

---

# INCENTIVES IN FEDERATED LEARNING WITH HETEROGENEOUS AGENTS

**Ariel D. Procaccia**

Harvard University  
arielpro@seas.harvard.edu

**Han Shao\***

University of Maryland, College Park  
hanshao@umd.edu

**Itai Shapira**

Harvard University  
itaishapira@g.harvard.edu

## ABSTRACT

Federated learning promises significant sample-efficiency gains by pooling data across multiple agents, yet incentive misalignment is an obstacle: each update is costly to the contributor but boosts every participant. We introduce a game-theoretic framework that captures heterogeneous data: an agent’s utility depends on who supplies each sample, not just how many. Agents aim to meet a PAC-style accuracy threshold at minimal personal cost. We show that uncoordinated play yields pathologies: pure equilibria may not exist, and the best equilibrium can be arbitrarily more costly than cooperation. To steer collaboration, we analyze the cost-minimizing contribution vector, prove that computing it is NP-hard, and derive a polynomial-time linear program that achieves a logarithmic approximation. Finally, pairing the LP with a simple pay-what-you-contribute rule—each agent receives a payment equal to its sample cost—yields a mechanism that is strategyproof and, within the class of contribution-based transfers, is unique.

## 1 INTRODUCTION

Federated learning (FL) is a collaborative training framework in which multiple agents—each holding a distinct dataset—jointly optimize a global model while keeping data local. Collaboration allows agents to tap into information spread across heterogeneous records—for example, a network of hospitals pooling imaging data from distinct patient demographics to detect rare conditions sooner. Although each agent could independently train a model, collaboration offers higher accuracy or comparable performance with significantly fewer examples [1], enhancing both individual and federation-wide welfare.

However, realizing these gains hinges on incentives. Contributing model updates incurs costs in compute, bandwidth, curation effort, and privacy risk, while the global model produced by collective learning is a non-excludable public good: once trained, every agent benefits from its accuracy regardless of individual effort. This asymmetry invites a classic free-rider dilemma [2–4]: as one agent’s data lifts others’ accuracy while the contributor alone incurs the cost, each participant is tempted to trim its share once its own target is met. The resulting free-riding slows training and can ultimately exhaust the data pool that makes FL viable.

Consequently, federated learning can be viewed as a strategic game: each agent picks a contribution level to maximize a private utility that trades its own labeling cost against the benefit of the joint model. Existing models of incentives in federated learning either assume each agent’s utility depends on the local labeled data distributions of all agents [5], or treat agents as homogeneous, which crucially implies that data are *exchangeable* and each agent’s utility depends solely on the total amount of data contributed across all agents [4, 6].

---

\*Work mostly done while at Harvard University.

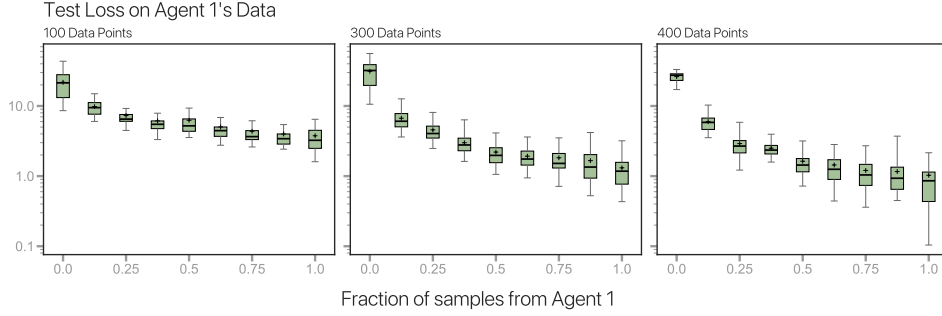


Figure 1: An agent’s expected loss falls as a larger share of a fixed training set comes from their *own* distribution. With a budget of  $m$  data points, we sample  $\lambda m$  from Agent 1 and  $(1 - \lambda)m$  from Agent 2 on FEMNIST [10], train a classifier, and repeat 100 times. For each  $m$ , Agent 1’s loss decreases monotonically in  $\lambda$ , confirming that data are *not* exchangeable—utilities depend on *who* contributes. Full details in Section A.1.

However, in typical federated learning scenarios, agents face *heterogeneous data distributions* and are primarily interested in improving performance on their own local distribution. This phenomenon—local data typically yielding greater marginal utility for local performance than data from other agents—has been well documented in the personalization and domain adaptation literature [7–9], and we empirically confirm that it is present in federated learning, as shown in Figure 1. We therefore model utilities that depend on *who* supplies each sample, not just how many.

In this work we ask: *how can incentives be aligned in this heterogeneous game?* We adopt a PAC accuracy objective—each agent wants their test error below  $\varepsilon$  with confidence  $1 - \delta$ —and study the game induced by this objective. Left on their own, agents settle into a contribution equilibrium that can be inefficient: some agents free-ride, others overspend, and total cost can explode relative to full cooperation. The central challenge is therefore to design contribution and transfer rules that coax self-interested, heterogeneously distributed agents into pooling enough data so that each meets its own accuracy target on its own distribution, while keeping the federation’s total cost near minimal.

**Our Results.** Section 2 builds a data-heterogeneous FL game with PAC thresholds: an agent’s utility depends on *who* supplies the samples, not just how many. Section 3 shows that decentralized play can fail badly—pure Nash equilibria may not exist and, when they do, their total cost can exceed the cooperative optimum by an unbounded factor (Price of Stability  $\rightarrow \infty$ ). Motivated by this gap, Section 4 assumes a planner with full information and full control. We prove that computing the cost-minimizing contribution vector is NP-hard (Theorem 2). Nonetheless, the PAC constraints admit a linear program (LP) whose cost is within a logarithmic factor of the optimum (Theorem 3); we use this LP allocation as the foundation of the mechanism design that follows. Section 5 drops the full-information assumption: agents now report private distributions. Pairing the LP allocation with a simple *pay-what-you-contribute* rule—each agent receives a payment equal to its sample cost—yields a mechanism that is strategyproof, and we derive conditions under which it uniquely satisfies these properties. Finally, while our main contributions are theoretical, we complement them with empirical validation and illustrative simulations in Section A that demonstrate the practical relevance of our assumptions and algorithms.

**Related Work.** Early work in FL centered on communication-efficient optimization and fairness under fully cooperative, i.i.d. or average-loss assumptions, with no pricing of data or modeling of strategic behavior. For instance, FedAvg optimizes average loss, while Agnostic FL and q-FFL reweight loss for worst-case or fairness gains, all assuming truthful reports [11–13]. Blum et al. [1] introduced a collaborative PAC model and showed that pooling across

$k$  heterogeneous tasks cuts sample complexity by  $\tilde{O}(\log k)$ , but participation is mandatory and incentives are ignored.

Subsequent work introduced incentives but frequently assumed that agents' data are exchangeable, so that model accuracy depends only on the total sample count contributed across all participants. Under this framework, various incentive strategies—such as transfer payments or reputational rewards—have been proposed to combat free-riding. For example, Karimireddy et al. [4] tie each agent's model quality to its data contribution (no external payments), whereas Murhekar et al. [6] design budget-balanced monetary transfers that implement welfare-maximizing equilibria; and [14–18] rely on reputation mechanisms and credit sharing. These designs treat data as interchangeable: model quality depends solely on the total sample count, so the marginal value of each sample ignores *who* provides it, overlooking realistic heterogeneity across data distributions.

A separate line of work studies heterogeneous data, aiming to capture how collaboration might form among agents with distinct interests. Donahue and Kleinberg [19], and Hasan [20] analyze coalition and Nash stability in model-sharing games and show that agents split into sub-coalitions when a global model biases some distributions. Blum et al. [5] further show that in a personalized PAC game with distribution-specific payoffs, pure Nash equilibria may not exist, and when they do, they can be arbitrarily inefficient, underscoring the fragility of cooperation absent incentives. These game-theoretic models emphasize that heterogeneity can severely complicate collaboration: individual incentives may fail to align with socially optimal pooling of data. Our work extends these papers by introducing a concrete utility model based on PAC-style threshold guarantees where each agent requires that the global model meet a distribution-specific accuracy threshold with high confidence.

Recent studies tackle truthfulness under heterogeneity: peer-prediction payments [21], budget-balanced truthful gradient schemes [22], and an optimal truthful mechanism for data sharing with interdependent valuations [23]; yet none compute minimum-cost allocations that meet each agent's welfare target.

## 2 MODEL

### 2.1 SETUP AND LEARNING PROTOCOL

We now make the collaborative game precise. Each agent selects how many examples to contribute; the federation pools these examples, trains a model, and each agent then evaluates the model on its own data distribution.

**Agents and data.** Consider  $k$  agents  $\mathcal{A} = \{1, \dots, k\}$  who wish to jointly learn a shared predictor but individually decide how much data to contribute. Let  $\mathcal{X}$  denote the instance space and  $\mathcal{Y}$  the label space. A hypothesis is a function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  mapping instances to predicted labels. Fix a hypothesis class  $\mathcal{H}$  with VC dimension  $d$ , and assume agents seek to approximate an unknown target function  $h^* \in \mathcal{H}$ . Each agent  $i \in \mathcal{A}$  has access to a local marginal data distribution  $\mathcal{D}_i$  over  $\mathcal{X}$  and can query labels for any data point drawn from this distribution. The collaborative learning process involves two stages:

**Stage 1 – Sample Collection.** Each agent  $i$  chooses a contribution level  $m_i \in \mathbb{N}$ , draws an i.i.d. *unlabeled* dataset  $U_i \sim \mathcal{D}_i^{m_i}$ , and queries the true labels of these samples, which are determined by the target function  $h^*$ . The labeled sets are pooled into a dataset

$$S = \bigcup_{i \in \mathcal{A}} \{ (x, h^*(x)) \mid x \in U_i \}.$$

Denote by  $\mathcal{P}(\mathcal{D}, \mathbf{m}, h^*)$  the distribution of this dataset, which is determined by the marginal data distributions  $\mathcal{D} = (\mathcal{D}_1, \dots, \mathcal{D}_k)$ , the contribution profile  $\mathbf{m} = (m_1, \dots, m_k)$ , and  $h^*$ .

**Stage 2 – Model Training.** A central server trains a model using Empirical Risk Minimization (ERM) method:

$$\text{ERM}(S) = \{h \in \mathcal{H} \mid \text{err}_S(h) = \min_{h' \in \mathcal{H}} \text{err}_S(h')\},$$

where  $\text{err}_S(h)$  is the empirical error of hypothesis  $h$  on dataset  $S$ . For any hypothesis  $h$ , marginal data distribution  $\mathcal{D}$ , and target function  $h^*$ , the generalization error is defined as

$$\text{err}_{\mathcal{D}, h^*}(h) = \mathbb{P}_{x \sim \mathcal{D}}(h(x) \neq h^*(x)) .$$

Given a labeled training set  $S$ , we define the generalization error of running ERM over  $S$  relative to  $(\mathcal{D}, h^*)$  as:

$$\text{err}_{\mathcal{D}, h^*}^{\text{ERM}}(S) = \max_{h \in \text{ERM}(S)} \text{err}_{\mathcal{D}, h^*}(h), \quad (1)$$

to reflect the generalization performance of an empirically optimal hypothesis under distribution  $\mathcal{D}$  with labeling function  $h^*$ . That is, if the ERM returns multiple minimizers, we take the worst generalization error, ensuring every bound holds under arbitrary tie-breaking.

Stages 1 and 2 define our collaborative game. Each agent  $i$  chooses an integer contribution level  $m_i$ . The server draws  $m_i$  i.i.d. samples from  $\mathcal{D}_i$  for each agent, pools the labeled data into  $S$ , and trains a single global ERM predictor  $\hat{h}(S)$  over  $\mathcal{H}$ . The next subsection specifies the payoffs from this global model and the cost of contributing data.

