with converged model parameters whose resulting training loss (and hence predictive performance) is closer to that of the server. Note that the importance coefficient  $r_{i,t}$  (4) measures the overall quality of the parameter updates/gradients uploaded/contributed by agent  $i$  over the entire training process till iteration  $t$ . Our main result below guarantees a notion of fairness that under some conditions on loss function  $\mathbf{F}$  and the server’s model parameters  $\mathbf{w}_t$ , if an agent  $i$  has a larger importance coefficient  $r_{i,t}$  and model parameters  $\mathbf{w}_{i,t-1}$  closer to that of the server (i.e.,  $\mathbf{w}_{t-1}$ ) than another agent by at least  $2\|\mathbf{v}_{i,t}\|$  in previous iteration  $t - 1$ , then it is rewarded with model parameters  $\mathbf{w}_{i,t}$  incurring smaller training loss  $\mathbf{F}(\mathbf{w}_{i,t})$  in iteration  $t$ :

**Theorem 2 (Fairness in Training Loss).** *Let  $\delta_{i,t} := \|\mathbf{w}_t - \mathbf{w}_{i,t}\|$ . Suppose that  $\mathbf{w}_t$  is near to a stationary point of  $\mathbf{F}$  for  $t \geq t^* \in \mathbb{Z}^+$  and some regularity conditions on  $\mathbf{F}$  hold. For all  $i, i' \in \mathcal{N}$  and  $t \geq t^*$ , if  $r_{i,t} \geq r_{i',t}$  and  $\delta_{i',t-1} - \delta_{i,t-1} \geq 2\|\mathbf{v}_{i,t}\|$ , then  $\mathbf{F}(\mathbf{w}_{i,t}) \leq \mathbf{F}(\mathbf{w}_{i',t})$ .*

Its proof is in Appendix A.3. Our experiments in Appendix B.3 will empirically verify the fairness guarantee in Theorem 2 (and fairness in test accuracy) without needing to impose its conditions.

## 4 Experiments and Discussion

### 4.1 Experimental Settings

**Datasets.** We perform extensive experiments on image classification datasets like MNIST [26] and CIFAR-10 [21] and text classification datasets like *movie review* (MR) [44] and *Stanford sentiment treebank* (SST) [20]. We use a 2-layer *convolutional neural network* (CNN) for MNIST [25], a 3-layer CNN for CIFAR-10 [22], and a text embedding CNN for MR and SST [20].

**Baselines.** We consider several existing FL baselines such as FedAvg [40],  $q$ -FFL[31], CFFL [37], and an *extended contribution index* (ECI) method from [54] utilizing validation accuracy-based SV

Table 1: Average test accuracy (%) achieved by the agents collaborating via our fair gradient reward mechanism with varying degrees of altruism  $\beta$  vs. tested baselines on all datasets. Each value in brackets denotes the highest test accuracy achieved by any agent.

No. Agents	MNIST						CIFAR-10			MR	SST
	10			20			10			5	5
	UNI	POW	CLA	UNI	POW	CLA	UNI	POW	CLA	POW	POW
Standalone	91 (91)	88 (92)	53 (92)	91 (91)	89 (92)	48 (90)	46 (47)	43 (49)	31 (44)	47(56)	31(34)
FedAvg	93 (94)	92 (94)	53 (93)	93 (93)	92 (94)	49 (92)	48 (48)	47 (50)	32 (47)	51(63)	33(35)
q-FFL	85 (91)	27 (45)	44 (64)	88 (91)	48 (53)	40 (59)	41 (46)	36 (36)	22 (28)	12(18)	23(25)
CFFL	90 (92)	85 (90)	34 (44)	91 (93)	88 (91)	39 (46)	39 (41)	35 (45)	22 (40)	44(53)	31(32)
ECI	94 (94)	92 (94)	53 (94)	94 (94)	92 (94)	49 (92)	49 (49)	47 (51)	31 (46)	56(61)	33(34)
DW	93 (94)	92 (94)	53 (93)	93 (93)	92 (94)	49 (92)	48 (48)	47 (50)	32 (47)	51(62)	33(35)
RR	94 (95)	<b>95 (95)</b>	64 (72)	94 (95)	94 (95)	50 (56)	47 (59)	49 (51)	26 (29)	<b>63(65)</b>	<b>36(36)</b>
Ours (EU)	94 (94)	94 (94)	54 (94)	94 (94)	94 (94)	49 (92)	49 (49)	49 (51)	32 (46)	54(59)	34(36)
Ours ( $\beta = 1$ )	96 (97)	94 (95)	74 ( <b>95</b> )	95 (96)	<b>96 (97)</b>	65 (93)	61 ( <b>62</b> )	60 ( <b>62</b> )	35 ( <b>54</b> )	62( <b>76</b> )	35(36)
Ours ( $\beta = 1.2$ )	94 (95)	<b>95 (95)</b>	<b>75 (95)</b>	96 (96)	<b>96 (97)</b>	65 (93)	61 ( <b>62</b> )	60 ( <b>62</b> )	35 ( <b>54</b> )	62(75)	34( <b>37</b> )
Ours ( $\beta = 1.5$ )	<b>97 (97)</b>	<b>95 (95)</b>	<b>75 (95)</b>	<b>96 (97)</b>	94 (95)	65 (93)	61 ( <b>62</b> )	59 ( <b>62</b> )	35 ( <b>54</b> )	62(74)	35( <b>37</b> )
Ours ( $\beta = 2$ )	96 (96)	<b>95 (96)</b>	73 (94)	<b>97 (97)</b>	95 (96)	<b>66 (95)</b>	<b>62 (62)</b>	<b>61 (62)</b>	<b>36 (54)</b>	62(75)	35(37)

and setting  $q_{i,t}$  for  $i \in \mathcal{N}$  in (5) to be proportional to the agents’ CIs. CFFL also utilizes the validation accuracy but is more efficient by using a leave-one-out approach instead of SV, while q-FFL aims at achieving egalitarian fairness by equalizing the local training losses of the agents. Furthermore, we implement simple FL baselines based on *round robin* (RR), *dataset weighted download* (DW), and *Euclidean distance* (EU). RR is commonly adopted in mechanism design to ensure fairness [6, 34] and also used in FL to schedule gradient downloads [51, 67]. For DW (EU),  $q_{i,t}$  for  $i \in \mathcal{N}$  in (5) are set to be proportional to the agents’ local dataset sizes (negative Euclidean distance of their unnormalized parameter updates from that of the server). We also include *standalone* agents as a baseline, i.e., each agent trains its CNN using only its local dataset without involving FL.

**Performance Metrics.** To measure fairness, we consider the *scaled Pearson correlation coefficient*<sup>3</sup>  $\rho := 100 \times \text{pearsonr}(\varphi, \xi) \in [-100, 100]$  between the test accuracies  $\varphi$  achieved by the agents when standalone [37] vs. that  $\xi$  achieved by them when collaborating via a gradient reward mechanism in FL after the entire training process has ended at iteration  $t = T$ . The corresponding experimental results will be reported in Sec. 4.2. To empirically verify the fairness guarantee in Theorem 2, we have also reported in Appendix B.3 results on the fairness metric  $\rho$  between the importance coefficients  $\varphi := (r_{i,T})_{i=1,\dots,N}$  (4) (i.e., measuring overall qualities of the parameter updates/gradients uploaded/contributed by the agents) vs. test accuracies (or negative training losses)  $\xi$  achieved by them. We consider other performance metrics like predictive performance (i.e., average and highest test accuracies achieved by the agents) and time overhead of the tested gradient reward mechanisms.

**Data Partitions among Agents.** We carefully construct two heterogeneous data partitions by varying the agents’ local dataset sizes and corresponding numbers of distinct classes. For **imbalanced dataset sizes** (POW), we follow a power law to partition the entire dataset among the agents. For MNIST, we partition the entire dataset of size  $\{3000, 6000, 12000\}$ , respectively, among  $\{5, 10, 20\}$  agents s.t. each agent has a randomly sampled local dataset of size 600 on average [40]. The size of the local dataset increases from the first to the last agent. Since the local dataset sizes vary significantly (i.e., superlinearly) among the agents, the agents with larger local datasets are expected to achieve better predictive performance. For **imbalanced class numbers** (CLA), we vary the number of distinct classes in the local datasets of the agents, while keeping their sizes fixed at 600. For this setting, we only consider MNIST and CIFAR-10 datasets and partition classes in a “linspace” manner as both contain 10 classes. To illustrate, for MNIST with 5 agents, agents 1, 2, 3, 4, 5 own local datasets with 1, 3, 5, 7, 10 classes, respectively; so, agent 1 (5) has a local dataset with 1 (10) class(es). Similarly, the agents with local datasets containing more classes are expected to achieve better predictive performance. We also include the simplest setting of the uniform/homogeneous data partition (UNI) where the agents are expected to achieve comparable predictive performance.

Additional details of the experimental settings are described in Appendix B.1.

## 4.2 Experimental Results

**Predictive Performance.** Table 1 shows results of the average and highest test accuracies achieved by the agents collaborating via our fair gradient reward mechanism vs. tested baselines on all

<sup>3</sup>The Pearson correlation coefficient has been applied to a similar use case in [19].

Table 2: Fairness metric  $\rho \in [-100, 100]$  achieved by our fair gradient reward mechanism with varying degrees of altruism  $\beta$  vs. tested baselines on all datasets. Higher value means greater fairness.

No. Agents	MNIST						CIFAR-10			MR	SST
	10			20			10			5	5
Data Partition	UNI	POW	CLA	UNI	POW	CLA	UNI	POW	CLA	POW	POW
FedAvg	-45.60	55.24	24.12	0.85	-32.58	40.83	18.47	97.48	98.75	48.68	57.50
q-FFL	-44.73	39.00	22.38	-22.01	38.71	48.07	-17.64	51.33	94.06	56.43	-75.92
CFFL	83.57	91.80	81.24	82.52	94.70	85.71	78.25	72.55	81.31	96.85	93.34
ECI	85.26	<b>99.83</b>	<b>99.98</b>	80.95	<b>99.41</b>	95.21	75.85	79.50	99.55	97.69	95.00
DW	89.15	98.93	65.34	86.94	99.63	35.21	-23.14	91.97	45.45	<b>99.20</b>	97.12
RR	83.77	71.17	-26.75	-18.64	25.47	95.86	30.67	0.70	90.67	44.16	-25.11
Ours (EU)	84.25	98.25	99.82	80.55	97.77	<b>99.97</b>	78.25	94.24	94.95	97.58	93.21
Ours ( $\beta = 1$ )	94.03	95.74	94.54	84.47	96.39	97.23	<b>98.80</b>	<b>98.78</b>	<b>99.89</b>	96.01	98.20
Ours ( $\beta = 1.2$ )	94.75	97.28	96.23	90.52	97.72	95.21	91.07	91.59	99.82	96.12	<b>98.47</b>
Ours ( $\beta = 1.5$ )	<b>96.34</b>	86.99	95.37	82.68	90.94	98.75	93.55	93.78	95.89	95.32	97.88
Ours ( $\beta = 2$ )	94.66	91.20	95.38	<b>96.90</b>	91.33	94.32	89.80	88.78	93.39	92.22	95.74

datasets. Our fair gradient reward mechanism generally outperforms the tested baselines on both metrics, especially for heterogeneous data partitions and on the MR dataset. On MNIST, for the CLA data partition among 10 agents, our fair gradient reward mechanism achieves average (highest) test accuracy of 75% (95%) at  $\beta = 1.5$ , while the best-performing ECI baseline achieves only that of 53% (94%). On CIFAR-10, for the CLA data partition among 10 agents, our fair gradient reward mechanism achieves average (highest) test accuracy of 36% (54%) at  $\beta = 2$ , while the best-performing DW baseline achieves only that of 32% (47%). On the MR dataset, our fair gradient reward mechanism achieves average (highest) test accuracy of 62% (76%) at  $\beta = 1$ , while the best-performing RR baseline achieves that of 63% (65%). Its better performance may be attributed to the adaptive re-weighting in the gradient aggregation step (1) via  $r_{i,t}$ , which can dynamically account for the heterogeneity in the agents’ local datasets [32]. While EU performs comparably to both FedAvg and ECI (i.e., difference in average test accuracies between EU vs. FedAvg/ECI is less than 3%), it does not perform better than our fair gradient reward mechanism (e.g., on MNIST, for the CLA data partition among 10 agents, the difference in average test accuracies between EU vs. our fair gradient reward mechanism at  $\beta = 1.5$  is more than 20%) because unlike cosine similarity, Euclidean distance fails to capture the directional difference between gradients, which is important since the negative gradients are pointing in the direction of lower loss. Importantly,  $q$ -FFL aims to equalize the local training losses w.r.t. the agent’s local datasets, which may be suboptimal for heterogeneous data partitions like POW and CLA. We provide further results in Appendix B.5 empirically comparing the predictive performances of our fair gradient reward mechanism vs.  $q$ -FFL.