## 2.2 PAC ACCURACY OBJECTIVE

Each player ultimately wants a model whose test error on its *own* distribution is low; if collaboration fails, the fallback is to train alone. We focus on a single, widely-used performance criterion to formalize this goal: a Probably Approximately Correct (PAC) accuracy threshold.<sup>1</sup> In standard realizable PAC learning, an agent wants—with probability at least  $(1 - \delta)$ —to keep its generalization error below a tolerance  $\varepsilon$ .<sup>2</sup> Here we adapt that notion to our federated setting, where (1) the learner uses ERM on the pooled samples and (2) each agent  $i$  draws from its own fixed marginal distribution  $\mathcal{D}_i$ . We then define the PAC accuracy objective as follows: agent  $i$  requires that, *for every* target hypothesis  $h^* \in \mathcal{H}$ ,

$$\Pr_{S \sim \mathcal{P}(\mathcal{D}, \mathbf{m}, h^*)} \left[ \text{err}_{\mathcal{D}_i, h^*}^{\text{ERM}}(S) \leq \varepsilon \right] \geq 1 - \delta. \quad (2)$$

We use  $a_i^{\varepsilon, \delta}(\mathbf{m})$  to denote a *binary* variable that is 1 if Equation 2 is satisfied, and 0 otherwise; we refer to this as agent  $i$ 's  $(\varepsilon, \delta)$ -*requirement* or *target*.

**Cost of Contributions.** Contributing data incurs costs. Computation, and privacy risk all translate into a per-sample monetary burden. We capture this with a linear cost: when agent  $i$  supplies  $m_i$  samples, it pays  $c_i m_i$ , where  $c_i > 0$  denotes the cost per sample. The parameter  $c_i$  can reflect labeling fees, extra compute, or other participation overhead.

**Utility Functions.** Each agent balances the model's benefit against the data cost. We normalize monetary units so that achieving the PAC goal is worth exactly one unit of payoff, yielding a simplified expression for utility:

$$u_i(\mathbf{m}) = a_i^{\varepsilon, \delta}(\mathbf{m}) - c_i m_i. \quad (3)$$

To ensure every agent's goal is attainable and worth pursuing, we assume *self-sufficiency*: every agent could, in the worst case, meet its own accuracy requirement by training on its own data alone. Concretely, let  $n_i^{\text{ind}}$  denote the smallest number of samples that agent  $i$  would need to label by itself (with no contributions from others) to satisfy the  $(\varepsilon, \delta)$ -requirement on  $\mathcal{D}_i$ . Under PAC accuracy, for all marginal distributions, it holds that

$$n_i^{\text{ind}} \leq O\left(\frac{d + \ln(1/\delta)}{\varepsilon}\right),$$

where  $d = \text{VCdim}(\mathcal{H})$ , the VC dimension of  $\mathcal{H}$ .

<sup>1</sup>Other metrics, such as expected error, can be treated analogously; see Section 6.

<sup>2</sup>The heterogeneous (different  $\varepsilon_i$  for each agent) case is analogous, as we discuss in Section 6; we present the common- $\varepsilon$  case for clarity.

<sup>3</sup>For more discussion about the applicability of these payoffs to real-world settings, see the discussion in Section B.

---

**Assumption 1.**  $c_i n_i^{\text{ind}} < 1$ , for every  $i$ .

This self-sufficiency implies that failing to meet the accuracy threshold results in strictly lower utility than meeting it (since an agent can always fall back on solo training to obtain the benefit, albeit at potentially high cost).

**Social Cost Optimization.** A central planner seeks to maximize *utilitarian social welfare*, namely  $\sum_{i \in \mathcal{A}} u_i(\mathbf{m})$ . Under [Assumption 1](#), this welfare objective coincides with minimizing the aggregate labeling cost subject to the same PAC constraints. Thus, from a global standpoint, the problem is to find the contribution profile that lets every agent meet its accuracy target at the lowest possible total cost.

$$\begin{aligned} \min_{\mathbf{m}} \quad & \mathbf{c}^\top \mathbf{m} \\ \text{s.t.} \quad & a_i^{\varepsilon, \delta}(\mathbf{m}) = 1 \quad \forall i \in \mathcal{A}. \end{aligned} \tag{4}$$

Here  $\mathbf{c}^\top \mathbf{m}$  is the federation’s total labeling cost, and the constraint is the PAC guarantee of [Equation 2](#): for every agent  $i$  and every target hypothesis  $h^*$ , the ERM model must, with probability at least  $1 - \delta$ , achieve error at most  $\varepsilon$  on  $\mathcal{D}_i$ . The planner therefore seeks the cheapest contribution profile that satisfies *all* agents’  $(\varepsilon, \delta)$ -requirements simultaneously.

**Warm-Up Example.** Consider a single agent whose marginal distribution is uniform over the  $n$  distinct points  $x_1, \dots, x_n \subset X$  and a hypothesis class that contains every possible labeling of those points. When  $\varepsilon < 1/n$ , satisfying the PAC accuracy threshold below  $\varepsilon$  with probability at least  $1 - \delta$  requires drawing enough samples (with replacement) so that every point appears at least once. With  $k > 1$  agents, we can see how collaboration can improve sample complexity: consider  $k$  agents whose distribution places most of its probability mass on a different point  $x_i$ , while still assigning every other point probability at least  $\varepsilon$ . By pooling data, agents supply one another with examples they rarely see, allowing each to collect fewer local samples while the group still satisfies the PAC threshold.

### 3 THE INEFFICIENCY OF EQUILIBRIA

With the cooperative optimum in hand ([Equation 4](#)), we now analyze the uncoordinated game in which each agent strategically chooses its sample size. We quantify the cost gap between equilibrium behavior and the social optimum and show that a Nash equilibrium (NE) can be arbitrarily more costly. This gap motivates central coordination, which we develop next, and it serves as the conceptual starting point for the rest of the paper.

A contribution profile  $\mathbf{m}$  is a NE if no agent can raise its utility by unilaterally changing its sample count. Since each agent can always fall back to its self-sufficient plan  $n_i^{\text{ind}}$ , an equilibrium exists only when every accuracy threshold is met; otherwise the under-served agent would deviate to  $(n_i^{\text{ind}}, \mathbf{m}_{-i})$  to raise its utility. Thus, at equilibrium no agent wants to cut its contribution below the threshold or add extra samples.

It is straightforward to show that pure NE may not exist in our framework ([Theorem 6](#)), whereas mixed NE always exists ([Theorem 5](#)). Even when a pure equilibrium exists, it can be highly inefficient. To illustrate this, consider a simple two-player setting. Let the instance space be  $\mathcal{X} = \{x_A, x_B\}$ . Two players, Alice and Bob, each place most of their probability mass on a different point: for a small  $\varepsilon \in (0, \frac{1}{4})$ ,

$$\begin{aligned} \mathcal{D}_{\text{Alice}}(x_A) &= 1 - 2\varepsilon, & \mathcal{D}_{\text{Alice}}(x_B) &= 2\varepsilon, \\ \mathcal{D}_{\text{Bob}}(x_A) &= 2\varepsilon, & \mathcal{D}_{\text{Bob}}(x_B) &= 1 - 2\varepsilon. \end{aligned}$$

Under the PAC objective, each agent requires that *both* points be correctly classified with high probability. If Alice and Bob each contribute one labeled example, then with probability  $(1 - 2\varepsilon)^2$  Alice draws  $x_A$  while Bob draws  $x_B$ . Thus  $(m_{\text{Alice}}, m_{\text{Bob}}) = (1, 1)$  meets both thresholds at total cost  $c_{\text{Alice}} + c_{\text{Bob}} = \Theta(1)$  and is optimal. If one free-rides, the other must keep sampling until seeing *both* points, requiring  $\Omega(1/\varepsilon)$  samples. Since neither player can unilaterally lower this cost, this profile is an NE with a cost of  $\Omega(1/\varepsilon)$ .

This example shows that some equilibria can be much costlier than the social optimum even though a cost-minimal equilibrium exists. How large can the best equilibrium’s cost be? We formalize it with the Price of Stability (PoS) [24]: the ratio of total samples in the least-cost equilibrium to those in the planner’s optimum. PoS captures the cost of decentralization. Even the best self-enforcing outcome can require far more samples than the cooperative minimum. In the following, we show that PoS can be arbitrarily large:

**Theorem 1** (PoS is unbounded). *For any  $\varepsilon \in (0, 1)$  there exists a sequence of instances of our problem in which the ratio of the best NE to the optimal solution approaches  $\Omega(\log(1/\varepsilon))$ .*

## 4 CENTRAL COORDINATION WITH FULL INFORMATION

The unbounded price of stability in Section 3 shows that self-directed contribution games may squander resources. A natural fix is to appoint a *central planner*, such as a regulator or platform operator, who can dictate how many labeled examples each agent contributes. Before addressing strategic issues, we begin with the technical question of whether the planner can compute the cost-minimizing contribution vector efficiently. In this section we show that the problem is NP-hard (Theorem 2), yet it admits an efficient approximation via an LP with logarithmic-type guarantees (Theorem 3). The next section returns to strategy under coordination and builds on this to design payments that align incentives.

To isolate the computational question and remove strategic frictions for now, in this section we grant the planner two powers. First, it has *full information*—it observes every marginal distribution  $\mathcal{D}_i$  and each per-sample cost  $c_i$ . Second, it has *complete control*—if an agent opts in, the planner can compel it to supply the prescribed number of samples  $m_i$ . These assumptions fit cases where data statistics are public—e.g., hospital demographics or published mobile-traffic summaries—and participation is contractual. With strategic frictions removed, the task is now purely computational: who should collect how many samples? We tackle that question here and relax the full-information assumption in Section 5.

### 4.1 COLLABORATIVE PAC SAMPLE-ALLOCATION PROBLEM

Given a contribution vector  $\mathbf{m}$ , let  $S$  be the pooled sample obtained by drawing  $m_i$  points independently from each  $\mathcal{D}_i$ . We call  $\mathbf{m}$  *feasible* if, with probability at least  $1 - \delta$ , every ERM hypothesis trained on  $S$  has error at most  $\varepsilon$  on every  $\mathcal{D}_i$ . The planner specifies only the draw counts from each distribution. Under these full-information assumptions, the task reduces to Equation 4. Our first question is thus computational: is this optimization tractable? We answer in the negative. The result below shows NP-hardness in  $|\mathcal{H}|$  even with a single agent. The full proof is in Section D.1.

**Theorem 2.** *Under the PAC-accuracy objective, determining whether a specified sample count  $m$  suffices to meet the  $(\varepsilon, \delta)$ -requirement is NP-hard with respect to the hypothesis-class size  $|\mathcal{H}|$ , even when there is only one agent.*

### 4.2 APPROXIMATION VIA LINEAR PROGRAMMING

Despite this hardness barrier, the structure of the PAC constraints admits an efficient relaxation. Solving this LP takes polynomial time and returns a contribution vector whose total cost is within a logarithmic factor of the optimum. More specifically, for a finite class  $\mathcal{H}$ , our task is to find a vector of sample counts  $\mathbf{m} = (m_i)_{i \in \mathcal{A}}$  that—with probability at least  $1 - \delta$ —forces ERM to discard every hypothesis  $h$  that is *bad* for some agent:

$$\exists i \in \mathcal{A} : \mathcal{D}_i(\{x : h(x) \neq h^*(x)\}) \geq \varepsilon.$$

For any ordered pair  $(h_1, h_2)$  satisfying this condition, the probability that *no* sample lands in the set  $\{x : h_1(x) \neq h_2(x)\}$  is  $\prod_{i \in \mathcal{A}} (1 - \mathcal{D}_i(\{h_1 \neq h_2\}))^{m_i}$ , which is log-linear in  $\mathbf{m}$ . We can therefore convert Equation 4 into an LP with polynomially many constraints, which can be solved efficiently. As an empirical check, Section A.2 validates this LP allocation on a finite hypothesis class by comparing to the true optimum  $\sum_i m_i^*$  across varying  $|\mathcal{H}|$ .

For infinite  $\mathcal{H}$ , we provide an approximate solution by solving the LP over a finite cover  $\overline{\mathcal{H}} \subset \mathcal{H}$ , whose size is polynomial in  $\frac{1}{\varepsilon}$  and  $\frac{1}{\delta}$  when the VC dimension  $d = \text{VCdim}(\mathcal{H})$  is

bounded. While the resulting solution ensures PAC constraints are satisfied for  $\overline{\mathcal{H}}$ , it may not suffice for the full class  $\mathcal{H}$ , as there may exist a hypothesis  $h \in \mathcal{H} \setminus \overline{\mathcal{H}}$  that is consistent but still incurs high error. Nevertheless, we show that scaling the solution by a factor of roughly  $d$  suffices to ensure the PAC objective is met for all of  $\mathcal{H}$ .

**Theorem 3** (Approximation via Linear Programming). *Given any  $\mathcal{H}$  and  $\varepsilon, \delta > 0$ :*

- For finite  $\mathcal{H}$ , the LP over  $(\mathcal{H}, \varepsilon, \delta)$  returns a  $\frac{\log(1/\delta) + \log |\mathcal{H}|}{\log(1/\delta)}$ -approximate solution to Equation 4.
- For infinite  $\mathcal{H}$ , running the LP over  $(\overline{\mathcal{H}}, \varepsilon, \delta')$  and multiplying it by  $d + \log(1/\delta'')$  returns a  $O(\frac{d^2(\log(c_{\max}kd/(c_{\min}\varepsilon\delta))^2)}{\log(1/\delta)})$ -approximate solution to Equation 4, where  $\overline{\mathcal{H}} \subset \mathcal{H}$  is a  $\gamma$ -cover of  $\mathcal{H}$ ,  $\gamma = \Theta(c_{\min}\varepsilon\delta/(c_{\max}k(d + \log(1/\delta))))$ ,  $d = \text{VCdim}(\mathcal{H})$ ,  $\delta'' = \frac{\delta}{4|\mathcal{H}|}$  and  $\delta' = \frac{\delta}{8(d + \log(2|\mathcal{H}|/\delta))}$ ,  $c_{\min} = \min_{i \in [k]} c_i$  and  $c_{\max} = \max_{i \in [k]} c_i$ .

## 5 MECHANISM DESIGN WITH APPROXIMATE SOLUTIONS

In Section 4 we studied the planner’s computational problem under full information and control. Here we keep prescriptive control but drop full information and reintroduce strategy under coordination. In the uncoordinated game of Section 3 agents act by choosing their own sample sizes. With the planner fixing sample counts, agents can only influence outcomes by misreporting their local distributions, which shifts the LP constraints from Section 4,<sup>4</sup> and thus the computed contribution vector.<sup>5</sup> Building on that LP, in this section, we design a mechanism that is strategyproof and unique within a broad contribution-based class, blocking this manipulation and addressing the inefficiency in Section 3. Together, these results highlight the contrast between uncoordinated and coordinated settings.

Specifically, we now regard the agents’ local marginal data distributions as *private information* that must be reported to the mechanism. To alleviate the resulting incentive issues, we allow payments to the agents. The sequence of events, therefore, is as follows:

1. Agents report their local marginal data distributions,  $\mathcal{D}^r = (\mathcal{D}_1^r, \dots, \mathcal{D}_k^r)$ .
2. The central planner computes a solution  $\mathbf{m}$  based on the reported distributions  $\mathcal{D}^r$ .
3. Agents contribute according to  $\mathbf{m}$ .
4. The central planner pays each agent  $i$  an amount  $p_i$ .

We make the standard assumption that agent utilities are quasi-linear, that is, the utility of agent  $i$  for contribution vector  $\mathbf{m}$  and payment  $p_i$  (from the mechanism to the agent) is  $u_i(\mathbf{m}) + p_i$ . We say that a mechanism (which computes agent contributions and payments) is *strategyproof* if an agent can never benefit from misreporting their local marginal distribution; in game-theoretic terms, it is a dominant strategy to report  $\mathcal{D}_i^r = \mathcal{D}_i$ .

The challenge now is twofold: computation and incentives. In other words, the question is this: *How can we tractably compute contributions and payments such that the resulting mechanism is strategyproof?*

### 5.1 THE PAY-WHAT-YOU-CONTRIBUTE MECHANISM

In our model, the answer to the foregoing question is surprisingly immediate. To compute the contribution vector  $\mathbf{m}$ , we use the approximation algorithm of Theorem 3. For payments, we use the simple *pay-what-you-contribute* (PWYC) scheme, that is, compensate each agent

<sup>4</sup>The information that agents must provide to the server is finite and explicit. Our LP needs, for each pair  $h_1, h_2 \in \mathcal{H}$ , the probability that agent  $i$ ’s data lie in the disagreement set  $\{x : h_1(x) \neq h_2(x)\}$ . If  $\mathcal{H}$  is finite, this is a vector with at most  $|\mathcal{H}|^2$  entries. If  $\mathcal{H}$  is infinite, Section D.2 replaces  $\mathcal{H}$  with a finite  $\gamma$ -cover whose size is polynomial in  $d = \text{VCdim}(\mathcal{H})$ . Agents then report disagreement masses only for pairs in the cover, so the information burden remains finite and explicit.

<sup>5</sup>Enforcing prescribed contributions without accessing raw data raises a verification challenge. Practical approaches include cryptographic proofs and auditable contribution records [25, 26].

for its contribution, up to a constant:

$$p_i(\mathbf{m}) = c_i \cdot m_i + C_i, \quad (5)$$

for constants  $C_1, \dots, C_k$ .

Why is PWYC strategyproof? The utility of agent  $i$  when reporting truthfully is

$$u_i(\mathbf{m}) = a_i^{\varepsilon, \delta}(\mathbf{m}) - c_i \cdot m_i + c_i \cdot m_i + C_i = 1 + C_i,$$

as its learning threshold is satisfied by the contribution vector  $\mathbf{m}$  computed by the approximation algorithm. Note that this is the maximum possible utility under this mechanism: for any  $\mathbf{m}'$ , the utility of agent  $i$  is either  $1 + C_i$  or  $C_i$ , depending on whether its learning threshold is satisfied by  $\mathbf{m}'$ .

## 5.2 ALTERNATIVES TO PAY WHAT YOU CONTRIBUTE?

While the PWYC mechanism is strategyproof, one may wonder whether it is possible to design more sophisticated mechanisms with the goal of, for example, minimizing payments. A natural candidate is the classic Vickrey-Clarke-Groves (VCG) mechanism, which in our setting computes the optimal solution  $\mathbf{m}^{\text{OPT}}$ , asks agents to contribute according to  $\mathbf{m}^{\text{OPT}}$ , and then pays each agent  $i$

$$p_i = \sum_{j \neq i} u_j(\mathbf{m}^{\text{OPT}}) + q_i(\mathcal{D}_{-i}) = k - 1 - \sum_{j \neq i} c_j \cdot m_j^{\text{OPT}} + q_i(\mathcal{D}_{-i}^r), \quad (6)$$

where  $q_i(\mathcal{D}_{-i})$  is a term independent of the report  $\mathcal{D}_i$  of agent  $i$ . By standard arguments [27], the VCG mechanism ensures that each agent reports truthfully.

An obstacle to using VCG directly is that computing the optimal contribution vector is computationally hard (Theorem 2). In the *algorithmic mechanism design* [28] literature, however, there are various mechanisms that overcome computational hardness by augmenting approximation algorithms with clever payment schemes, including ones inspired by VCG [29, 30]. Is there such a rich mechanism design space in our setting? We give a partial negative answer to this question, showing that the PWYC mechanism is, in a qualified sense, unique.<sup>6</sup>

We start by defining a class of approximate algorithms and a class of “easy-to-compute” payment rules. We first introduce a *local obliviousness* property of contribution solutions, which helps distinguish approximate solutions from exact optima.

**Definition 1** (Locally Oblivious Approximations). *Given an approximation algorithm APPROX and a solution  $\mathbf{m}$ , we say that the algorithm is locally oblivious at  $\mathbf{m}$  if, for any neighbor  $\mathbf{m}'$  with  $\|\mathbf{m}' - \mathbf{m}\|_1 = 1$ , there exist distributions  $\mathcal{D} = (\mathcal{D}_1, \dots, \mathcal{D}_k)$  and  $\mathcal{D}' = (\mathcal{D}'_1, \dots, \mathcal{D}'_k)$  such that:*

- $\mathbf{m} = \text{APPROX}(\mathcal{D})$  and  $\mathbf{m}' = \text{APPROX}(\mathcal{D}'_i, \mathcal{D}_{-i})$  for all  $i \in [k]$ ,
- Both  $\mathbf{m}$  and  $\mathbf{m}'$  are feasible for  $\mathcal{D}$  and for each  $(\mathcal{D}'_i, \mathcal{D}_{-i})$ , for all  $i \in [k]$ .

Local obliviousness highlights a key difference between approximate solutions and exact optima. Intuitively, it reflects a certain slack in approximation: for any solution  $\mathbf{m}$  and any neighbor, we can construct pairs of distribution vectors where the approximation algorithm outputs  $\mathbf{m}$  and its neighbor respectively, and both solutions remain feasible across all these instances. For example, we show that the approximation algorithm introduced in Theorem 3 is locally oblivious at all  $\mathbf{m}$  satisfying  $m_1, m_2 \geq 2|\mathcal{H}| \log |\mathcal{H}|$  in the two-agent setting.

**Lemma 1.** *For any  $H$  greater than a universal constant  $C$ , there exists a PAC learning instance of  $(\mathcal{H}, \varepsilon, \delta)$  in the two-agent setting with  $|\mathcal{H}| = H$  such that the approximation algorithm introduced in Theorem 3 is oblivious at all  $\mathbf{m}$  with  $m_1, m_2 \geq 2|\mathcal{H}| \log |\mathcal{H}|$ .*

<sup>6</sup>This uniqueness is unfortunate, since a budget-balanced mechanism would be preferable. Even in standard mechanism design settings, particularly the celebrated VCG mechanism, exact budget balance is generally impossible, so the planner must make net payments to align incentives. Still, PWYC need not put all costs on the planner. By choosing constants  $C_i$  in Equation 5, the planner can recover part of the cost while preserving strategyproofness.

In general, a payment mechanism maps the reported distributions  $\mathcal{D}^r$  to a payment vector  $\mathbf{p}$ . However, it is often unclear how to effectively utilize the reported distributions directly. To address this, we consider a class of mechanisms that are easy to compute in practice—those that depend on the distributions only through the resulting contribution solution  $\mathbf{m}$ .

**Definition 2** (Contribution-Based Mechanisms). *Given an algorithm  $APPROX$ , a payment mechanism  $\mathbf{p}$  is contribution-based if there exists  $\mathbf{f} : \mathbb{N}^k \rightarrow \mathbb{R}^k$  and  $\mathbf{q}$  such that*

$$p_i(\mathcal{D}^r) = f_i(APPROX(\mathcal{D}^r)) + q_i(\mathcal{D}_{-i}^r),$$

*Thus agent  $i$ 's own report affects  $p_i$  only via the computed contribution vector.*

We now show that PWYC is the only strategyproof contribution-based payment mechanism when approximate solutions are used for two agents (see proof in [Section E](#)):

**Theorem 4.** *For any strategyproof contribution-based payment mechanism  $\mathbf{f}$ , approximation algorithm  $APPROX$ , and connected  $M \subset \mathbb{N}^k$ , if  $APPROX$  is oblivious at all  $\mathbf{m} \in M$ , there exist constants  $C_1, C_2$  such that for all  $\mathbf{m} \in M$ :*

$$f_i(\mathbf{m}) = c_i \cdot m_i + C_i, \forall i = 1, 2,$$

**Corollary 1.** *For any  $H$  greater than a universal constant  $C$ , there exists a PAC learning instance of  $(\mathcal{H}, \varepsilon, \delta)$  in the two-agent setting with  $|\mathcal{H}| = H$  such that when applying the approximation algorithm introduced in [Theorem 3](#) to compute the contribution solution, the strategyproof contribution-based payment  $\mathbf{f}$  must satisfy*

$$f_i(\mathbf{m}) = c_i \cdot m_i + C_i, \forall i = 1, 2,$$

*for all  $\mathbf{m}$  with  $m_1, m_2 \geq 2|\mathcal{H}| \log |\mathcal{H}|$  for some constants  $C_1, C_2$ .*

[Corollary 1](#) follows by combining [Theorem 4](#) with [Lemma 1](#). [Section F](#) extends the same approximation to a broader class of objectives.