**Fairness.** To empirically verify the fairness guarantee in Theorem 2, Table 2 shows results on the fairness metric  $\rho$  achieved by our fair gradient reward mechanism vs. tested baselines on all datasets. From Table 2, our fair gradient reward mechanism achieves a high degree of fairness of above 80, while the commonly used FedAvg performs suboptimally s.t. it produces the lowest degree of fairness of  $-45.6$ . On MNIST, for the POW data partition among 10/20 agents and the CLA data partition among 10 agents, ECI outperforms our fair gradient reward mechanism, albeit at a much higher time overhead of over 100 times and with additional information from an auxiliary dataset. CFFL underperforms our fair gradient reward mechanism and ECI as it adopts the leave-one-out approach which seems less accurate than SV in valuing the contributions of the agents [19]. Both  $q$ -FFL and RR promote egalitarian fairness instead of our notion of fairness via SV and hence do not perform optimally. DW achieves high degrees of fairness only for the POW data partition because it uses the agents’ local dataset sizes to determine their gradient rewards. Fig. 3 illustrates an intuitive trend of the predictive performances achieved by 10 agents collaborating via our fair gradient reward mechanism for homogeneous and heterogeneous data partitions among the agents on MNIST and CIFAR-10: For the UNI data partition, all agents achieve comparable predictive performance. Their predictive performances vary more (most) for the POW (CLA) data partition, hence demonstrating that our fair gradient reward mechanism can distinguish the contributions of the agents and reward them with sparsified gradients fairly.

We have performed an additional experiment to understand our fair gradient reward mechanism for homogeneous and heterogeneous data partitions among 3 agents on MNIST and CIFAR-10 where for POW and CLA, agent 1 (3) uploads/contributes parameter updates/gradients of lowest (highest) quality over the entire training process. Fig. 4 shows how  $r_{i,t}$  for agent  $i = 1, 3$  varies over iterations  $t$ : Interestingly, for the CLA data partition, though agent 3 (brown solid line) is initially mistaken to



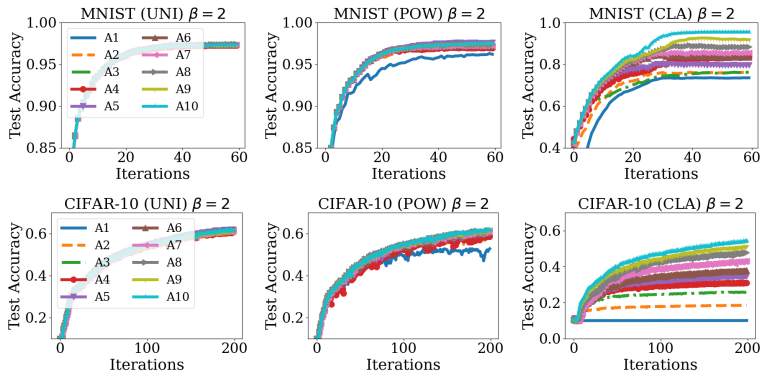


Figure 3: Test accuracy achieved by agent  $i = 1, \dots, 10$  (abbreviated to  $A_i$ ) collaborating via our fair gradient reward mechanism at  $\beta = 2$  for the UNI (left), POW (middle), and CLA (right) data partitions among the 10 agents on MNIST (top) and CIFAR-10 (bottom). Their predictive performances vary least, more, and most for the respective UNI, POW, and CLA data partitions.

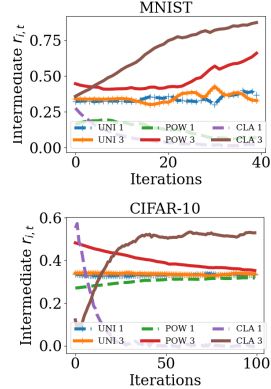


Figure 4: Graphs of  $r_{i,t}$  (4) for agent  $i = 1, 3$  vs. iteration  $t$  for UNI, POW, and CLA data partitions among 3 agents on MNIST and CIFAR-10.

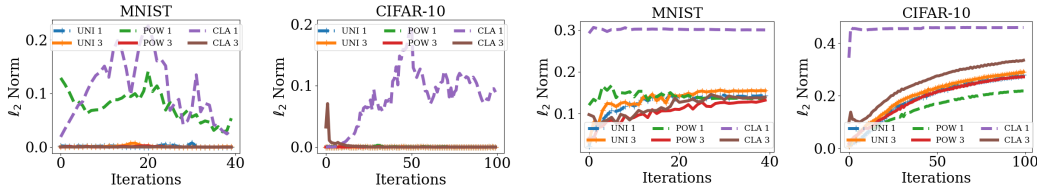


Figure 5: Graphs of  $\ell_2$  distance between down-loaded  $v_{i,t}$  (5) of agent  $i = 1, 3$  and aggregated layer's model parameters of agent  $i = 1, 3$   $u_{N,t}$  (1) vs. iteration  $t$  for UNI, POW, and CLA and that of the server vs. iteration  $t$  for UNI, POW, and CLA data partitions among 3 agents on MNIST (left) POW, and CLA data partitions among 3 agents on MNIST (left) and CIFAR-10 (right).

provide a low contribution, the dynamic update of  $r_{3,t}$  (4) allows its true contribution to be recognized quickly. Fig. 5 (Fig. 6) shows how the  $\ell_2$  distance between the downloaded sparsified gradient  $v_{i,t}$  (5) of agent  $i = 1, 3$  and aggregated parameter update/gradient  $u_{N,t}$  (1) (last layer's model parameters of agent  $i = 1, 3$  and that of the server) varies over iterations  $t$ : In particular, for the CLA data partition, agent  $i = 1$  ( $i = 3$ ) who uploads/contributes parameter updates/gradients of lowest (highest) quality over the entire training process downloads  $v_{i,t}$  as reward that is further from (closer to)  $u_{N,t}$ , hence training last layer's model parameters to be further from (closer to) that of the server. Such results further validate that in Fig. 2 previously.

Lastly, Fig. 7 confirms that for the CLA data partition among 10 agents on MNIST, increasing the degree of altruism  $\beta$  leads to all agents downloading higher-quality gradient rewards  $v_{i,t}$  (5) and thus incurring smaller training loss. In particular, agent 1 (abbreviated to A1 and represented by a blue solid line) who uploads/contributes parameter updates/gradients of lowest quality over the entire training process benefits most as  $\beta$  increases, as explained previously in Sec. 3.4. Additional results w.r.t. test loss are reported in Appendix B.4.

**Time Overhead.** Table 3 compares the time overhead (seconds) of our fair gradient reward mechanism vs. tested baselines on all datasets; the ratio between the time overhead vs. training time is given in brackets. Our fair gradient reward mechanism is much more efficient than ECI and CFFL which also consistently achieve fairness. In particular, our fair gradient reward mechanism incurs a small time overhead of at most  $0.4 \times$  of the training time, while ECI incurs a significant time overhead of up to  $140 \times$  of the training time due to the calculation of the CI incurring  $\mathcal{O}(2^N)$  time, even with the permutation sampling-based approximation [39, 56] for 10/20 agents. CFFL incurs at most  $2 \times$  of the training time (i.e., 5-6 times longer than ours) from the additional validation in each iteration.

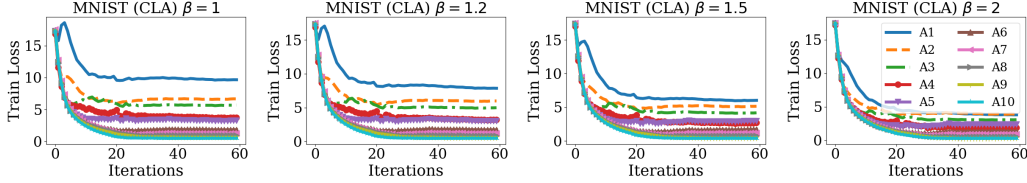


Figure 7: Training losses incurred by agent  $i = 1, \dots, 10$  (abbreviated to  $A_i$ ) collaborating via our fair gradient reward mechanism with varying degrees of altruism  $\beta = 1.0, 1.2, 1.5, 2$  for the CLA data partition on MNIST.

Table 3: Time overhead (seconds) of our fair gradient reward mechanism vs. tested baselines on all datasets. Each value in brackets denotes the ratio between the time overhead vs. training time.

No. Agents	MNIST			CIFAR-10		MR	SST
	5	10	20	5	10	5	5
FedAvg	1.17 (7e-3)	1.05 (1e-2)	4.29 (1e-2)	1.66 (7e-3)	7.41 (1e-2)	1.3 (1e-4)	1.31 (6e-4)
q-FFL	6.14 (4e-2)	4.97 (5e-2)	91.20 (0.3)	97.28 (0.4)	58.94 (7e-2)	90.01 (8e-3)	82.85 (4e-2)
CFFL	32.15 (0.2)	21.79 (0.3)	500.03 (1.6)	570.12 (2.0)	302.44 (0.4)	479.12 (0.2)	487.71 (2e-1)
ECI	2377.33 (16)	11937.80 (141)	23749.06 (74)	3571.75 (15)	58835.83 (84)	422.85 (4e-2)	801.20 (0.4)
DW	<b>0.89 (6e-3)</b>	<b>0.79 (9e-3)</b>	<b>1.60 (5e-3)</b>	<b>1.21 (5e-3)</b>	<b>5.29 (7e-3)</b>	<b>0.99 (1e-5)</b>	0.98 (5e-4)
RR	<b>0.89 (6e-3)</b>	0.82 (9e-3)	<b>1.60 (5e-3)</b>	3.31 (1e-2)	5.41 (7e-3)	1.01 (5e-4)	<b>0.99 (5e-4)</b>
Ours (EU)	<b>0.89 (6e-3)</b>	0.81 (9e-3)	1.61 (5e-3)	1.22 (5e-3)	5.33 (7e-3)	1.01 (5e-4)	<b>0.99 (5e-4)</b>
Ours (Cosine)	6.34 (4e-2)	4.94 (5e-2)	94.30 (0.3)	98.39 (0.4)	54.94 (7e-2)	89.81 (8e-3)	82.87 (4e-2)

**Hyperparameters.** We find that  $\alpha \in [0.8, 1]$  (i.e., relative weight on  $r_{i,t-1}$  in (4)),  $\beta \in [1, 2]$  (i.e., degree of altruism in (5)) and  $\Gamma \in [0.1, 1]$  (i.e., normalization coefficient in (1)) are effective in achieving competitive predictive performance and fairness. In our experiments, we set  $\alpha = 0.95$ ,  $\beta = [1, 1.2, 1.5, 2]$ , and  $\Gamma = 0.5$  for MNIST,  $\Gamma = 0.15$  for CIFAR-10, and  $\Gamma = 1$  for SST and MR.