## 6 DISCUSSION

We conclude by discussing the implications of our results, examining key modeling assumptions, and outlining connections to broader theory.

**Beyond the PAC accuracy objective.** The  $(\varepsilon, \delta)$ -guarantee controls the *tail*: with probability at least  $1 - \delta$  the generalization error does not exceed  $\varepsilon$ . A complementary member of the same threshold family bounds the *expected* error by  $\varepsilon$ , trading worst-case assurance for an average-case criterion that may better reflect practical risk tolerance. Most of our analysis carries over unchanged: [Section F.1](#) develops approximation schemes for computing optimal contribution vectors under expected error, and [Section F.2](#) proves that these solutions continue to satisfy the local obliviousness property.

Taking a broader view of utilities, we focus on *threshold payoffs*. With quasi-linear utility, maximizing utility equals minimizing cost subject to the threshold. This replaces the intractable dependence of each agent's loss on its contribution with a tractable optimization. Future work can explore richer models that trade off cost and payoff more flexibly.

**Common accuracy threshold.** For clarity, we assume a common accuracy target  $\varepsilon$  across agents. This simplifying assumption streamlines notation and highlights the structural properties of the contribution game. However, our results do not rely on uniformity. Each agent's objective can be parameterized by a distinct  $\varepsilon_i$ , and all definitions, equilibrium analyses, and approximation guarantees extend naturally to this heterogeneous setting.

**Linear costs.** Linear per-sample costs simplify analysis and are common in the literature, but they can be restrictive and not strictly necessary. Permitting any convex, non-decreasing cost function captures realistic scenarios with increasing marginal cost—for instance, later data points may be more expensive to collect. The strategic conclusions are unchanged: a NE still exists and pure equilibria need not. Meanwhile, the *linear* program in our approximation algorithm becomes a convex program that remains tractable. Thus, the incentive mechanism we derive continues to work even when agents' data-generation costs rise with effort, at the price of solving slightly more involved convex optimizations.

---

## REFERENCES

- [1] Avrim Blum, Nika Haghtalab, Ariel D Procaccia, and Mingda Qiao. Collaborative pac learning. *Advances in Neural Information Processing Systems*, 30, 2017.
- [2] Dejun Yang, Guoliang Xue, Xi Fang, and Jian Tang. Incentive mechanisms for crowd-sensing: Crowdsourcing with smartphones. *IEEE/ACM transactions on networking*, 24(3):1732–1744, 2015.
- [3] Abrar Ahmed and Bong Jun Choi. Frimfl: A fair and reliable incentive mechanism in federated learning. *Electronics*, 12(15):3259, 2023.
- [4] Sai Praneeth Karimireddy, Wenshuo Guo, and Michael I Jordan. Mechanisms that incentivize data sharing in federated learning. *arXiv preprint arXiv:2207.04557*, 2022.
- [5] Avrim Blum, Nika Haghtalab, Richard Lanus Phillips, and Han Shao. One for one, or all for all: Equilibria and optimality of collaboration in federated learning. In *International Conference on Machine Learning*, pages 1005–1014. PMLR, 2021.
- [6] Aniket Murhekar, Zhuowen Yuan, Bhaskar Ray Chaudhury, Bo Li, and Ruta Mehta. Incentives in federated learning: Equilibria, dynamics, and mechanisms for welfare maximization. *Advances in Neural Information Processing Systems*, 36, 2024.
- [7] Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations for domain adaptation. *Advances in neural information processing systems*, 19, 2006.
- [8] Ayan Kumar Bhunia, Shuvojit Ghose, Amandeep Kumar, Pinaki Nath Chowdhury, Aneeshan Sain, and Yi-Zhe Song. Metahttr: Towards writer-adaptive handwritten text recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15830–15839, 2021.
- [9] Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.
- [10] Sebastian Caldas, Peter Wu, Tian Li, Jakub Konečný, H. Brendan McMahan, Virginia Smith, and Ameet Talwalkar. LEAF: A benchmark for federated settings. *CoRR*, abs/1812.01097, 2018.
- [11] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- [12] Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic Federated Learning. In *Proceedings of the 36th International Conference on Machine Learning*, pages 4615–4625. PMLR, May 2019. ISSN: 2640-3498.
- [13] Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. Fair Resource Allocation in Federated Learning, February 2020. arXiv:1905.10497 [cs].
- [14] Jierui Lin, Min Du, and Jian Liu. Free-riders in Federated Learning: Attacks and Defenses, November 2019. arXiv:1911.12560 [cs].
- [15] Jiawen Kang, Zehui Xiong, Dusit Niyato, Shengli Xie, and Junshan Zhang. Incentive Mechanism for Reliable Federated Learning: A Joint Optimization Approach to Combining Reputation and Contract Theory. *IEEE Internet of Things Journal*, 6(6): 10700–10714, December 2019.
- [16] Yunus Sarikaya and Ozgur Ercetin. Motivating Workers in Federated Learning: A Stackelberg Game Perspective, August 2019. arXiv:1908.03092 [cs].

- 
- [17] Ningning Ding, Zhixuan Fang, and Jianwei Huang. Incentive Mechanism Design for Federated Learning with Multi-Dimensional Private Information. In *2020 18th International Symposium on Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks (WiOPT)*, pages 1–8, June 2020.
  - [18] Yann Fraboni, Richard Vidal, and Marco Lorenzi. Free-rider Attacks on Model Aggregation in Federated Learning, February 2021. arXiv:2006.11901 [cs].
  - [19] Kate Donahue and Jon Kleinberg. Optimality and Stability in Federated Learning: A Game-theoretic Approach, June 2021. arXiv:2106.09580 [cs].
  - [20] Cengiz Hasan. Incentive Mechanism Design for Federated Learning: Hedonic Game Approach, May 2021. arXiv:2101.09673 [cs].
  - [21] Jinlong Pang, Jiaheng Wei, Chen Qian, and Yang Liu. Incentivizing data collection from heterogeneous clients in federated learning. *arxiv*, 2022.
  - [22] Dimitar Chakarov, Nikita Tsoy, Kristian Minchev, and Nikola Konstantinov. Incentivizing truthful collaboration in heterogeneous federated learning. *arXiv preprint arXiv:2412.00980*, 2024.
  - [23] Mengjing Chen, Yang Liu, Weiran Shen, Yiheng Shen, Pingzhong Tang, and Qiang Yang. A mechanism design approach for multi-party machine learning. In *International Workshop on Frontiers in Algorithmics*, pages 248–268. Springer, 2022.
  - [24] Elliot Anshelevich, Anirban Dasgupta, Jon Kleinberg, Éva Tardos, Tom Wexler, and Tim Roughgarden. The price of stability for network design with fair cost allocation. *SIAM Journal on Computing*, 38(4):1602–1623, 2008.
  - [25] Juan Ma, Hao Liu, Mingyue Zhang, and Zhiming Liu. Vpfl: Enabling verifiability and privacy in federated learning with zero-knowledge proofs. *Knowledge-Based Systems*, 299:112115, 2024.
  - [26] K Naveen Kumar, Ranjeet Ranjan Jha, C Krishna Mohan, and Ravindra Babu Tallamraju. Fortifying federated learning towards trustworthiness via auditable data valuation and verifiable client contribution. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 4999–5009, 2025.
  - [27] Noam Nisan. Introduction to mechanism design (for computer scientists). In Noam Nisan, Tim Roughgarden, Éva Tardos, and Vijay Vazirani, editors, *Algorithmic Game Theory*, chapter 9. Cambridge University Press, 2007.
  - [28] Noam Nisan and Amir Ronen. Algorithmic mechanism design. *Games and Economic Behavior*, 35(1–2):166–196, 2001.
  - [29] Daniel Lehmann, Liadan Ita O’callaghan, and Yoav Shoham. Truth revelation in approximately efficient combinatorial auctions. *Journal of the ACM (JACM)*, 49(5):577–602, 2002.
  - [30] Shahar Dobzinski. Better mechanisms for combinatorial auctions via maximal-in-range algorithms? *ACM SIGecom Exchanges*, 7(1):30–33, 2007.
  - [31] Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
  - [32] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik. Emnist: an extension of mnist to handwritten. *Proceedings of the IEEE*, 4322, 2017.
  - [33] William Shakespeare et al. *William Shakespeare: the complete works*. Barnes & Noble Publishing, 1989.

- 
- [34] Nikhil Vyas, Depen Morwani, Rosie Zhao, Mujin Kwun, Itai Shapira, David Brandfonbrener, Lucas Janson, and Sham Kakade. Soap: Improving and stabilizing shampoo using adam. *arXiv preprint arXiv:2409.11321*, 2024.
- [35] Michael R Garey and David S Johnson. *Computers and intractability*, volume 29. wh freeman New York, 2002.
- [36] David Haussler. Sphere packing numbers for subsets of the boolean n-cube with bounded vapnik-chervonenkis dimension. *Journal of Combinatorial Theory, Series A*, 69(2):217–232, 1995.

## A EMPIRICAL VALIDATION AND ILLUSTRATIVE SIMULATIONS

**Overview.** Our main contributions are theoretical. This appendix complements them with two empirical studies that illustrate the modeling assumptions and the practical significance of the planner and mechanism. First, we show that data are not exchangeable, since an agent benefits more from data drawn from its own distribution. Second, we validate the planner’s LP allocation on a finite hypothesis class by comparing its total cost to the true optimum found by search.

**Dataset Selection and Motivation.** We build on the LEAF benchmark [10], which provides realistic federated datasets naturally partitioned by user. We selected two datasets, **FEMNIST** and **Shakespeare**, to cover both vision and textual tasks. These specific datasets were chosen because they have the highest average number of datapoints per user among the LEAF collection [10]. FEMNIST, built on Extended MNIST [31, 32], consists of handwritten character images partitioned by writer. Shakespeare is constructed from *The Complete Works of William Shakespeare* [33, 11], partitioned by speaking role, with each role in each play treated as a distinct agent.

**Dataset Construction.** We construct two-agent subsets from each dataset by selecting the two users with the most data points. For FEMNIST, these are writers f0261\_06 and f0289\_10, jointly contributing 1,047 samples. For Shakespeare, we select the roles Hamlet (from *Hamlet*) and Iago (from *Othello*), together providing 112,116 samples. Additionally, we create a variant scenario with one top agent and a random convex combination of five other agents forming the second agent. We illustrate qualitative differences between the FEMNIST agents in Figure 2.

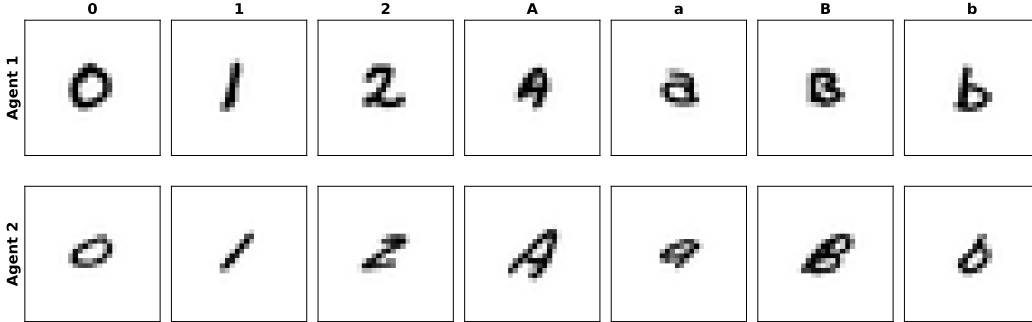


Figure 2: Samples from selected character classes for two FEMNIST agents, illustrating distinct handwriting styles.

### A.1 EFFECT OF DATA COMPOSITION ON AGENT-SPECIFIC PERFORMANCE

**Experimental Setup.** For varying total data size  $m$  and fraction  $\lambda \in [0, 1]$  from Agent 1, we construct datasets  $D^{(m, \lambda)}$  comprising  $\lambda m$  samples from Agent 1 and  $(1 - \lambda)m$  samples from Agent 2. Models trained on these datasets are evaluated separately on each agent’s test set. We use standard cross-entropy loss and the SOAP optimizer [34]. To capture performance variability, experiments are repeated 100 times with different random seeds. See Table 1 for detailed parameters. We employed task-appropriate model architectures aligned with prior federated learning works [11, 10, 6]. For FEMNIST (image classification), the model is a lightweight CNN. For Shakespeare (next-character prediction), we use a recurrent neural network with an LSTM backbone.

**Results and Discussion.** These experiments confirm the intuitive expectations about data volume and personalization, and quantify their effects. (i) For any fixed share  $\lambda$ , increasing the budget  $m$  lowers both agents’ losses—more data helps everyone (Assumption 2). (ii) For any fixed  $m$ , raising  $\lambda$  cuts Agent 1’s loss: when the training set includes a higher

Table 1: Summary of experimental parameters.

Dataset	Agent 1	Agent 2	# Seeds	$m$ values	Learning rate	Epochs
FEMNIST	f0261_06	f0289_10	100	10 – 450	0.01	50
Shakespeare	Hamlet	Iago	100	10,000 – 45,000	0.01	25

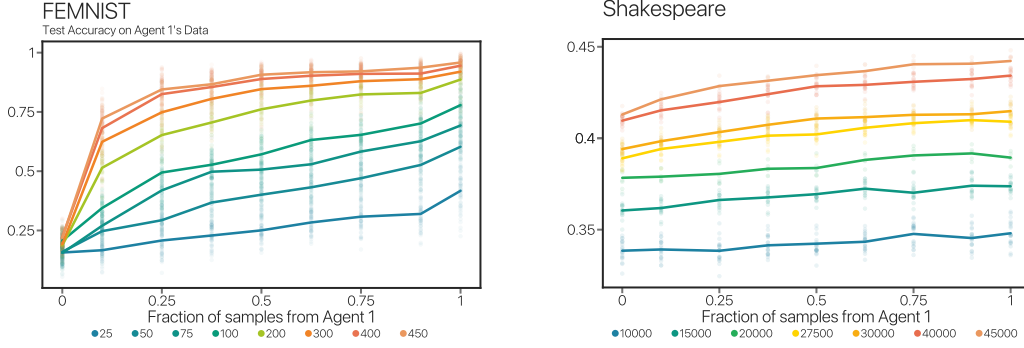


Figure 3: Agent 1’s test accuracy as a function of the fraction of samples contributed from its own distribution (Agent 1), compared against random convex combinations from other agents (Agent 2). Each curve represents different total dataset sizes ( $m$ ). The monotonic trend supports the weak-monotonicity [Assumption 2](#).

proportion of one writer’s style, the model is better tuned to that style and thus generalizes better on that agent’s test set. These trends mirror prior personalization findings in federated learning and domain adaptation [\[9\]](#).

#### A.2 PLANNER VALIDATION ON A FINITE HYPOTHESIS CLASS

We empirically validate the planner’s LP allocation from [Section 4.2](#) by comparing its total prescribed samples to the true optimum obtained by discrete search over integer allocations. We report the ratio

$$r = \frac{\sum_i m_i^{\text{LP}}}{\sum_i m_i^*}$$

under uniform per-sample costs  $c_i = 1$ .

**Finite-class construction.** We focus on FEMNIST and construct a finite hypothesis class  $\mathcal{H}$  by sampling checkpoints from a compact CNN with two convolutional blocks and a linear classifier. For each agent, we estimate pairwise disagreements  $\mathcal{D}_i(\{h_1 \neq h_2\})$  for all  $h_1, h_2 \in \mathcal{H}$  using unlabeled data, i.e., the fraction of examples on which the two checkpoints predict different labels. These estimates instantiate the constraints in [Section 4.2](#), yielding the LP-based contribution vector  $\mathbf{m}^{\text{LP}}$ .

**Experimental setup.** We assume  $c_i = 1$  for all agents. Given  $\mathbf{m}^{\text{LP}}$ , we perform an independent discrete search over integer allocations and, by repeated evaluation of the checkpoints in  $\mathcal{H}$ , find the smallest  $\mathbf{m}^*$  that satisfies every agent’s  $(\varepsilon, \delta)$  target. We sweep the size  $|\mathcal{H}|$  of the sampled class and record  $r$ .

**Results.** The observed values of  $r$  across  $|\mathcal{H}|$  are reported in [Table 2](#). The LP allocation achieves logarithmic-type approximation factors consistent with [Theorem 3](#).

## B ON BINARY PAYOFFS

We model each agent’s payoff as 1 if its PAC requirement holds and 0 otherwise. With this choice, *each agent* solves a cost-minimization problem subject to its threshold:

$$\min_{m_i \in \mathbb{N}} m_i \quad \text{s.t.} \quad a_i^{\varepsilon, \delta}(\mathbf{m}) = 1.$$

Table 2: LP total versus the optimal total on FEMNIST for a finite hypothesis class  $\mathcal{H}$ .

$ \mathcal{H} $	$r = \frac{\sum_i m_i^{\text{LP}}}{\sum_i m_i^*}$
5	1.67
10	3.19
15	3.94
20	3.57
25	4.24
27	3.99
50	4.20

We believe this individual cost-minimization task reflects real-world scenarios,<sup>7</sup> and by standard Lagrange multiplier theory, there exists a multiplier  $\lambda > 0$  such that any optimal solution  $m_i^*$  to the constrained problem is also optimal for a penalized objective of the form

$$\max_{m_i} a_i(\mathbf{m}) - \lambda m_i$$

Hence the binary case, which made the exposition lighter, is not a fundamental limitation but a modeling convenience that extends to a broader family of monotone utilities that trade off accuracy against cost with weight  $\lambda$ .

## C EXISTENCE AND EFFICIENCY OF EQUILIBRIA

In this section, we study the strategic behavior of self-interested agents in our model.

### C.1 EXISTENCE OF EQUILIBRIA

We impose weak *monotonicity*—more data from any agent never hurts another. The assumption guarantees equilibrium existence.

**Assumption 2** (Monotonicity). *For every agent  $i$ , utility is weakly increasing in any agent’s contribution. Let  $\mathbf{m}$  and  $\mathbf{m}'$  be two contribution profiles with  $m'_j \geq m_j$  for every agent  $j$ . Then*

$$u_i(\mathbf{m}') \geq u_i(\mathbf{m}) \quad \forall i \in \mathcal{A}.$$

**Theorem 5** (Existence of Nash Equilibrium). *A Nash equilibrium exists.*

*Proof.* Since every agent  $i \in \mathcal{A}$  is self-sufficient,  $i$  can, by contributing  $n_i^{\text{ind}} < \infty$  samples on their own, satisfy their objective. Since additional samples supplied by other agents never reduce any agent’s probability of meeting its objective, the strategy  $m_i = n_i^{\text{ind}}$  guarantees utility  $1 - c_i n_i^{\text{ind}} \geq 0$  for every profile  $\mathbf{m}_{-i}$ . In contrast, any action with  $c_i m_i > 1$  yields utility  $1 - c_i m_i < 0$  irrespective of the others’ choices, so every such action is strictly dominated by  $n_i^{\text{ind}}$ . Thus agent  $i$ ’s undominated actions lie in the finite set  $\{0, 1, \dots, n_i^{\text{ind}}\}$ . With each player restricted to finitely many pure strategies, the game is finite, and Nash’s existence theorem for finite games guarantees at least one (possibly mixed) Nash equilibrium.  $\square$

We demonstrate, however, that pure Nash equilibria do not always exist:

**Theorem 6** (Non-existence of Pure Nash Equilibrium). *There exists a PAC learning setting in which no pure Nash equilibrium exists.*

*Proof.* Consider an instance space  $\mathcal{X} = \{x_1, x_2, x_3\}$  and a hypothesis class  $\mathcal{H}$  that contains all possible labeling. Let agents have data distributions: Define the data distributions of the three agents by the point-probabilities

$$\begin{aligned} \mathcal{D}_1(x_1) &= \frac{1}{3}, \mathcal{D}_1(x_2) = \frac{2}{3}, \mathcal{D}_1(x_3) = 0, \\ \mathcal{D}_2(x_1) &= 0, \mathcal{D}_2(x_2) = \frac{1}{3}, \mathcal{D}_2(x_3) = \frac{2}{3}, \\ \mathcal{D}_3(x_1) &= \frac{2}{3}, \mathcal{D}_3(x_2) = 0, \mathcal{D}_3(x_3) = \frac{1}{3}. \end{aligned}$$

<sup>7</sup>For example, hospitals may aim to minimize their data contribution while ensuring that the global model achieves at least 99% accuracy on their local patient population, being indifferent between 99% and 99.5% once the target is met.

Fix the PAC parameters  $\varepsilon = \frac{1}{3}$  and  $\delta = \frac{2}{3}$ . Each agent seeks accuracy at least  $\frac{2}{3}$  with probability at least  $\frac{1}{3}$ . Hence every agent  $i$  requires that, with probability at least  $1 - \delta = \frac{1}{3}$ , the learned classifier incurs error at most  $\varepsilon$  on  $\mathcal{D}_i$ .

In this setup, no agent is incentivized to contribute more than one sample. Checking possible pure strategies, one observes: At  $\mathbf{m} = (1, 1, 0)$ , agent 1 can deviate to 0 without losing accuracy, incentivizing deviation. And, at  $\mathbf{m} = (0, 1, 0)$ , the contribution level is insufficient for agent 3. Thus, no pure equilibrium exists.  $\square$

## C.2 EFFICIENCY OF EQUILIBRIA: PRICE OF STABILITY

**Theorem 1** (PoS is unbounded). *For any  $\varepsilon \in (0, 1)$  there exists a sequence of instances of our problem in which the ratio of the best NE to the optimal solution approaches  $\Omega(\log(1/\varepsilon))$ .*

*Proof.* Fix  $\varepsilon \in (0, \frac{1}{2})$  and  $\delta \in (0, 1)$ . We construct an instance whose PoS satisfies

$$\text{PoS} = \frac{\log(1/\varepsilon) + \log(1/\delta)}{\log(1/\delta)} = \Omega(\log(1/\varepsilon)),$$

so the ratio diverges as  $\varepsilon \rightarrow 0$ . Let  $n > \frac{1}{2\varepsilon}$  and set

$$\mathcal{X} = \{x_1, \dots, x_n, y, z\}.$$

Define the hypothesis class

$$\mathcal{H} = \{h_0\} \cup \{h_i : i \in [n]\},$$

where

$$h_0(x) = 0 \quad (\forall x \in \mathcal{X}), \quad h_i(x) = \begin{cases} 1 & \text{if } x \in \{x_i, z\}, \\ 0 & \text{otherwise.} \end{cases}$$

Thus  $y$  is always labeled 0, while the label of  $z$  determines whether all  $x_1, \dots, x_n$  are 0 ( $h_0$ ) or exactly one of them is 1 ( $h^i$ ).