## 5 Conclusion and Future Work

In this paper, we have described a novel *cosine gradient Shapley value* (CGSV) (Sec. 3.2) to fairly evaluate the expected marginal contribution of each agent’s uploaded model parameter update/gradient in FL without needing an auxiliary validation dataset and present an efficient approximation of CGSV with a bounded error (Sec. 3.3). Based on the approximate CGSV, we have designed a novel training-time fair gradient reward mechanism (Sec. 3.4) by exploiting the trick of sparsifying the aggregated parameter update/gradient downloaded from the server as reward to each agent such that its resulting quality is commensurate to that of the agent’s uploaded/contributed parameter update/gradient. Consequently, an agent who uploads/contributes higher-quality parameter updates/gradients over the entire training process should eventually be rewarded with converged model parameters whose resulting training loss (and hence predictive performance) is closer to that of the server, as demonstrated in our fairness guarantee (Sec. 3.5). We have empirically demonstrated the effectiveness of our fair gradient reward mechanism on multiple benchmark datasets in terms of fairness, predictive performance, and time overhead (Sec. 4). In particular, our fair gradient reward mechanism is much more efficient than several existing FL baselines since it requires only slight calculations by the server.

Our proposed fair gradient reward mechanism also provides practitioners the flexibility to trade off between fairness and equality in gradient rewards via a hyperparameter  $\beta$  controlling the degree of altruism (Sec. 3.4). For future work, it would be interesting to consider the notion of fairness when there are some adversaries. We would also consider generalizing our work and fairness guarantee to other types of CML (e.g., model fusion [16, 17, 24]) and collaborative Bayesian optimization [53].

## Acknowledgments and Disclosure of Funding

This research is supported by the National Research Foundation, Singapore under its AI Singapore Programme (Award No: AISG2-RP-2020-018). Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not reflect the views of National Research Foundation, Singapore. Xinyi Xu is supported by the Institute for Infocomm Research of Agency for Science, Technology and Research (A\*STAR).

## References

- [1] A. Agarwal, M. Dahleh, and T. Sarkar. A marketplace for data: An algorithmic solution. In *Proc. EC*, 2019.
- [2] D. Alistarh, T. Hoeffler, M. Johansson, S. Khirirat, N. Konstantinov, and C. Renggli. The convergence of sparsified gradient methods. In *Proc. NeurIPS*, 2018.
- [3] M. Bahir, M. And, B. Peleg, M. Maschler, and B. Peleg. A characterization, existence proof and dimension bounds for the kernel of a game. *Pacific Journal of Mathematics*, 18(2), 1966.
- [4] E. Balkanski and Y. Singer. Mechanisms for fair attribution. In *Proc. EC*, 2015.
- [5] G. Chalkiadakis, E. Elkind, and M. Wooldridge. Computational aspects of cooperative game theory. In R. J. Brachman, W. W. Cohen, and T. G. Dietterich, editors, *Synthesis Lectures on Artificial Intelligence and Machine Learning*. Morgan & Claypool Publishers, 2011.
- [6] C. Chen, W. Wang, and B. Li. Round-robin synchronization: Mitigating communication bottlenecks in parameter servers. In *Proc. IEEE INFOCOM*, 2019.
- [7] C. Y. Chen, J. Choi, D. Brand, A. Agrawal, W. Zhang, and K. Gopalakrishnan. ADaComp: Adaptive residual gradient compression for data-parallel distributed training. In *Proc. AAAI*, 2018.
- [8] M. Chen, B. Mao, and T. Ma. Fedsta: A staleness-aware asynchronous federated learning algorithm with non-iid data. *Future Generation Computer Systems*, 120:1–12, 2021.
- [9] Z. Chen, Z. Liu, K. L. Ng, H. Yu, Y. Liu, and Q. Yang. A gamified research tool for incentive mechanism design in federated learning. In Q. Yang, L. Fan, and H. Yu, editors, *Federated Learning*, volume 12500 of *Lecture Notes in Computer Science*, pages 168–175. Springer, Cham, 2020.
- [10] M. Cong, H. Yu, X. Weng, and S. Yiu. A game-theoretic framework for incentive mechanism design in federated learning. In Q. Yang, L. Fan, and H. Yu, editors, *Federated Learning*, volume 12500 of *Lecture Notes in Computer Science*, pages 205–222. Springer, Cham, 2020.
- [11] V. Conitzer and T. Sandholm. Computing shapley values, manipulating value division schemes, and checking core membership in multi-issue domains. In *Proc. AAAI*, 2004.
- [12] W. Dai, Y. Zhou, N. Dong, H. Zhang, and Eric P. Xing. Toward understanding the impact of staleness in distributed machine learning. In *Proc. ICLR*, 2019.
- [13] J. M. Drazen, S. Morrissey, D. Malina, M. B. Hamel, and E. W. Champion. The importance — and the complexities — of data sharing. *New England Journal of Medicine*, 375(12):1182–1183, 2016.
- [14] A. Ghorbani and J. Y. Zou. Data Shapley: Equitable valuation of data for machine learning. In *Proc. ICML*, 2019.
- [15] A. Hard, K. Rao, R. Mathews, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage. Federated learning for mobile keyboard prediction. arXiv:1811.03604, 2018.
- [16] Q. M. Hoang, T. N. Hoang, B. K. H. Low, and C. Kingsford. Collective model fusion for multiple black-box experts. In *Proc. ICML*, pages 2742–2750, 2019.
- [17] T. N. Hoang, C. T. Lam, B. K. H. Low, and P. Jaillet. Learning task-agnostic embedding of multiple black-box experts for multi-task model fusion. In *Proc. ICML*, pages 4282–4292, 2020.
- [18] Y. Hu, D. Niu, J. Yang, and S. Zhou. Fdml: A collaborative machine learning framework for distributed features. In *Proc. SIGKDD*, page 2232–2240, 2019.
- [19] R. Jia, D. Dao, B. Wang, F. A. Hubis, N. Hynes, N. M. Gürel, B. Li, C. Zhang, D. Song, and C. J. Spanos. Towards efficient data valuation based on the Shapley value. In *Proc. AISTATS*, 2019.

- [20] Y. Kim. Convolutional neural networks for sentence classification. In *Proc. EMNLP*, 2014.
- [21] A. Krizhevsky. Learning multiple layers of features from tiny images. Master’s thesis, Department of Computer Science, University of Toronto, 2009.
- [22] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In *Proc. NeurIPS*, 2012.
- [23] H. M. Krumholz and J. Waldstreicher. Toward fairness in data sharing. *New England Journal of Medicine*, 375(5):405–407, 2016.
- [24] C. T. Lam, T. N. Hoang, B. K. H. Low, and P. Jaillet. Model fusion for personalized learning. In *Proc. ICML*, pages 5948–5958, 2021.
- [25] Y. LeCun, B. Boser, J. Denker, D. Henderson, R. Howard, W. Hubbard, and L. Jackel. Hand-written digit recognition with a back-propagation network. In *Proc. NeurIPS*, 1990.
- [26] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 1998.
- [27] D. Leroy, A. Coucke, T. Lavril, T. Gisselbrecht, and J. Dureau. Federated learning for keyword spotting. In *Proc. ICASSP*, 2019.
- [28] H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein. Visualizing the loss landscape of neural nets. In *Proc. NeurIPS*, 2018.
- [29] Q. Li, Z. Wen, Z. Wu, S. Hu, N. Wang, Y. Li, X. Liu, and B. He. A survey on federated learning systems: Vision, hype and reality for data privacy and protection. *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [30] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020.
- [31] T. Li, M. Sanjabi, A. Beirami, and V. Smith. Fair resource allocation in federated learning. In *Proc. ICLR*, 2020.
- [32] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang. On the convergence of FedAvg on non-iid data. In *Proc. ICLR*, 2020.
- [33] Y. Lin, Y. Wang, S. Han, W. J. Dally, and H. Mao. Deep gradient compression: Reducing the communication bandwidth for distributed training. In *Proc. ICLR*, 2018.
- [34] R. Lipton, E. Markakis, E. Mossel, and A. Saberi. On approximately fair allocations of indivisible goods. In *Proc. EC*, 2004.
- [35] E. Lorch. Visualizing Deep Network Training Trajectories with PCA. In *Proc. ICML*, 2016.
- [36] L. Lyu, Y. Li, K. Nandakumar, J. Yu, and X. Ma. How to democratise and protect AI: Fair and differentially private decentralised deep learning. *IEEE Transactions on Dependable and Secure Computing*, 2020.
- [37] L. Lyu, X. Xu, Q. Wang, and H. Yu. Collaborative fairness in federated learning. In Q. Yang, L. Fan, and H. Yu, editors, *Federated Learning*, volume 12500 of *Lecture Notes in Computer Science*, pages 189–204. Springer, Cham, 2020.
- [38] L. Lyu, J. Yu, K. Nandakumar, Y. Li, X. Ma, J. Jin, H. Yu, and K. S. Ng. Towards fair and privacy-preserving federated deep models. *IEEE Transactions on Parallel and Distributed Systems*, 31(11):2524–2541, 2020.
- [39] S. Maleki, L. Tran-Thanh, G. Hines, T. Rahwan, and A. Rogers. Bounding the estimation error of sampling-based Shapley value approximation with/without stratifying. arXiv:1306.4265, 2013.
- [40] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Proc. AISTATS*, 2017.

- [41] B. McMahan and D. Ramage. Federated learning: Collaborative machine learning without centralized training data. <https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>. Accessed: 2021-10-08.
- [42] T. P. Michalak, P. L. Szczepański, T. Rahwan, A. Chrobak, S. Brânzei, M. Wooldridge, and N. R. Jennings. Implementation and computation of a value for generalized characteristic function games. *ACM Transactions on Economics and Computation*, 2(4), 2014.
- [43] D. Ng, X. Lan, M. M.-S. Yao, W. P. Chan, and M. Feng. Federated learning: a collaborative effort to achieve better medical imaging models for individual sites that have small labelled datasets. *Quantitative Imaging in Medicine and Surgery*, 11(2), 2020.
- [44] B. Pang and L. Lee. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proc. ACL*, 2005.
- [45] R. Pascanu, T. Mikolov, and Y. Bengio. On the difficulty of training recurrent neural networks. In *Proc. ICML*, 2013.
- [46] R. Raskar, P. Vepakomma, T. Swedish, and A. Sharan. Data markets to support AI for all: Pricing, valuation and governance. arXiv:1905.06462, 2019.
- [47] A. Richardson, A. Filos-Ratsikas, and B. Faltings. Budget-bounded incentives for federated learning. In Q. Yang, L. Fan, and H. Yu, editors, *Federated Learning*, volume 12500 of *Lecture Notes in Computer Science*, pages 176–188. Springer, Cham, 2020.
- [48] A. Richardson, A. Filos-Ratsikas, and B. Faltings. Incentivizing and rewarding high-quality data via influence functions. In *Proc. ICML Workshop on Incentives in Machine Learning*, 2020.
- [49] N. Rieke, J. Hancox, W. Li, F. Milletari, H. R. Roth, S. Albarqouni, S. Bakas, M. N. Galtier, B. A. Landman, K. Maier-Hein, S. Ourselin, M. Sheller, R. M. Summers, A. Trask, D. Xu, M. Baust, and M. J. Cardoso. The future of digital health with federated learning. *npj Digital Medicine*, 3(119), 2020.
- [50] L. S. Shapley. A value for  $n$ -person games. In H. W. Kuhn and A. W. Tucker, editors, *Contributions to the Theory of Games*, volume 2, pages 307–317. Princeton Univ. Press, 1953.
- [51] R. Shokri and V. Shmatikov. Privacy-preserving deep learning. In *Proc. ACM SIGSAC*, 2015.
- [52] R. H. L. Sim, Y. Zhang, M. C. Chan, and B. K. H. Low. Collaborative machine learning with incentive-aware model rewards. In *Proc. ICML*, 2020.
- [53] R. H. L. Sim, Y. Zhang, B. K. H. Low, and P. Jaillet. Collaborative Bayesian optimization with fair regret. In *Proc. ICML*, pages 9691–9701, 2021.
- [54] T. Song, Y. Tong, and S. Wei. Profit allocation for federated learning. In *Proc. IEEE International Conference on Big Data*, 2019.
- [55] S. Tay, X. Xu, C. S. Foo, and B. K. H. Low. Incentivizing collaboration in machine learning via synthetic data rewards. In *Proc. AAAI*, 2022.
- [56] T. Wang, J. Rausch, C. Zhang, R. Jia, and D. Song. A principled approach to data valuation for federated learning. In Q. Yang, L. Fan, and H. Yu, editors, *Federated Learning*, volume 12500 of *Lecture Notes in Computer Science*, pages 153–167. Springer, Cham, 2020.
- [57] K. Wei, J. Li, C. Ma, M. Ding, and H. V. Poor. Differentially private federated learning: Algorithm, analysis and optimization. In *Federated Learning Systems: Towards Next-Generation AI*, pages 51–78. Springer International Publishing, Cham, 2021.
- [58] E. Winter. The Shapley value. In R. Aumann and S. Hart, editors, *Handbook of Game Theory with Economic Applications*, volume 3, chapter 53, pages 2025–2054. Elsevier B.V., 2002.
- [59] X. Xu and L. Lyu. A reputation mechanism is all you need: Collaborative fairness and adversarial robustness in federated learning. In *Proc. ICML Workshop on Federated Learning for User Privacy and Data Confidentiality*, 2021.

- [60] X. Xu, Z. Wu, C. S. Foo, and B. K. H. Low. Validation free and replication robust volume-based data valuation. In *Proc. NeurIPS*, 2021.
- [61] Z. Yan, D. Xiao, M. Chen, J. Zhou, and W. Wu. Dual-way gradient sparsification for asynchronous distributed deep learning. In *Proc. ICPP*, 2020.
- [62] Q. Yang, Y. Liu, T. Chen, and Y. Tong. Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology*, 10(2), 2019.
- [63] Q. Yang, Y. Liu, Y. Cheng, Y. Kang, T. Chen, and H. Yu. *Federated Learning*. Morgan & Claypool Publishers, 2019.
- [64] H. P. Young. Monotonic solutions of cooperative games. *International Journal of Game Theory*, 14(2):65–72, 1985.
- [65] H. Yu, Z. Liu, Y. Liu, T. Chen, M. Cong, X. Weng, D. Niyato, and Q. Yang. A fairness-aware incentive scheme for federated learning. In *Proc. AIES*, 2020.
- [66] J. Zhang, Y. Wu, and R. Pan. Incentive mechanism for horizontal federated learning based on reputation and reverse auction. In *Proc. TheWebConf*, page 947–956, 2021.
- [67] S. Zhang, A. Choromanska, and Y. LeCun. Deep learning with elastic averaging SGD. In *Proc. NeurIPS*, 2015.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]** Section 3 formalizes our fairness framework and solution while Section 4 provides the theoretical guarantees.
  - (b) Did you describe the limitations of your work? **[Yes]** See Section 4 where we explicitly point out that the noise in estimation and generalizability are potential limitations.
  - (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? **[Yes]** The assumptions are stated and interpreted when the results are presented.
  - (b) Did you include complete proofs of all theoretical results? **[Yes]** See Appendix A for complete proof of all theoretical results.
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[Yes]**
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]** See Section 5 for training details and Appendix B for additional information on hyperparameters.
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[N/A]**
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]** See Appendix B for computational resources used.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? **[Yes]** See Section 5 under the paragraph for Datasets.
  - (b) Did you mention the license of the assets? **[N/A]**
  - (c) Did you include any new assets either in the supplemental material or as a URL? **[No]**

- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

## A Theoretical Results

### A.1 Fairness Properties of CGSV

For  $\mathcal{S} \subseteq \mathcal{N} \setminus \{i\}$ , let  $\Delta_{\mathcal{S},i} := \nu(\mathcal{S} \cup i) - \nu(\mathcal{S})$ , the following properties are satisfied by CGSV:

- (P1) Null Player [58]:  $\Delta_{\mathcal{S},i} = 0, \forall \mathcal{S} \subseteq \mathcal{N} \setminus \{i\} \implies \phi_i = 0$ .
- (P2) Symmetry [58]:  $\Delta_{\mathcal{S},i} = \Delta_{\mathcal{S},i'}, \forall \mathcal{S} \subseteq \mathcal{N} \setminus \{i, i'\} \implies \phi_i = \phi_{i'}$ .
- (P3) Strict Desirability [3]:  $\Delta_{\mathcal{S},i} \geq \Delta_{\mathcal{S},i'}, \forall \mathcal{S} \subseteq \mathcal{N} \setminus \{i, i'\}$  and  $\exists \mathcal{S}' \subseteq \mathcal{N} \setminus \{i, i'\}$  s.t.  $\Delta_{\mathcal{S}',i} > \Delta_{\mathcal{S}',i'} \implies \phi_i > \phi_{i'}$ .
- (P4) Coalitional Monotonicity [64, Equ.(4)]:  $\nu(\mathcal{S}) \geq \nu'(\mathcal{S})$  for some  $\mathcal{S} \subseteq \mathcal{N}$  and  $\nu(\mathcal{S}') = \nu'(\mathcal{S}') \forall \mathcal{S}' \subseteq \mathcal{N} \mathcal{S}' \neq \mathcal{S} \implies \phi_i(\nu) \geq \phi_i(\nu') \forall i \in \mathcal{S}$ .
- (P5) Individual Monotonicity [64, Equ.(5)]:  $\forall i, \nu(), \nu'()$ ,  $\nu(\mathcal{S}) \geq \nu'(\mathcal{S}) \forall \mathcal{S}$  containing  $i$  and  $\nu(\mathcal{S}') = \nu'(\mathcal{S}') \forall \mathcal{S}'$  not containing  $i \implies \phi_i(\nu) \geq \phi_i(\nu')$ .

(P1) can be intuitively understood as if  $i$  adds zero value to the group, then the corresponding CGSV will be zero. This is to prevent agents who wish to exploit the system by uploading randomly generated gradients instead of actual gradients. Note that in a high-dimensional space, the cosine similarity between a random gradient and an actual gradient is likely to be close to zero.

(P2) and (P3) provide a comparative relationship between any pair of agents  $i, i'$ . In the simplest case as in (P2),  $i, i'$  provide exactly identical contributions, then their corresponding CGSVs are equal. On the other hand, if  $i$  consistently provides more than  $i'$  as in (P3), then the CGSV for  $i$  is higher to correctly reflect this relation. Therefore, these two properties ensure the agents who contribute more by uploading better gradients are properly recognized (with higher CGSV). This is crucial in our designed reward mechanism which follows such relations in the CGSV.

(P4) states that if a group of agents in  $\mathcal{S}$  collectively do better, while all other groups  $\mathcal{S}' \neq \mathcal{S}$  stay the same, then the agents in  $\mathcal{S}$  do not lose. In particular, applying (P4) repeatedly gives an equivalent result regarding a single agent  $i$ , which we call *individual monotonicity* as in (P5).

(P5) takes the perspective of agent  $i$  while all other agents do not change, and agent  $i$  makes better contributions and improves (or at least does not hurt) all the coalitions  $i$  is in, then agent  $i$  does not lose. Consequently, it implies an incentive for the agents to make better contributions which could increase their CGSVs, which in turn correspond to better rewards in our mechanism.

### A.2 Proof of Theorem 1

Theorem 1 relies on the following equivalent formulation of  $\phi_i$ ,

$$\phi_i = \underbrace{\sum_{\mathcal{S} \subseteq \mathcal{N} \setminus \{i\}} A_{\mathcal{S}} \nu(\mathcal{S})}_{\text{additive error}} + \underbrace{[\sum_{\mathcal{S} \subseteq \mathcal{N} \setminus \{i\}} B_{\mathcal{S}}]}_{\text{multiplicative factor } L_i} \psi_i \quad (6)$$

where  $A_{\mathcal{S}}, B_{\mathcal{S}}$  are constants specific to  $\mathcal{S}$ . Note that Theorem 1 provides an upper bound for the additive error in (6) so we approximate  $\phi_i \approx L_i \psi_i$ . Further, by recalling a property of CGSV (invari-

ance under linear transformation), we can avoid explicitly calculating the multiplicative factor  $L_i$  via normalization of  $\psi_i$  if all  $L_i$  are approximately equal, as we show in Lemma 1.

$$\phi_i = \underbrace{\sum_{S \subseteq \mathcal{N} \setminus \{i\}} A_S \nu(S)}_{\text{additive error}} + \underbrace{\left[ \sum_{S \subseteq \mathcal{N} \setminus \{i\}} B_S \right]}_{\text{multiplicative factor } L_i} \psi_i$$

Before the proof, we first show the derivation (reproduced above). The intuitive idea is that the approximation  $\psi_i := \cos(\mathbf{u}_i, \mathbf{u}_{\mathcal{N}})$  appears in the summation of  $\phi_i$  repeatedly, so we collect all its coefficients into the multiplicative factor  $L_i$  and collect everything else as the additive error.

$$\begin{aligned} \phi_i &= \frac{1}{N} \sum_{S \subseteq \mathcal{N} \setminus \{i\}} \frac{1}{\binom{N-1}{|S|}} \nu(S \cup \{i\}) - \nu(S) \\ &= \frac{1}{N} \sum_{S \subseteq \mathcal{N} \setminus \{i\}} \frac{1}{\binom{N-1}{|S|}} \underbrace{(\cos(\mathbf{u}_{S \cup \{i\}}, \mathbf{u}_{\mathcal{N}}) - \cos(\mathbf{u}_S, \mathbf{u}_{\mathcal{N}}))}_A \end{aligned}$$

We focus on  $A$  and leave the rest unchanged.

$$\begin{aligned} A &= \frac{\langle \mathbf{u}_{S \cup \{i\}}, \mathbf{u}_{\mathcal{N}} \rangle}{\|\mathbf{u}_{S \cup \{i\}}\| \times \|\mathbf{u}_{\mathcal{N}}\|} - \frac{\langle \mathbf{u}_S, \mathbf{u}_{\mathcal{N}} \rangle}{\|\mathbf{u}_S\| \times \|\mathbf{u}_{\mathcal{N}}\|} \\ &= \frac{\langle \mathbf{u}_{S \cup \{i\}}, \mathbf{u}_{\mathcal{N}} \rangle}{\Gamma_1 \Gamma_{\mathcal{N}}} - \frac{\langle \mathbf{u}_S, \mathbf{u}_{\mathcal{N}} \rangle}{\Gamma_2 \Gamma_{\mathcal{N}}} \\ &= \frac{1}{\Gamma_1 \Gamma_2 \Gamma_{\mathcal{N}}} \langle \Gamma_2 \mathbf{u}_{S \cup \{i\}} - \Gamma_1 \mathbf{u}_S, \mathbf{u}_{\mathcal{N}} \rangle \\ &= \frac{(\Gamma_2 - \Gamma_1) \langle \mathbf{u}_S, \mathbf{u}_{\mathcal{N}} \rangle}{\Gamma_1 \Gamma_2 \Gamma_{\mathcal{N}}} - \frac{r_i \Gamma_2 \langle \mathbf{u}_i, \mathbf{u}_{\mathcal{N}} \rangle}{\Gamma_1 \Gamma_2 \Gamma_{\mathcal{N}}} \\ &= \frac{\Gamma_2 - \Gamma_1}{\Gamma_1} \nu(S) - \frac{r_i \Gamma_2}{\Gamma_1} \cos(\mathbf{u}_i, \mathbf{u}_{\mathcal{N}}) \end{aligned}$$

where  $\Gamma_1 := \|\mathbf{u}_{S \cup \{i\}}\|$ ,  $\Gamma_2 := \|\mathbf{u}_S\|$  and  $\Gamma_{\mathcal{N}} = \|\mathbf{u}_{\mathcal{N}}\|$ . Substitute this back with  $A_S = \frac{1}{N} \frac{1}{\binom{N-1}{|S|}} \frac{\Gamma_2 - \Gamma_1}{\Gamma_1}$  and  $B_S = \frac{1}{N} \frac{1}{\binom{N-1}{|S|}} \frac{r_i \Gamma_2}{\Gamma_1}$  to complete the derivation.