Alice and Bob have marginal distributions

$$\begin{aligned} \mathcal{D}_{\text{Alice}}(x_i) &= 2\varepsilon \quad \forall i \in [n], & \mathcal{D}_{\text{Alice}}(y) &= \mathcal{D}_{\text{Alice}}(z) = 0, \\ \mathcal{D}_{\text{Bob}}(z) &= \varepsilon, \quad \mathcal{D}_{\text{Bob}}(y) &= 1 - \varepsilon, & \mathcal{D}_{\text{Bob}}(x_i) &= 0 \quad \forall i \in [n]. \end{aligned}$$

In every NE, Bob contributes 0 samples, because any  $h \in \mathcal{H}$  has error at most  $\varepsilon$  on his distribution, so his  $(\varepsilon, \delta)$ -requirement is already satisfied. In the worst case, where  $h_0$  is the target function, Alice must choose  $m$  so that she samples every point  $x_1, \dots, x_n$  with probability at least  $1 - \delta$ . Let  $m_{\text{NE}}$  be the smallest integer for which this condition holds.

In the minimum-cost contribution vector, Bob supplies  $m_{\text{Bob}}^{\text{opt}} = \frac{\log(1/\delta)}{\log(1-\varepsilon)}$  samples, which include  $z$  with probability  $1 - \delta$ . If the true target is  $h_0$  no further data are needed. Otherwise the target is some  $h_i$ ; Alice then has to sample only until she sees the single point  $x_i$ . Let  $m_{\text{Alice}}^{\text{OPT}}$  be the least integer guaranteeing this with probability  $1 - \delta$ . Sampling until one designated point appears is far cheaper than sampling until all  $n$  points appear. Hence the optimal cost is

$$\|\mathbf{m}^{\text{opt}}\|_1 = m_{\text{Bob}}^{\text{opt}} + m_{\text{Alice}}^{\text{opt}} \leq \frac{3 \log(1/\delta)}{2\varepsilon}.$$

Taking the ratio,

$$\text{PoS} = \frac{\|\mathbf{m}^{\text{NE}}\|_1}{\|\mathbf{m}^{\text{OPT}}\|_1} = \frac{\log(1/\varepsilon) + \log(1/\delta)}{\log(1/\delta)} = \Omega(\log(1/\varepsilon)),$$

which grows without bound as  $\varepsilon \rightarrow 0$ .  $\square$

## D DETAILS AND PROOFS FOR SECTION 4

### D.1 PROOF OF THE NP-HARDNESS RESULT (THEOREM 2)

We now formally address the computational complexity of deciding whether a given number of samples is sufficient to satisfy the  $(\varepsilon, \delta)$ -PAC learning guarantee for all agents. Given a data domain  $\mathcal{X}$ , a hypothesis class  $\mathcal{H}$  of size  $n$ , distributions  $\mathcal{D} = (\mathcal{D}_1, \dots, \mathcal{D}_k)$  for each agent, an accuracy parameter  $\varepsilon > 0$ , a confidence parameter  $\delta \in [0, 1)$ , and an integer  $m$ , we define the *Collaborative PAC Sample-Allocation Problem* as follows:

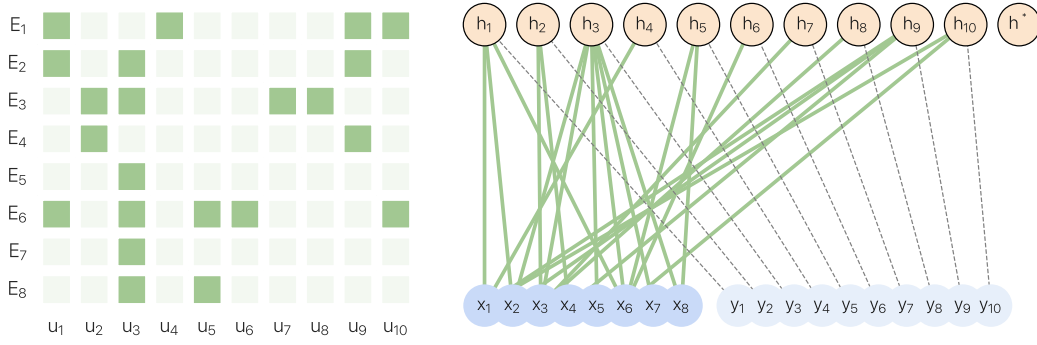


Figure 4: Visualization of the SET COVER reduction used in Theorem 2. **Left:** Incidence matrix of the original Set Cover instance. Each row is a subset  $E_j$  and each column an element  $u_i$ ; a **square** means  $u_i \in E_j$ . **Right:** The corresponding bipartite graph construction: each node  $x_j$  corresponds to subset  $E_j$  and is connected to the hypotheses  $h_i$  for which  $u_i \in E_j$ . A size- $m$  set cover on the left corresponds to an  $m$ -sample training set that forces ERM to output  $h^*$  on the right.

**Definition 3** (Collaborative PAC Sample-Allocation Problem). *Given domain  $\mathcal{X}$ , hypothesis class  $\mathcal{H}$ , distributions  $\mathcal{D} = (\mathcal{D}_1, \dots, \mathcal{D}_k)$ , accuracy parameter  $\varepsilon > 0$ , confidence parameter  $\delta \in [0, 1)$ , and integer  $m$ , decide whether drawing  $m$  samples i.i.d. from each  $\mathcal{D}_i$  suffices to guarantee, with probability at least  $1 - \delta$ , an ERM hypothesis of error at most  $\varepsilon$  for every agent  $i \in \mathcal{A}$ .*

In other words, if each agent contributes exactly  $m_i$  samples from their distribution, will the pooled dataset ensure  $\varepsilon$ -accurate performance on all distributions simultaneously? We now examine the computational complexity of this question. We establish the following computational hardness result:

**Theorem 2.** *Under the PAC-accuracy objective, determining whether a specified sample count  $m$  suffices to meet the  $(\varepsilon, \delta)$ -requirement is NP-hard with respect to the hypothesis-class size  $|\mathcal{H}|$ , even when there is only one agent.*

*Proof.* We reduce from the NP-complete SET COVER problem [35]: given a universe  $U = \{u_1, \dots, u_n\}$ , a family of subsets  $\mathcal{E} = \{E_1, \dots, E_r\}$  with  $\bigcup_{i \in [r]} E_i = U$ , and an integer  $m$ , decide whether at most  $m$  subsets cover  $U$ .

We map an instance  $(U, \mathcal{E}, m)$  of SET COVER to an instance of the Collaborative PAC Sample-Allocation Problem as follows: let the data domain be  $\mathcal{X} = \mathcal{X}_1 \cup \mathcal{X}_2$  where

$$\mathcal{X}_1 = \{x_1, \dots, x_r\}, \quad \mathcal{X}_2 = \{y_1, \dots, y_n\}.$$

Each point  $x_j \in \mathcal{X}_1$  corresponds directly to subset  $E_j \in \mathcal{E}$ , and each point  $y_i \in \mathcal{X}_2$  corresponds to element  $u_i \in U$ . Define the hypothesis class

$$\mathcal{H} = \{h^*, h_1, \dots, h_n\},$$

where, for each  $i \in [n]$  and  $j \in [k]$ :

$$h_i(x_j) = \begin{cases} 1 & \text{if } u_i \in E_j \\ 0 & \text{otherwise} \end{cases}, \quad h_i(y_j) = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}, \quad h^*(x_j) = 0, \quad h^*(y_j) = 0.$$

Set the distribution  $\mathcal{D}$  to be uniform over  $\mathcal{X}$ . Choose  $0 < \varepsilon < 1/(r + n)$  so that a single classification error violates the  $\varepsilon$ -accuracy goal, and set  $\delta$  with  $1 - \delta \geq (|\mathcal{X}|^{-1})^m$ .

If a set cover of size  $p \leq m$  exists, say  $\{E_{j_1}, \dots, E_{j_p}\}$ , consider the sample  $S = \{x_{j_1}, \dots, x_{j_p}\}$ . For every competing hypothesis  $h_i \neq h^*$  there is a subset  $E_{j_q}$  containing  $u_i$ , hence  $h_i(x_{j_q}) = 1 \neq 0 = h^*(x_{j_q})$ . Each  $h_i$  therefore incurs at least one error on  $S$ , so ERM selects  $h^*$ . Thus  $m$  samples suffice to meet the agent's  $(\varepsilon, \delta)$ -requirement.

Conversely, assume that  $m$  samples drawn i.i.d. from  $\mathcal{D}$  suffice to guarantee, with probability at least  $1 - \delta$ , an ERM hypothesis achieving error at most  $\varepsilon$  for any possible true hypothesis.

Due to our choice of  $\varepsilon$ , the ERM hypothesis must correctly classify all points in  $\mathcal{X}$ . Consider the worst-case scenario where the true hypothesis is  $h^*$ . To distinguish  $h^*$  from all other hypotheses  $h_i$ , the training set must include points from  $\mathcal{X}_1$  that differentiate  $h^*$  from each  $h_i$ . Specifically, for each element  $u_i \in U$ , the training set must contain at least one point  $x_j \in \mathcal{X}_1$  with  $h_i(x_j) = 1$  (meaning  $u_i \in E_j$ ), ensuring  $h_i$  is eliminated by ERM. Therefore, the set of points in the minimal training set corresponds directly to subsets forming a valid set cover of size at most  $m$ .

Notably, if any other hypothesis  $h_i \neq h^*$  were the true target, fewer samples would trivially suffice. Thus, the case where  $h^*$  is the true hypothesis indeed represents the worst-case.  $\square$

## D.2 PROOF OF THE APPROXIMATION GUARANTEE (THEOREM 3)

**Theorem 3** (Approximation via Linear Programming). *Given any  $\mathcal{H}$  and  $\varepsilon, \delta > 0$ :*

- For finite  $\mathcal{H}$ , the LP over  $(\mathcal{H}, \varepsilon, \delta)$  returns a  $\frac{\log(1/\delta) + \log |\mathcal{H}|}{\log(1/\delta)}$ -approximate solution to Equation 4.
- For infinite  $\mathcal{H}$ , running the LP over  $(\overline{\mathcal{H}}, \varepsilon, \delta')$  and multiplying it by  $d + \log(1/\delta'')$  returns a  $O\left(\frac{d^2(\log(c_{\max}kd/(c_{\min}\varepsilon\delta))^2)}{\log(1/\delta)}\right)$ -approximate solution to Equation 4, where  $\overline{\mathcal{H}} \subset \mathcal{H}$  is a  $\gamma$ -cover of  $\mathcal{H}$ ,  $\gamma = \Theta(c_{\min}\varepsilon\delta/(c_{\max}k(d + \log(1/\delta))))$ ,  $d = \text{VCdim}(\mathcal{H})$ ,  $\delta'' = \frac{\delta}{4|\overline{\mathcal{H}}|}$  and  $\delta' = \frac{\delta}{8(d + \log(2|\overline{\mathcal{H}}|/\delta))}$ ,  $c_{\min} = \min_{i \in [k]} c_i$  and  $c_{\max} = \max_{i \in [k]} c_i$ .

*Proof.* We start with the case where the hypothesis class  $\mathcal{H}$  is finite, the extension to infinite  $\mathcal{H}$  follows with a standard covering argument.

**Proof for finite  $\mathcal{H}$ .** For any pair of hypotheses  $(h_1, h_2) \in \mathcal{H}$ , define the *disagreement regions* between  $h_1$  and  $h_2$  as

$$\text{DIS}(h_1, h_2) = \{x \in \mathcal{X} \mid h_1(x) \neq h_2(x)\}.$$

This region represents the set of points where  $h_1$  and  $h_2$  disagree. For each agent  $i$ , let  $\mathcal{D}_i(\text{DIS}(h_1, h_2))$  be the probability mass that  $\mathcal{D}_i$  places on that region.