*Proof Sketch of Theorem 1.* The high-level idea is that, with cosine similarity, the approximation  $\psi_i := \cos(\mathbf{u}_i, \mathbf{u}_{\mathcal{N}})$  is a significant component of the actual CGSV  $\phi_i$  by a multiplicative factor. Because the multiplicative coefficient  $\frac{1}{\binom{N-1}{|S|}}$  becomes very small with a large  $N$ , so it reduces the effect of the terms *not* involving  $\mathbf{u}_i$ . While it also reduces the effect of the terms involving  $\mathbf{u}_i$ , the idea is that if the actual contribution from  $\mathbf{u}_i$  is relatively large (by the assumption  $|\langle r_i \mathbf{u}_i, \mathbf{u}_{\mathcal{N}} \rangle| \geq \frac{1}{T}$ ), then we ensure the error is small relatively. Note in the theorem, we have absorbed  $r_i$  into the constant  $\frac{1}{T}$ .  $\square$

*Proof of Theorem 1.* Notice the summation enumerates the same list of terms for both, so we minimize  $A_S \nu(S)$  relative to  $B_S \psi_i$ . Specifically, we examine the pairwise ratio between the two corresponding terms  $\frac{\Gamma_2 - \Gamma_1}{\Gamma_1} \cos(\mathbf{u}_S, \mathbf{u}_{\mathcal{N}})$  and  $\frac{r_i \Gamma_2}{\Gamma_1} \cos(\mathbf{u}_i, \mathbf{u}_{\mathcal{N}})$  in the summation as follows:

$$\begin{aligned} \frac{(|\Gamma_2 - \Gamma_1| \cos(\mathbf{u}_S, \mathbf{u}_{\mathcal{N}}))}{|r_i \Gamma_2 \cos(\mathbf{u}_i, \mathbf{u}_{\mathcal{N}})|} &= \frac{|\Gamma_2 - \Gamma_1|}{|\Gamma_2|} \frac{|\langle \mathbf{u}_S, \mathbf{u}_{\mathcal{N}} \rangle|}{|r_i \langle \mathbf{u}_i, \mathbf{u}_{\mathcal{N}} \rangle|} \\ &= |\Gamma_2 - \Gamma_1| \frac{|\langle \mathbf{u}_S, \mathbf{u}_{\mathcal{N}} \rangle|}{|\Gamma_2|} \frac{1}{|r_i \langle \mathbf{u}_i, \mathbf{u}_{\mathcal{N}} \rangle|} \\ &\leq \Gamma \sqrt{\frac{\|\mathbf{u}_S\|^2 \|\mathbf{u}_{\mathcal{N}}\|^2}{\Gamma_2^2}} \frac{1}{r_i |\langle \mathbf{u}_i, \mathbf{u}_{\mathcal{N}} \rangle|} \\ &\leq \Gamma \Gamma_{\mathcal{N}} \frac{1}{r_i |\langle \mathbf{u}_i, \mathbf{u}_{\mathcal{N}} \rangle|} \\ &\leq T \Gamma^2 \end{aligned}$$



We bound  $|\Gamma_2 - \Gamma_1| \leq \Gamma$  with the gradient normalization constant  $\Gamma$  by triangle inequality.

We bound  $\langle \mathbf{u}_S, \mathbf{u}_N \rangle / \Gamma_2$  using Cauchy-Schwarz inequality, bound  $\Gamma_N$  with  $\Gamma$  as  $\mathbf{u}_N$  is a convex sum of vectors each with norm  $\Gamma$ , and use the assumption to bound  $1/r_i |\langle \mathbf{u}_i, \mathbf{u}_N \rangle|$ .

The above inequality bounds error-to-approximation ratio, i.e.,  $|A_S \nu(S)| / |B_S \psi_i|$  is bounded by  $I\Gamma^2$  for every coalition  $\mathcal{S}$  in the summation, which implies

$$\frac{|\sum_{\mathcal{S} \subseteq \mathcal{N} \setminus \{i\}} A_S \nu(S)|}{|\sum_{\mathcal{S} \subseteq \mathcal{N} \setminus \{i\}} B_S \psi_i|} \leq I\Gamma^2.$$

This result is useful because the error-to-approximation ratio is consistently bounded regardless of the normalization on  $L_i \psi_i$ , such as linear scaling we conduct subsequently.

If we additionally assume,

$$\frac{r_i \Gamma}{\Gamma_1} = \frac{r_i \|\mathbf{u}_i\|}{\|\sum_{i' \in \mathcal{S} \cup \{i\}} r_{i'} \mathbf{u}_{i'}\|} \leq 1,$$

then the error term  $\phi_i - L_i \psi_i \leq I\Gamma^2$  before normalization. Its proof is by showing  $|\sum_{\mathcal{S} \subseteq \mathcal{N} \setminus \{i\}} B_S \psi_i \leq 1|$  and rearranging the terms in the previous inequality.

Note  $|\psi_i| \leq 1$  and by the assumption  $r_i \Gamma / \Gamma_1 \leq 1$ , we first have  $|r_i \Gamma \psi_i / \Gamma_1| \leq 1$ , so  $\sum_{\mathcal{S} \subseteq \mathcal{N} \setminus \{i\}} r_i \Gamma \psi_i / \Gamma_1 \leq \sum_{\mathcal{S} \subseteq \mathcal{N} \setminus \{i\}} 1 = \text{number of terms in the summation}$ . Putting back in the coefficients we can show that  $\sum_{\mathcal{S} \subseteq \mathcal{N} \setminus \{i\}} 1/N \times 1/\binom{N-1}{|\mathcal{S}|} < 1$  because the enumeration  $\mathcal{S} \subseteq \mathcal{N} \setminus \{i\}$  is not exhaustive while the coefficients are specified to have sum is 1 when the enumeration is exhaustive.

This additional assumption excludes degenerate cases where multiple agents (with approximately equal  $r_i$ 's) upload gradients in opposite directions and counteract each other which results in a net gradient vector approximately equal to a zero vector. Such cases are unlikely as the gradient vectors are calculated based on randomly selected mini-batches, and these gradient vectors are in a high dimension. □

In order to recall the property that CGSV is invariant under linear transformation, we require that all  $L_i$ 's are approximately equal. To show this, we specify an assumption to exclude the degenerate cases by requiring  $\mathbf{u}_{\mathcal{S} \cup \{i\}}$  and  $\mathbf{u}_{\mathcal{S} \cup \{j\}}$  are lower bounded linearly in  $\Gamma$ . This assumption stipulates that  $\mathbf{u}_{\mathcal{S} \cup \{i\}}$  and  $\mathbf{u}_{\mathcal{S} \cup \{j\}}$  are away from zero vectors, and have norms of the same magnitude of  $\mathbf{u}_i$  and  $\mathbf{u}_j$ .

**Lemma 1 (Closeness of  $L_i$ ).** Assume  $\exists M > 0$ , s.t.  $\forall \mathcal{S} \subseteq \mathcal{N} \setminus \{i, i'\}$ ,  $\min(\|\mathbf{u}_{\mathcal{S} \cup \{i\}}\|, \|\mathbf{u}_{\mathcal{S} \cup \{i'\}}\|) \geq M\Gamma$ , then

$$\max_{i, i' \in \mathcal{N}} L_i - L_{i'} \leq \sum_{s \in \mathcal{N} \setminus \{i, i'\}} \frac{1}{\binom{N-1}{|s|}} \frac{2}{M^2 \Gamma}.$$

*Proof of Lemma 1.* Due to symmetry,  $\mathbf{u}_i = \mathbf{u}_{i'} \implies L_i = L_{i'}$ . We only need to consider  $\mathbf{u}_i \neq \mathbf{u}_{i'}$ .

We consider the terms by grouping the coalitions  $\mathcal{S}$  into three types: 1)  $i \notin \mathcal{S} \wedge i' \notin \mathcal{S}$ ; 2)  $i \in \mathcal{S} \oplus i' \in \mathcal{S}$ ; 3)  $i \in \mathcal{S} \wedge i' \in \mathcal{S}$ . We need not consider 3) as the summation for  $i$  is over  $\mathcal{S} \subseteq \mathcal{N} \setminus \{i\}$ .

For 2), let  $\mathcal{S} \subseteq \mathcal{N} \setminus \{i, i'\}$  and  $\mathcal{S}_i = \mathcal{S} \cup \{i\}$ ,  $\mathcal{S}_{i'} = \mathcal{S} \cup \{i'\}$ . We can see that  $\mathcal{S}_i, \mathcal{S}_{i'}$  constitute a pair of symmetric case for two terms in the summation of  $L_i$  and  $L_{i'}$  respectively. In particular, for  $L_i$ , the term in summation is  $1/\binom{N-1}{|\mathcal{S}_{i'}|} \times 1/\|\mathbf{u}_{\mathcal{S}_{i'} \cup \{i\}}\|$ . Since  $\mathbf{u}_{\mathcal{S}_i \cup \{i'\}} = \mathbf{u}_{\mathcal{S} \cup \{i, i'\}} = \mathbf{u}_{\mathcal{S}_{i'} \cup \{i\}}$  and  $|\mathcal{S}_i| = |\mathcal{S}| + 1 = |\mathcal{S}_{i'}|$ , these two symmetric terms cancel out and the 2) type coalitions contribute exactly 0 to  $L_i - L_{i'}$ .

Now for 1) the coalitions  $\mathcal{S} \subseteq \mathcal{N} \setminus \{i, i'\}$ , we bound the sum of terms as follows,

$$\begin{aligned}
L_i - L_j &= \sum_{\mathcal{S} \subseteq \mathcal{N} \setminus \{i, i'\}} \frac{1}{\binom{N-1}{|\mathcal{S}|}} \left[ \frac{1}{\|\mathbf{u}_{\mathcal{S} \cup \{i\}}\|} - \frac{1}{\|\mathbf{u}_{\mathcal{S} \cup \{i'\}}\|} \right] \\
&= \sum_{\mathcal{S} \subseteq \mathcal{N} \setminus \{i, i'\}} \frac{1}{\binom{N-1}{|\mathcal{S}|}} \frac{\|\mathbf{u}_{\mathcal{S} \cup \{i'\}}\| - \|\mathbf{u}_{\mathcal{S} \cup \{i\}}\|}{\|\mathbf{u}_{\mathcal{S} \cup \{i\}}\| \times \|\mathbf{u}_{\mathcal{S} \cup \{i'\}}\|} \\
&\leq \sum_{\mathcal{S} \subseteq \mathcal{N} \setminus \{i, i'\}} \frac{1}{\binom{N-1}{|\mathcal{S}|}} \frac{2\Gamma}{M^2\Gamma^2} \\
&\leq \sum_{\mathcal{S} \subseteq \mathcal{N} \setminus \{i, i'\}} \frac{1}{\binom{N-1}{|\mathcal{S}|}} \frac{2}{M^2\Gamma}.
\end{aligned}$$

The first inequality is because the numerator is upper bounded by  $2\Gamma$  due to triangle inequality, and the denominator is lower bounded by  $M^2\Gamma^2$  using the assumption. This error decreases quickly with more agents due to the coefficient  $\binom{N-1}{|\mathcal{S}|}^{-1}$ .  $\square$

### A.3 Proof of Theorem 2

This proof includes an intermediate step of showing  $\delta_{i',t} \geq \delta_{i,t}$ . First observe the following inequalities using the triangle inequality:

$$\delta_{i,t} \leq \delta_{i,t-1} + \|\mathbf{v}_{i,t}\| \quad \text{and} \quad \delta_{i',t} \geq \delta_{i',t-1} - \|\mathbf{v}_{i',t}\|. \quad (7)$$

From the condition

$$\delta_{i',t-1} - \delta_{i,t-1} \geq 2\|\mathbf{v}_{i,t}\|,$$

we have

$$\delta_{i',t-1} - \delta_{i,t-1} \geq 2\|\mathbf{v}_{i,t}\| \geq \|\mathbf{v}_{i,t}\| + \|\mathbf{v}_{i',t}\| \quad (8)$$

The inequality  $\|\mathbf{v}_{i,t}\| \geq \|\mathbf{v}_{i',t}\|$  follows directly by applying  $r_{i,t} \geq r_{i',t}$  to (5) and observing  $\text{mask}(\cdot)$  retains the largest components in magnitude and making the rest zeros.

Rearranging (8) gives

$$\delta_{i',t-1} - \|\mathbf{v}_{i',t}\| \geq \delta_{i,t-1} + \|\mathbf{v}_{i,t}\|.$$

Connecting both inequalities in (7) gives

$$\delta_{i',t} \geq \delta_{i',t-1} - \|\mathbf{v}_{i',t}\| \geq \delta_{i,t-1} + \|\mathbf{v}_{i,t}\| \geq \delta_{i,t}.$$

Subsequently, we use  $\delta_{i',t} \geq \delta_{i,t}$  and some regularity conditions of  $\mathbf{F}(\cdot)$  to establish  $\mathbf{F}(\mathbf{w}_{i,t}) \leq \mathbf{F}(\mathbf{w}_{i',t})$ . Specifically, we assume  $\mathbf{F}(\cdot)$  is both  $L$ -smooth and  $\mu$ -strongly convex with  $L \leq \mu$ .

We first recall the respective definitions for  $\mu$ -strongly convex and  $L$ -smooth functions.

**Definition 2 ( $L$ -Smooth  $\mathbf{F}$ ).** If  $\mathbf{F}$  is  $L$ -smooth, then  $\forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}$ ,

$$\mathbf{F}(\mathbf{w}) \leq \mathbf{F}(\mathbf{w}') + \nabla \mathbf{F}(\mathbf{w}')^T (\mathbf{w} - \mathbf{w}') + \frac{L}{2} \|\mathbf{w} - \mathbf{w}'\|^2.$$

**Definition 3 ( $\mu$ -Strongly Convex  $\mathbf{F}$ ).** If  $\mathbf{F}$  is  $\mu$ -strongly convex, then  $\forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}$ ,

$$\mathbf{F}(\mathbf{w}) \geq \mathbf{F}(\mathbf{w}') + \nabla \mathbf{F}(\mathbf{w}')^T (\mathbf{w} - \mathbf{w}') + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}'\|^2.$$

From  $L$ -smoothness, we have

$$\mathbf{F}(\mathbf{w}_{i,t}) \leq \underbrace{\mathbf{F}(\mathbf{w}_{\mathcal{N},t}) + \nabla \mathbf{F}(\mathbf{w}_{\mathcal{N},t})^T (\mathbf{w}_{i,t} - \mathbf{w}_{\mathcal{N},t})}_{R_L} + \frac{L}{2} \delta_{i,t}^2.$$

From  $\mu$ -strong convexity, we have

$$\mathbf{F}(\mathbf{w}_{i',t}) \geq \underbrace{\mathbf{F}(\mathbf{w}_{\mathcal{N},t}) + \nabla \mathbf{F}(\mathbf{w}_{\mathcal{N},t})^T (\mathbf{w}_{i',t} - \mathbf{w}_{\mathcal{N},t})}_{R_\mu} + \frac{\mu}{2} \delta_{i',t}^2.$$

In order to prove  $\mathbf{F}(\mathbf{w}_{i,t}) \leq \mathbf{F}(\mathbf{w}_{i',t})$ , it suffices to prove  $R_L \leq R_\mu$  or equivalently  $R_L - R_\mu \leq 0$ .

$$R_L - R_\mu = \underbrace{\nabla \mathbf{F}(\mathbf{w}_{\mathcal{N},t})^\top (\mathbf{w}_{i,t} - \mathbf{w}_{i',t})}_{R_1} + \underbrace{\frac{1}{2}(L\delta_{i,t}^2 - \mu\delta_{i',t}^2)}_{R_2}.$$

With  $L \leq \mu$  and  $\delta_{i,t} \leq \delta_{i',t}$ , we have

$$R_2 = \frac{1}{2}(L\delta_{i,t}^2 - \mu\delta_{i',t}^2) \leq \frac{L}{2}(\delta_{i,t}^2 - \delta_{i',t}^2) \leq 0.$$

We formalize  $\mathbf{w}_{\mathcal{N},t}$  being near to a stationary point by specifying an upper bound on the gradient:

$$\|\nabla \mathbf{F}(\mathbf{w}_{\mathcal{N},t})\| \leq \frac{L|\delta_{i,t}^2 - \delta_{i',t}^2|}{2\|\mathbf{w}_{i,t} - \mathbf{w}_{i',t}\|}.$$

We have the following:

$$\begin{aligned} |R_1| &\triangleq |\nabla \mathbf{F}(\mathbf{w}_{\mathcal{N},t})^\top (\mathbf{w}_{i,t} - \mathbf{w}_{i',t})| \leq \|\nabla \mathbf{F}(\mathbf{w}_{\mathcal{N},t})\| \times \|\mathbf{w}_{i,t} - \mathbf{w}_{i',t}\| \\ &\leq \frac{L|\delta_{i,t}^2 - \delta_{i',t}^2|}{2} \\ &\leq |R_2| \end{aligned}$$

where the first inequality is by Cauchy-Schwarz, the second inequality is by substituting the above upper bound and the last inequality is due to taking absolute values of two negative values.

Finally, since  $|R_1| \leq |R_2|$  and  $R_2 \leq 0$ , we get  $R_1 + R_2 \leq 0$  and hence  $R_L + R_\mu \triangleq R_1 + R_2 \leq 0$ .

## B Experimental Results

### B.1 Experimental Settings

**Additional Details.** For CIFAR-10, we follow power law to randomly partition total {10000, 20000} examples among {5, 10} agents respectively. For MR (SST), we follow power law to randomly partition 9596 (8544) examples among 5 agents. We provide the training hyper-parameters used for different datasets in Table 4.

Table 4: Framework-independent hyper-parameters. Batch size  $B$ , initial step-size  $\eta$ , step-size exponential decay  $\gamma$ , total iterations  $T$ . Note for experiments with more than 5 agents for MNIST and CIFAR-10,  $\eta$  is 0.25 and 0.025, respectively.

Dataset	$B$	$\eta$ ( $\gamma$ )	$T$
MNIST	32	0.15 (0.977)	60
CIFAR-10	128	0.015 (0.977)	200
MR	128	5e-5 (0.977)	100
SST	256	1e-5 (0.977)	100

**Experiment Hardware and Software.** All experiments are conducted on a server with 16 cores (Intel(R) Xeon(R) CPU E5-2683 v4 @ 2.10GHz), 256 GB RAM and 4 GPUs (GeForce GTX 1080 Ti). Our implementation mainly uses PyTorch, torchtext, torchvision and some auxiliary packages such as Numpy, Pandas and Matplotlib. The specific versions and package requirements are provided together with the source code. To reduce the impact of randomness in the experiments, we adopt several measures: fix the model initializations (we initialize model weights and save them for future experiments); fix all the random seeds; and invoke the deterministic behavior of PyTorch. As a result, given the same model initialization, our implementation is expected to produce consistent results on the same machine over experimental runs.

### B.2 5-Agent Case for MNIST and CIFAR-10

For completion, we include the accuracy and fairness results under the consistent setting as the main paper for the 5-agent case for MNIST and CIFAR-10 for the three data partitions in Table 5 and Table 6, respectively.

Table 5: Average test accuracies (%) of all the agents for all baselines and our method with varying degrees of altruism  $\beta$ , on all four datasets. Values in the bracket denote the highest test accuracies among all the agents.

Data Partition	MNIST $N = 5$			CIFAR-10 $N = 5$		
	UNI	POW	CLA	UNI	POW	CLA
<i>Standalone</i>	91(91)	87(94)	50(91)	44(46)	42(52)	29(44)
<i>FedAvg</i>	93(93)	91(95)	50(92)	46(47)	46(52)	30(45)
<i>q-FFL</i>	82(85)	59(78)	49(84)	31(32)	31(34)	19(24)
<i>CFFL</i>	24(39)	21(37)	27(28)	44(45)	40(49)	26(43)
<i>ECI</i>	93(94)	94(95)	52(92)	46(47)	44(44)	30(41)
<i>DW</i>	93(93)	91(95)	50(92)	46(47)	46(52)	30(45)
<i>RR</i>	94(95)	95(95)	67(73)	40(45)	49(57)	23(32)
<i>Ours (EU)</i>	94(94)	93(95)	50(92)	47(48)	48(52)	30(45)
<i>Ours (<math>\beta = 1</math>)</i>	<b>96(97)</b>	<b>96(97)</b>	72(93)	<b>57(57)</b>	<b>56(57)</b>	<b>31(48)</b>
<i>Ours (<math>\beta = 1.2</math>)</i>	<b>96(97)</b>	<b>96(97)</b>	<b>73(94)</b>	<b>57(57)</b>	<b>56(57)</b>	<b>31(48)</b>
<i>Ours (<math>\beta = 1.5</math>)</i>	<b>96(97)</b>	<b>96(97)</b>	<b>76(94)</b>	<b>57(57)</b>	<b>56(57)</b>	<b>31(48)</b>
<i>Ours (<math>\beta = 2</math>)</i>	<b>97(97)</b>	<b>96(97)</b>	<b>79(94)</b>	<b>57(57)</b>	<b>56(58)</b>	<b>31(48)</b>

Table 6: Fairness metric  $\rho$  for all baselines and our method with varying degrees of altruism  $\beta$ , on MNIST and CIFAR-10 with 5 agents.  $\rho$  is computed between **standalone test accuracies** and **final test accuracies**. The higher the values, the better in terms of fairness in rewards.