We want to guarantee that with probability at least  $1 - \delta$ , the ERM solution has error  $\leq \varepsilon$  for all agents simultaneously. Concretely, if the true labeler is  $h^* \in \mathcal{H}$ , then *any* other hypothesis  $h_2$  that has  $\mathcal{D}_i(\text{DIS}(h^*, h_2)) > \varepsilon$  (meaning it is “bad” for agent  $i$ ) must be *eliminated* by at least one sample from that disagreement region.

The probability that  $h_2$  is never eliminated, given  $h^*$  is the true labeler, is exactly the probability that *no sample* from any agent  $i$  ever lands in  $\text{DIS}(h^*, h_2)$ . Since each agent  $i$  contributes  $m_i$  i.i.d. points from  $\mathcal{D}_i$ ,

$$\Pr(h_2 \text{ is not eliminated} \mid h^* = h^*) = \prod_{i=1}^k \left(1 - \mathcal{D}_i(\text{DIS}(h^*, h_2))\right)^{m_i}.$$

Hence the requirements  $\Pr(h_2 \text{ not eliminated} \mid h^*) \leq \frac{\delta}{H}$  are log-linear constraints:

$$\sum_{i=1}^k m_i \log[1 - \mathcal{D}_i(\text{DIS}(h^*, h_2))] \leq \log\left(\frac{\delta}{H}\right).$$

We impose such constraints for *all* pairs  $(h^*, h_2)$  where  $\mathcal{D}_i(\text{DIS}(h^*, h_2)) > \varepsilon$  for agent  $i$ . By union bounding across  $\leq H$  possible “bad” hypotheses  $h_2$  (for each  $h^*$ ) or  $\leq H^2$  pairs overall, we ensure that with probability  $\geq 1 - \delta$ , any truly bad hypothesis (in the sense of having error  $> \varepsilon$ ) is eliminated.

Hence, we obtain a linear constraint in terms of  $m_i$ , which allows us to approximate the elimination probability for each hypothesis pair  $(h_1, h_2)$  by solving the following linear program:

$$\min_{\mathbf{m} \in \mathbb{N}^k} \quad \mathbf{c}^\top \mathbf{m}$$

subject to

$$\sum_{i=1}^k m_i \log(1 - \mathcal{D}_i(\text{DIS}(h_1, h_2))) \leq \log \frac{\delta}{H} \quad \forall (h_1, h_2) \in \mathcal{H}^2 \text{ s.t. } \exists i \in [k], \mathcal{D}_i(\text{DIS}(h_1, h_2)) > \varepsilon. \quad (7)$$

One can then round the fractional solution up to integer counts  $m_i$ . Let  $\mathbf{m}^*$  be the solution of the linear program. For any agent  $i$ , we want to show:

$$\forall h^* \in \mathcal{H}, \quad \Pr_{U_j \sim \mathcal{D}_j^{m_j^*}, j \in [k]} \left( \text{err}_{\mathcal{D}_i}^{\text{ERM}} \left( \bigcup_{j \in [k]} U_j \times h^*(U_j) \right) > \varepsilon \right) \leq \delta.$$

This holds because the linear program enforces that, for each pair  $(h_1, h_2) \in \mathcal{H}^2$ , the probability of failing to eliminate any incorrect hypothesis  $h_2$  (given  $h^* = h_1$ ) is bounded by  $\delta/H$ . By applying a union bound over all hypothesis pairs, we achieve the desired bound.

To show that this solution is  $\frac{\log(1/\delta) + \log H}{\log(1/\delta)}$ -approximate optimal, let  $\mathbf{m}$  be the true optimal solution. For each pair  $(h_1, h_2)$  such that  $\mathcal{D}_i(\text{DIS}(h_1, h_2)) > \varepsilon$ , if we take  $\mathbf{m}$  samples, we have:

$$\prod_{i=1}^k (1 - \mathcal{D}_i(\text{DIS}(h_1, h_2)))^{m_i} \leq \delta.$$

Taking logarithms we obtain:

$$\sum_{i=1}^k m_i \log(1 - \mathcal{D}_i(\text{DIS}(h_1, h_2))) \leq \log \delta.$$

Now, given a sample size of  $\frac{\log(1/\delta) + \log H}{\log(1/\delta)} \cdot \mathbf{m}$ , we have:

$$\sum_{i=1}^k \frac{\log(1/\delta) + \log H}{\log(1/\delta)} \cdot m_i \log(1 - \mathcal{D}_i(\text{DIS}(h_1, h_2))) \leq \log \frac{\delta}{H},$$

which satisfies the constraint in [Equation 7](#).

Next, we show how to lift the guarantee to infinite classes using  $\gamma$ -covers and apply our results to this finite cover.

**Proof for infinite  $\mathcal{H}$ .** For any  $\gamma \in (0, 1)$ , we say  $\overline{\mathcal{H}}$  is a  $\gamma$ -cover of  $\mathcal{H}$  under  $\mathcal{D}$  if for all  $h \in \mathcal{H}$ , there exists an  $\bar{h} \in \overline{\mathcal{H}}$  satisfying that  $\max_i \mathcal{D}_i(\text{DIS}(h, \bar{h})) \leq \gamma$ .

**Lemma 2.** Let  $d = \text{VCdim}(\mathcal{H})$  denote the VC dimension of  $\mathcal{H}$ . There is a  $\gamma$ -cover  $\overline{\mathcal{H}} \subset \mathcal{H}$  under  $\mathcal{D}$  of size  $\left(\frac{41k}{\gamma}\right)^d$ .

*Proof of Lemma 2.* In learning theory, for any distribution  $\mathcal{D}$ , a subset  $\overline{\mathcal{H}} \subseteq \mathcal{H}$  is a  $\gamma$ -cover of  $\mathcal{H}$  under  $\mathcal{D}$  if, for every  $h \in \mathcal{H}$ , there exists  $\bar{h} \in \overline{\mathcal{H}}$  such that  $\mathcal{D}(\text{DIS}(h, \bar{h})) \leq \gamma$ . Haussler's sphere-packing bound guarantees [\[36\]](#) the existence of such a cover with

$$|\overline{\mathcal{H}}| \leq \left(\frac{41}{\gamma}\right)^d.$$

Construct a  $\frac{\gamma}{k}$ -cover of  $\mathcal{H}$  under the averaged distribution  $\frac{1}{k} \sum_{i=1}^k \mathcal{D}_i$ . For any  $h \in \mathcal{H}$  and its representative  $\bar{h}$  in this cover,

$$\max_i \mathcal{D}_i(\text{DIS}(\bar{h}, h)) \leq \sum_{i=1}^k \mathcal{D}_i(\text{DIS}(\bar{h}, h)) \leq \gamma,$$

so the same set is a  $\gamma$ -cover under  $\mathcal{D} = (\mathcal{D}_1, \dots, \mathcal{D}_k)$ . Substituting  $\frac{\gamma}{k}$  into Haussler's bound yields

$$|\overline{\mathcal{H}}| \leq \left(\frac{41k}{\gamma}\right)^d.$$

□

Let  $\mathbf{m}^{\overline{\mathcal{H}}, \delta'}$  denote the solution to the LP ([Equation 7](#)) given hypothesis class  $\overline{\mathcal{H}}$  and confidence parameter  $\delta'$ . Let  $c_{\min} = \min_{i \in [k]} c_i$  and  $c_{\max} = \max_{i \in [k]} c_i$ .

**Lemma 3.** By choosing  $\gamma = \Theta(c_{\min}\varepsilon\delta/(c_{\max}k(d + \log(1/\delta))))$  and  $\mathbf{m} = \mathbf{m}^{\overline{\mathcal{H}}, \delta/2}$ , for any target hypothesis  $h^* \in \mathcal{H}$  and agent  $i$ , with probability at least  $1 - \delta$ , any consistent hypothesis  $\bar{h} \in \overline{\mathcal{H}}$  will satisfy

$$\text{err}_{\mathcal{D}_i, h^*}(\bar{h}) \leq \varepsilon.$$

*Proof of Lemma 3.* Since each agent can achieve the PAC learning objective with at most  $O((d + \log(1/\delta))/\varepsilon)$  data points individually, we restrict attention to contribution vectors satisfying  $\|\mathbf{m}\|_1 \leq O(\frac{c_{\max}k(d + \log(1/\delta))}{c_{\min}\varepsilon})$  (otherwise we can replace the solution with  $O((d + \log(1/\delta))/\varepsilon)\mathbf{1}$ ).

For any  $\mathbf{m}$  satisfying this constraint, any  $h^* \in \mathcal{H}$ , and any agent  $i$ , let the labeled sample be

$$S^{h^*} = \bigcup_{j \in [k]} \{(x, h^*(x)) \mid x \in U_j\}.$$

Define

$$\overline{\text{err}} = \max \left\{ \text{err}_{\mathcal{D}_i}(\bar{h}) \mid \bar{h} \in \overline{\mathcal{H}}, \text{err}_{S^{h^*}}(\bar{h}) = \min_{h' \in \overline{\mathcal{H}}} \text{err}_{S^{h^*}}(h') \right\},$$

the worst-case error when running ERM over  $\overline{\mathcal{H}}$ . We can show that  $\overline{\text{err}}$  is small with high probability. More specifically, assuming a  $(\varepsilon, \delta/2)$ -PAC guarantee for  $\overline{\mathcal{H}}$ , we have

$$\begin{aligned} & \Pr_{S^{h^*}} \left[ \exists \bar{h} : \text{err}_{\mathcal{D}_i}(\bar{h}) > \varepsilon, \text{err}_{S^{h^*}}(\bar{h}) = \min_{h' \in \overline{\mathcal{H}}} \text{err}_{S^{h^*}}(h') \right] \\ & \leq \Pr_{S^{h^*}} \left[ \exists \bar{h} : \text{err}_{\mathcal{D}_i}(\bar{h}) > \varepsilon, \text{err}_{S^{h^*}}(\bar{h}) = 0, \text{err}_{S^{h^*}}(\bar{h}^*) = 0 \right] + \Pr_{S^{h^*}} \left[ \text{err}_{S^{h^*}}(\bar{h}^*) \neq 0 \right] \\ & \leq \Pr_{S^{h^*}} \left[ \exists \bar{h} : \text{err}_{\mathcal{D}_i}(\bar{h}) > \varepsilon, \text{err}_{S^{h^*}}(\bar{h}) = 0 \right] + \Pr_{S^{h^*}} \left[ \text{err}_{S^{h^*}}(\bar{h}^*) \neq 0 \right] \\ & \leq \frac{\delta}{2} + (1 - (1 - \gamma)^{\sum_j m_j}) \leq \frac{\delta}{2} + (1 - (1 - \gamma)^{\frac{c_{\max}k(d + \log(1/\delta))}{c_{\min}\varepsilon}}). \end{aligned}$$

Choosing  $\gamma = c_{\min}\varepsilon\delta/(c_{\max}k(d + \log(1/\delta)))$  yields

$$\frac{\delta}{2} + \left( 1 - (1 - \gamma)^{\frac{c_{\max}k(d + \log(1/\delta))}{c_{\min}\varepsilon}} \right) \leq \delta,$$

so for any  $h^*$  and agent  $i$ , with probability at least  $1 - \delta$  we can always find a hypothesis whose error is at most  $\varepsilon$ .  $\square$

**Lemma 4.** By choosing  $\gamma = \Theta(c_{\min}\varepsilon\delta/(c_{\max}k(d + \log(1/\delta))))$ , the solution  $(d + \log(1/\delta'')) \cdot \mathbf{m}^{\overline{\mathcal{H}}, \delta'}$  is sufficient to achieve  $(\varepsilon, \delta)$ -PAC accuracy objective (Equation 2) for  $\mathcal{H}$ , where  $\delta'' = \frac{\delta}{4|\overline{\mathcal{H}}|}$  and  $\delta' = \frac{\delta}{8(d + \log(2|\overline{\mathcal{H}}|/\delta))}$ .