Data Partition	MNIST $N = 5$			CIFAR-10 $N = 5$		
	UNI	POW	CLA	UNI	POW	CLA
<i>FedAvg</i>	-18.6	25.47	95.01	18.47	97.48	98.75
<i>q-FFL</i>	26.46	47.26	96.07	5.53	33.25	97.60
<i>CFFL</i>	30.76	18.06	-23.04	66.21	63.35	-13.94
<i>ECI</i>	37.18	62.13	96.41	85.43	97.86	98.45
<i>DW</i>	-33.1	99.35	12.11	80.64	99.17	99.90
<i>RR</i>	-47.5	94.84	81.36	74.43	-23.7	97.17
<i>Ours (EU)</i>	71.63	70.15	91.57	<b>96.36</b>	<b>99.71</b>	<b>99.91</b>
<i>Ours (<math>\beta = 1</math>)</i>	<b>83.53</b>	<b>99.57</b>	<b>98.62</b>	85.32	95.04	99.70
<i>Ours (<math>\beta = 1.2</math>)</i>	75.84	99.46	97.67	78.35	95.81	99.73
<i>Ours (<math>\beta = 1.5</math>)</i>	76.92	<b>99.57</b>	95.37	81.05	95.56	99.72
<i>Ours (<math>\beta = 2</math>)</i>	21.16	-33.99	97.72	99.22	99.89	99.97

### B.3 Fairness Comparison Based on CGSV

As these FL-based variants all use gradients as the communication medium, we can accordingly adapt our CGSV approximation and the moving averaging formulation as in (4). Specifically, we compute the moving average  $r_{i,t}$  for each framework respectively as the (cumulative) contribution of  $i$  and use it for fairness evaluation. We calculate the fairness metric  $\rho$  by considering two types of reward  $\xi$ 's: *final test accuracies* in Table 7, and *negative training losses* in Table 8. In comparison, the fairness results in Table 2 are computed between the standalone test accuracies and final test accuracies.

Table 7: Fairness metric  $\rho$  for all baselines and our method with varying degrees of altruism  $\beta$ , on all four datasets.  $\rho$  is computed between  $(r_{i,T})_{i=1,\dots,N}$  and **final test accuracies** where  $T$  denotes the last iteration as in Table 4. The higher the values, the better in terms of fairness in rewards.

No. Agents	MNIST						CIFAR-10			MR	SST
	10			20			10			5	5
Data Partition	UNI	POW	CLA	UNI	POW	CLA	UNI	POW	CLA	POW	POW
<i>FedAvg</i>	-10.54	45.13	33.89	77.80	62.19	55.45	56.21	45.77	77.18	45.54	-7.05
<i>q-FFL</i>	-86.23	-15.75	67.98	-54.10	-35.45	12.54	61.78	36.55	-58.45	49.53	-93.94
<i>CFFL</i>	85.89	38.64	15.80	-44.25	5.62	-16.38	41.76	-54.63	45.32	19.00	19.45
<i>ECI</i>	44.17	85.74	91.15	79.03	89.69	91.81	80.67	89.57	<b>97.30</b>	84.58	93.12
<i>DW</i>	0.22	88.64	-40.55	65.21	91.77	73.14	-10.25	87.25	9.27	73.53	72.07
<i>RR</i>	3.17	75.12	78.15	-3.33	86.43	91.88	43.04	-17.31	84.95	-15.5	-6.43
<i>Ours (EU)</i>	49.44	79.85	83.78	55.27	91.63	85.35	<b>89.89</b>	93.02	87.19	87.88	93.21
<i>Ours (<math>\beta = 1</math>)</i>	<b>90.49</b>	94.68	77.27	<b>93.20</b>	92.89	<b>94.58</b>	82.66	92.75	94.89	96.01	94.31
<i>Ours (<math>\beta = 1.2</math>)</i>	89.74	<b>96.28</b>	82.95	90.65	<b>96.99</b>	93.42	78.99	<b>93.41</b>	95.96	<b>96.12</b>	90.00
<i>Ours (<math>\beta = 1.5</math>)</i>	80.23	91.40	<b>93.60</b>	90.89	90.31	93.66	76.42	92.89	89.52	95.32	95.11
<i>Ours (<math>\beta = 2</math>)</i>	82.36	91.03	88.45	74.88	87.44	91.20	72.71	89.63	84.92	85.01	<b>97.49</b>

### B.4 Empirical Validation of Theorem 2 via Test Loss

In addition to the results from Appendix B.3, we perform experiments to further validate Theorem 2 by considering the *test* loss (instead of training loss) as the reward, i.e., a lower test loss corresponds

Table 8: Fairness metric  $\rho$  for all baselines and our method with varying degrees of altruism  $\beta$ , on all four datasets.  $\rho$  is computed between  $(r_{i,T})_{i=1,\dots,N}$  and **negative training losses** where  $T$  denotes the last iteration as in Table 4. The higher the values, the better in terms of fairness in rewards.

No. Agents	MNIST						CIFAR-10			MR	SST
	10			20			10			5	5
Data Partition	UNI	POW	CLA	UNI	POW	CLA	UNI	POW	CLA	POW	POW
FedAvg	68.68	86.73	96.01	83.05	87.35	84.03	55.84	83.57	95.33	87.26	73.53
q-FFL	48.04	60.71	48.68	5.78	30.44	-10.48	12.65	88.00	-55.00	<b>99.39</b>	94.39
CFFL	-52.81	-6.07	-57.21	-48.42	13.13	-7.90	3.55	-11.76	41.22	33.30	44.33
ECI	71.79	92.97	82.10	82.10	78.01	58.35	<b>84.90</b>	85.81	93.81	95.10	82.75
DW	68.12	<b>95.13</b>	95.59	59.05	72.13	85.23	49.31	90.17	95.36	88.87	73.12
RR	46.52	87.84	<b>96.65</b>	31.99	92.73	91.20	16.13	85.28	<b>97.02</b>	87.97	73.56
Ours (EU)	71.59	87.32	96.11	82.14	86.08	84.00	61.53	83.31	95.12	88.73	73.23
Ours ( $\beta = 1$ )	<b>91.64</b>	92.26	96.84	<b>89.47</b>	<b>94.78</b>	<b>96.82</b>	84.85	<b>94.59</b>	90.13	90.44	91.92
Ours ( $\beta = 1.2$ )	90.36	91.94	91.19	88.87	93.28	96.41	83.84	90.94	90.16	90.05	<b>97.54</b>
Ours ( $\beta = 1.5$ )	91.03	93.33	92.27	88.21	92.11	91.39	84.43	90.51	90.33	89.54	89.84
Ours ( $\beta = 2$ )	85.02	88.61	94.88	87.51	90.09	92.36	78.95	88.84	90.65	88.72	94.64

to a better reward. Figure 8 demonstrates the same consistent trend as with training losses. This demonstrates the generalizability of Theorem 2 that the agents who upload better gradients have receive better models (i.e., with lower test losses).

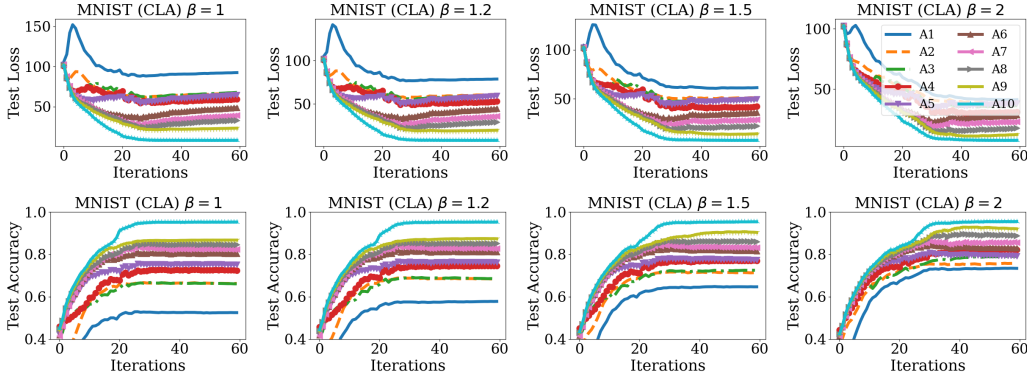


Figure 8: Test *losses* (first row) and test *accuracies* (second row) of the final models vs. the altruism degree  $\beta$  for MNIST CLA. From left to right,  $\beta = [1.0, 1.2, 1.5, 2]$ . A higher  $\beta$  leads to better performance for agents with lower contributions.

## B.5 Additional Comparison with q-FFL

Since q-FFL sets out to achieve a different notion of fairness than ours, we perform more in-depth comparison to examine the effects these two algorithms have on the agents’ final models. We plot the final performance in terms of test accuracy, test loss and train loss of all 5 agents for MNIST and CIFAR-10 under the three types of data partitions in Figures 9 and 10. And Figures 11 and 12 show the corresponding results for 10 agents.

We observe that in all scenarios, our algorithm performs noticeably better in terms of the final test accuracy. However, it may be due to that in q-FFL each agent is interested in performing well on their *own* local/private test sets (of the same distribution of their local train set). We do not investigate that use case. Instead, our scenario is that all the agent are interested in one common objective (represented by the same test set on which the test accuracy and test loss is computed).

Specifically, we observe from the second row of Figure 9, q-FFL ‘under’-optimizes agent 4 while our algorithm fairly evaluates and rewards all the agents (the increasing trend of test accuracy and decreasing trend of train loss). Moreover, we observe from the second row of Figure 10, q-FFL ‘equalizes’ the performance of the agents in terms of test accuracy and test loss, while our algorithm fairly rewards the agents (the increasing trend of test accuracy and decreasing trend of test and train losses).

In summary, despite the similarity in the keyword terminology, namely fairness, our algorithm is fundamentally different from q-FFL in that in our setting, all agents share one learning objective while in q-FFL each agent has their own learning objective (which may differ considerably from others’).

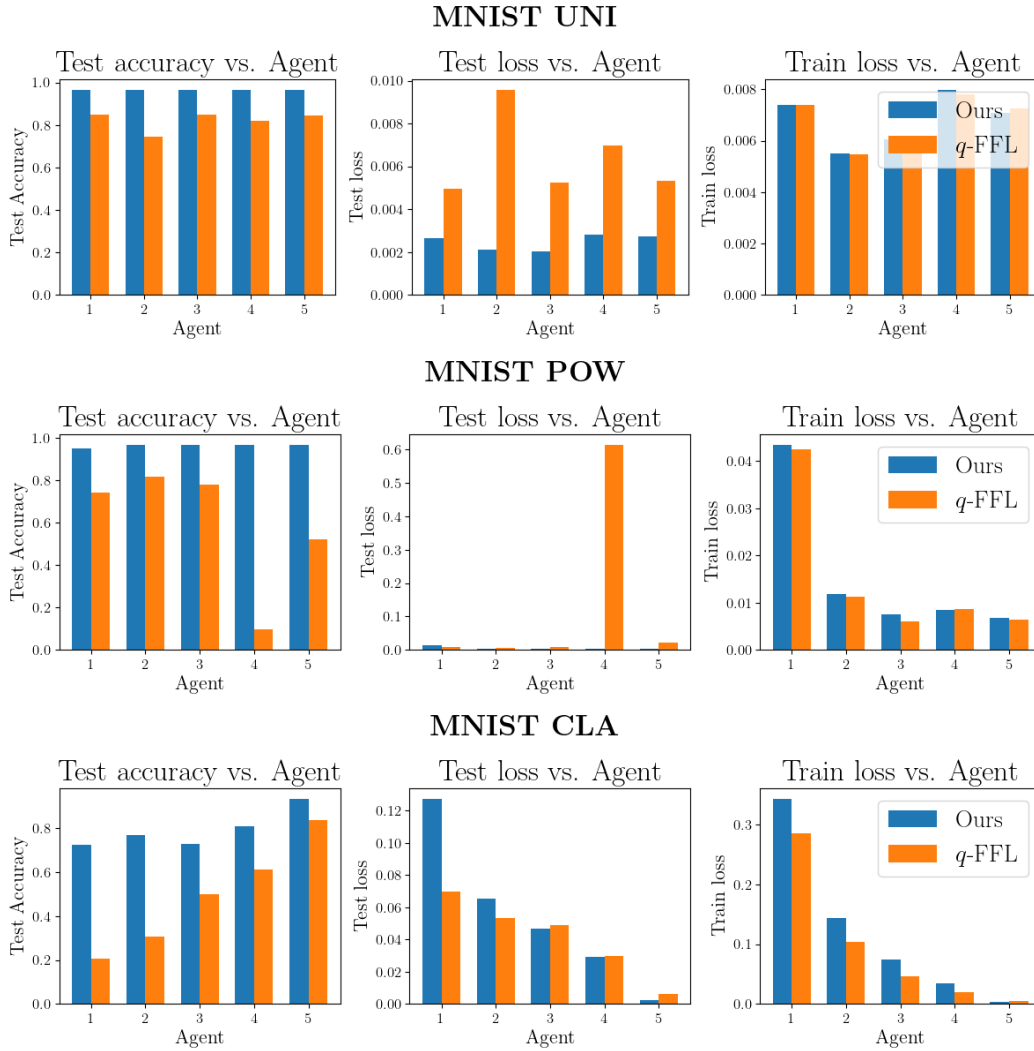


Figure 9: Comparison of final performance across agents between our method and q-FFL.

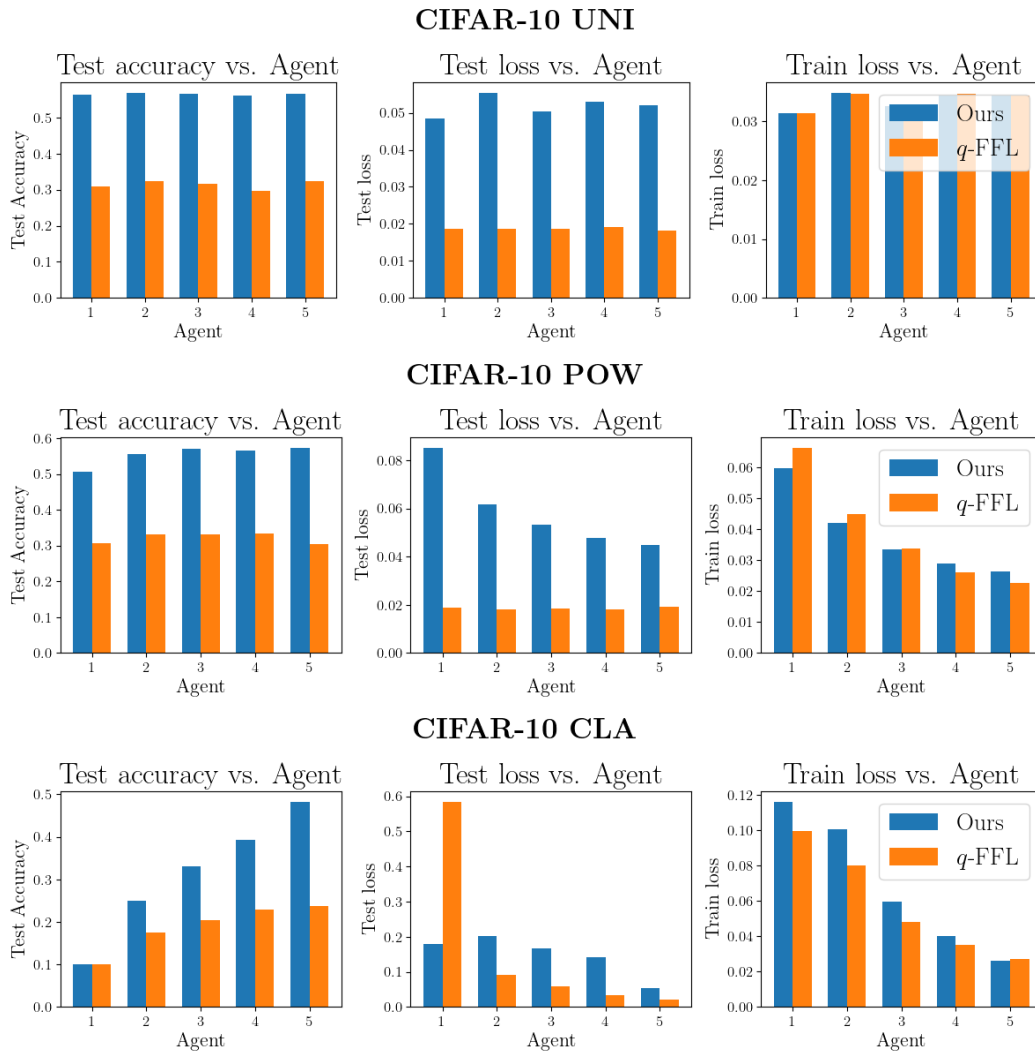


Figure 10: Comparison of final performance across agents between our method and q-FFL.

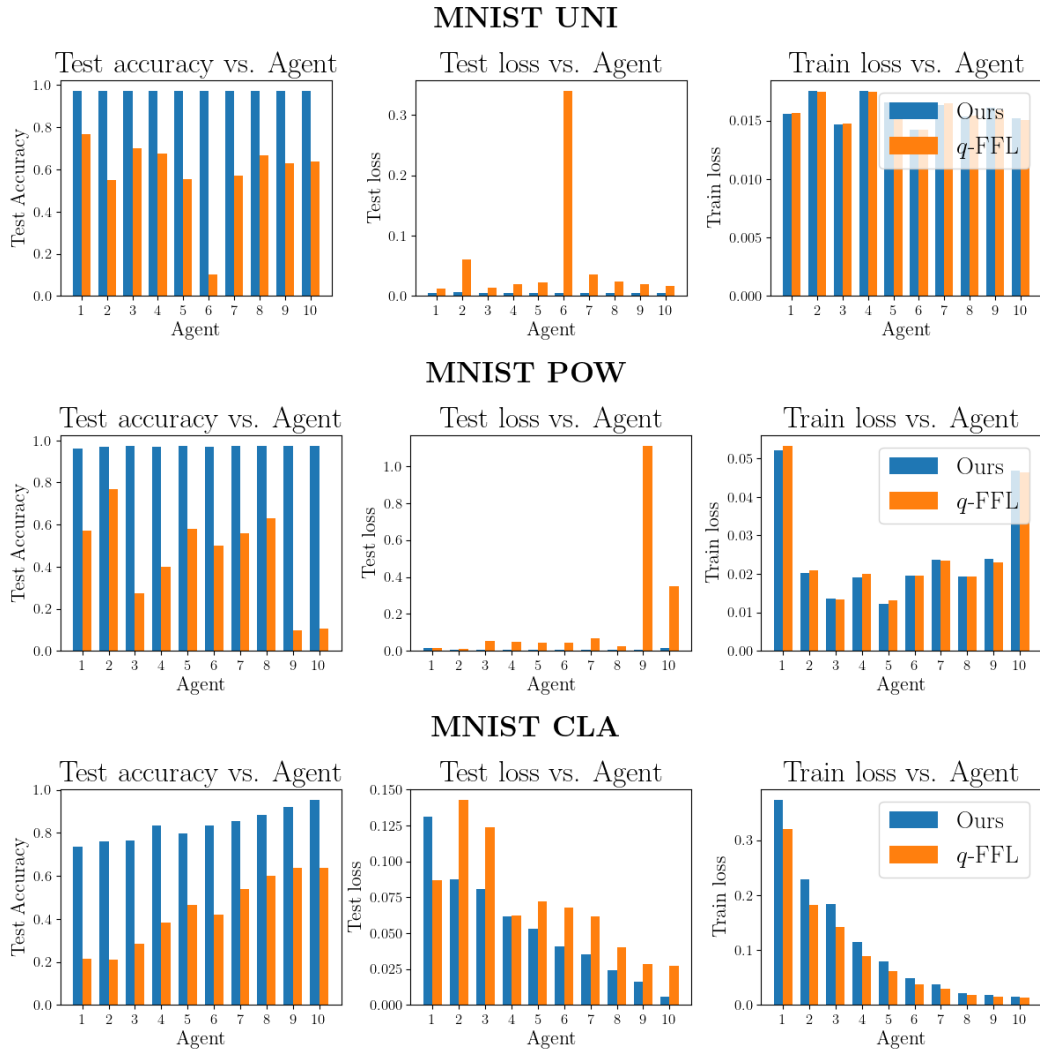


Figure 11: Comparison of final performance across agents between our method and q-FFL.



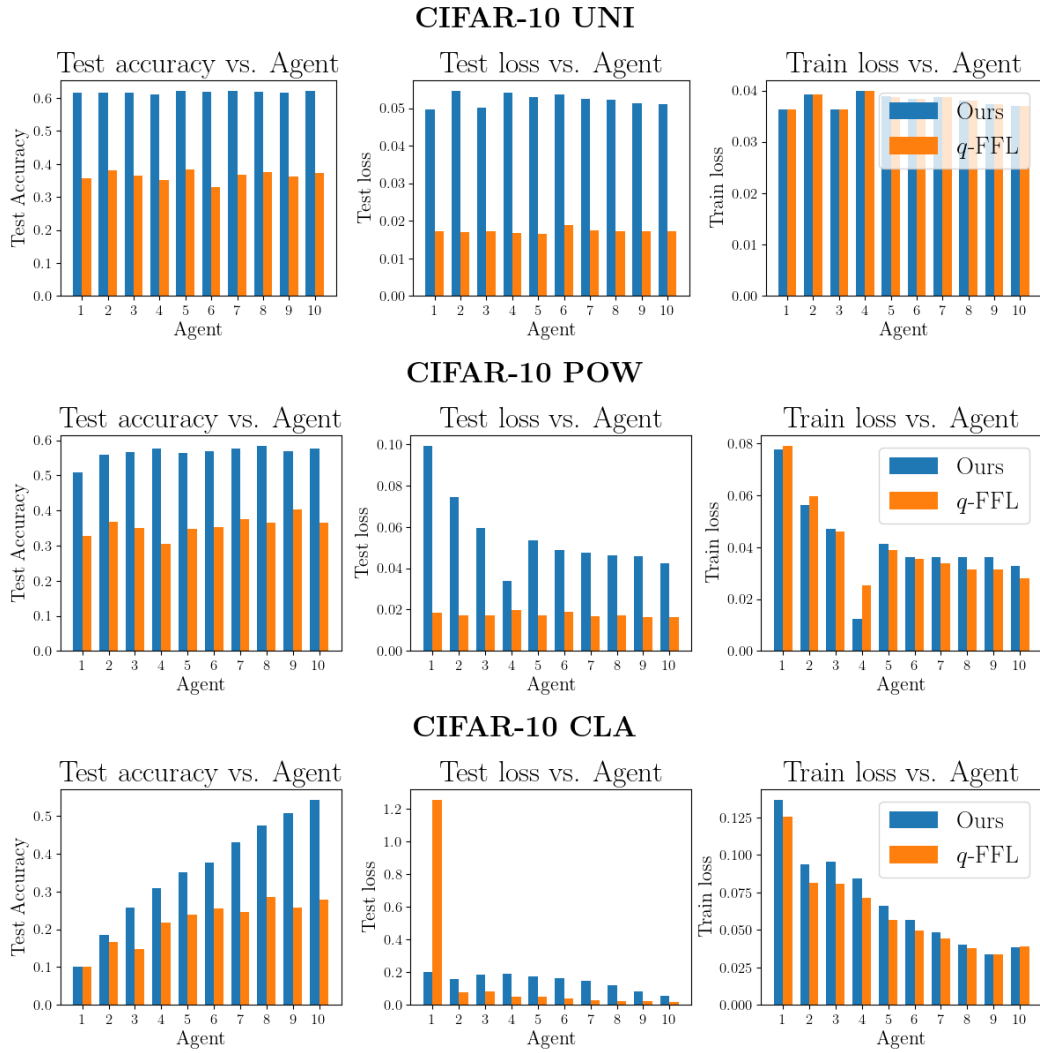


Figure 12: Comparison of final performance across agents between our method and q-FFL.