*Proof of Lemma 4.* Similar to Lemma 3, we again restrict attention to contribution vectors satisfying  $\|\mathbf{m}\|_1 \leq O(\frac{c_{\max}k(d + \log(1/\delta))}{c_{\min}\varepsilon})$ . This guarantee that for any target  $h^*$ , with probability at least  $1 - \frac{\delta}{2}$ ,  $\bar{h}^*$  is consistent. Hence, we can view  $\bar{h}^*$  as our target hypothesis. We now prove that for any target hypothesis  $\bar{h}^* \in \overline{\mathcal{H}}$  and any representative hypothesis  $\bar{h} \in \overline{\mathcal{H}}$  satisfying  $\exists i \in [k], \mathcal{D}_i(\text{DIS}(\bar{h}^*, \bar{h})) > \varepsilon$ , with probability at least  $1 - \frac{\delta'}{2|\overline{\mathcal{H}}|}$ , any  $h$  in the  $\gamma$ -ball of  $\bar{h}$  will be eliminated.

Let  $\beta_i = \mathcal{D}_i(\text{DIS}(\bar{h}^*, \bar{h}))$  denote the probability mass of  $\text{DIS}(\bar{h}^*, \bar{h})$  under agent  $i$ 's data distribution and let  $A_\gamma = \{i \in [k] \mid \beta_i \geq 4\gamma\}$  denote the set of agents whose probability mass of the disagreement region is at least  $4\gamma$ . The approximation solution  $\mathbf{m}^{\overline{\mathcal{H}}, \delta'}$  guarantees that

$$\prod_{i=1}^k (1 - \beta_i)^{m_i^{\overline{\mathcal{H}}, \delta'}} \leq \frac{\delta'}{|\overline{\mathcal{H}}|}.$$

If  $\beta_i < 4\gamma$ , we have

$$(1 - \beta_i)^{m_i^{\overline{\mathcal{H}}, \delta'}} > (1 - 4\gamma)^{m_i^{\overline{\mathcal{H}}, \delta'}} \geq (1 - 4\gamma)^{\delta'/\gamma} \geq e^{-(4 \log 4) \delta'},$$

where we adopt the inequality  $(1 - 1/x)^x \geq \frac{1}{4}$  for all  $x \geq 2$ . Hence, we have

$$e^{-(4 \log 4)k\delta'} \cdot \prod_{i:\beta_i \geq 4\gamma} (1 - \beta_i)^{m_i^{\overline{\mathcal{H}}, \delta'}} \leq \frac{\delta'}{|\overline{\mathcal{H}}|}.$$

By rearranging terms, we have

$$\prod_{i:\beta_i \geq 4\gamma} (1 - \beta_i)^{m_i^{\overline{\mathcal{H}}, \delta'}} \leq \frac{\delta'}{|\overline{\mathcal{H}}|} \cdot e^{(4 \log 4)k\delta'} \leq \frac{2\delta'}{|\overline{\mathcal{H}}|},$$

when  $\delta' = O(\frac{1}{k})$ . That is to say, even if we only collect  $m_i$  samples from agent  $i \in A_\gamma$  and don't collect from  $i \notin A_\gamma$ , we will still be able to obtain one sample from the  $\text{DIS}(\overline{h}^*, \overline{h})$  with probability at least  $1 - \frac{2\delta'}{|\overline{\mathcal{H}}|}$ .

For any distribution  $\mathcal{D}$ , let  $\mathcal{D}_{|\overline{h}^*, \overline{h}}$  denote the distribution restricted to  $\text{DIS}(\overline{h}^*, \overline{h})$ , i.e.,

$$\mathcal{D}_{|\overline{h}^*, \overline{h}}(x) = \begin{cases} \frac{\mathcal{D}(x)}{\mathcal{D}(\text{DIS}(\overline{h}^*, \overline{h}))} & \text{if } x \in \text{DIS}(\overline{h}^*, \overline{h}), \\ 0 & \text{otherwise.} \end{cases}$$

Consider  $\mathcal{D} = \frac{\sum_{i \in A_\gamma} m_i \mathcal{D}_i}{\sum_{i \in A_\gamma} m_i}$  being a mixture of  $\{\mathcal{D}_i | i \in A_\gamma\}$ , the disagreement region between  $\overline{h}^*$  and  $h$  under distribution  $\mathcal{D}_{|\overline{h}^*, \overline{h}}$  is at least  $\frac{3}{4}$ . This is because

$$\mathcal{D}_{|\overline{h}^*, \overline{h}}(\text{DIS}(\overline{h}^*, h)) = 1 - \mathcal{D}_{|\overline{h}^*, \overline{h}}(\text{DIS}(\overline{h}, h)) \geq 1 - \frac{\gamma}{4\gamma} = \frac{3}{4}.$$

Hence, by standard PAC learning guarantee, if we obtain  $n = O(d + \log(1/\delta''))$  samples from  $\mathcal{D}_{|\overline{h}^*, \overline{h}}$ , then with probability  $1 - \delta''$ , all  $h$ 's in the  $\gamma$ -ball of  $\overline{h}$  are not consistent.

The current approximation solution  $\mathbf{m}^{\overline{\mathcal{H}}, \delta'}$  can only guarantee that with probability  $1 - \frac{2\delta'}{|\overline{\mathcal{H}}|}$ , we obtain one sample from  $\mathcal{D}_{|\overline{h}^*, \overline{h}}$ , i.e.,

$$\prod_{i:\beta_i \geq 4\gamma} (1 - \beta_i)^{m_i} \leq \frac{2\delta'}{|\overline{\mathcal{H}}|}.$$

When we increase  $\mathbf{m}^{\overline{\mathcal{H}}, \delta'}$  by  $n$  times, then with probability at least  $1 - \frac{2n\delta'}{|\overline{\mathcal{H}}|}$ , we obtain  $n$  samples from  $\mathcal{D}_{|\overline{h}^*, \overline{h}}$ .

Hence, by setting  $\delta'' = \frac{\delta}{4|\overline{\mathcal{H}}|}$  and  $\delta' = \frac{\delta}{8(d + \log(2|\overline{\mathcal{H}}|/\delta))}$ , solving the LP for  $\overline{\mathcal{H}}, \delta'$  to obtain an approximate solution  $\mathbf{m}^{\overline{\mathcal{H}}, \delta'}$  and multiplying it by  $d + \log(1/\delta'')$ , the solution  $(d + \log(1/\delta'')) \cdot \mathbf{m}^{\overline{\mathcal{H}}, \delta'}$  is sufficient to achieve  $(\varepsilon, \delta)$ -PAC accuracy objective (Equation 2) for  $\mathcal{H}$ .  $\square$

The solution  $\mathbf{m}^{\overline{\mathcal{H}}, \delta'}$  returned by the approximation algorithm for the  $(\varepsilon, \delta')$ -PAC objective for  $\overline{\mathcal{H}}$  is at most a factor

$$\begin{aligned} \frac{\log(|\overline{\mathcal{H}}|/\delta')}{\log(1/\delta)} &= \frac{\log(8/\delta) + \log(d + \log(2|\overline{\mathcal{H}}|/\delta)) + d \log(k(d + \log(1/\delta)))/(\varepsilon\delta)}{\log(1/\delta)} \\ &= O\left(\frac{d(\log k + \log d + \log(1/\varepsilon) + \log(1/\delta) + \log(c_{\max}/c_{\min}))}{\log(1/\delta)}\right) \end{aligned}$$

larger than the optimal solution for the  $(\varepsilon, \delta)$ -PAC objective for  $\overline{\mathcal{H}}$ . Since the PAC objective for  $\overline{\mathcal{H}}$  is easier than the PAC objective for  $\mathcal{H}$ , we have  $(d + \log(1/\delta'')) \cdot \mathbf{m}^{\overline{\mathcal{H}}, \delta'}$  is a  $O\left(\frac{d^2(\log k + \log d + \log(1/\varepsilon) + \log(1/\delta) + \log(c_{\max}/c_{\min}))^2}{\log(1/\delta)}\right)$ -approximation solution for  $(\varepsilon, \delta)$ -PAC objective for  $\mathcal{H}$ .  $\square$

## E DETAILS AND PROOFS FOR SECTION 5

**Lemma 1.** *For any  $H$  greater than a universal constant  $C$ , there exists a PAC learning instance of  $(\mathcal{H}, \varepsilon, \delta)$  in the two-agent setting with  $|\mathcal{H}| = H$  such that the approximation algorithm introduced in [Theorem 3](#) is oblivious at all  $\mathbf{m}$  with  $m_1, m_2 \geq 2|\mathcal{H}| \log |\mathcal{H}|$ .*

*Proof of Lemma 1.* Consider  $\mathcal{X} = \{1, 2, \dots, n\}$  as our domain and the hypothesis class  $\mathcal{H}$  of all singletons (i.e. functions  $h_i$  that label exactly one  $x \in \mathcal{X}$  as positive and all others as negative) plus the all-negative function. Thus  $H := |\mathcal{H}| = n + 1$ .

We now show that that for any  $\mathbf{m}$  satisfying  $m_1, m_2 \geq 2|\mathcal{H}| \log |\mathcal{H}|$  and any neighbor  $\mathbf{m}'$  with  $\|\mathbf{m}' - \mathbf{m}\|_1 = 1$ , we can always find distributions  $\mathcal{D}_1, \mathcal{D}_1', \mathcal{D}_2, \mathcal{D}_2'$  satisfying that

- $\mathbf{m} = \text{APPROX}(\mathcal{D}_1, \mathcal{D}_2)$  and  $\mathbf{m}' = \text{APPROX}(\mathcal{D}_1', \mathcal{D}_2) = \text{APPROX}(\mathcal{D}_1, \mathcal{D}_2')$ .
- Contributions  $\mathbf{m}$  and  $\text{APPROX}(\mathcal{D}_1', \mathcal{D}_2)'$  are both feasible for  $(\mathcal{D}_1, \mathcal{D}_2)$ ,  $(\mathcal{D}_1', \mathcal{D}_2)$ , and  $(\mathcal{D}_1, \mathcal{D}_2')$ .

**Construction.** We construct a pair of distributions  $(\mathcal{D}_1, \mathcal{D}_2)$  by selecting  $p_1, p_2, q_1, q_2$  satisfying

$$p_1 \leq p_2 \leq \frac{1}{2n}, \quad q_2 \leq q_1 \leq \frac{1}{2n}, \quad p_1 + p_2 \leq \frac{1}{n}, \quad q_1 + q_2 \leq \frac{1}{n}.$$

Then let

$$p_3 = \dots = p_n = \frac{1 - (p_1 + p_2)}{n - 2}, \quad q_3 = \dots = q_n = \frac{1 - (q_1 + q_2)}{n - 2}.$$

Hence let  $\mathcal{D}_1 = (p_1, \dots, p_n)$  and  $\mathcal{D}_2 = (q_1, \dots, q_n)$ . Both are well-defined distributions. We denote by  $\mathcal{P}$  the family of all such pairs  $(\mathcal{D}_1, \mathcal{D}_2)$  with different choices of  $p_1, p_2, q_1, q_2$  satisfying the above constraints.

The approximate-optimal solution  $\text{APPROX}(\mathcal{D}_1, \mathcal{D}_2)$  is found by solving a linear program derived from disagreement regions. For example, the disagreement region between the all-negative hypothesis  $h_0$  and a singleton  $h_1$  that is positive on point 1 has measure  $p_1$  in  $\mathcal{D}_1$  (and  $q_1$  in  $\mathcal{D}_2$ ). The “binding constraints” come from pairs  $(h_0, h_1)$  and  $(h_0, h_2)$  (the singletons with points 1 or 2). By contrast, any singleton on point  $i \geq 3$  yields a disagreement measure at least  $p_3$ , which is larger, and thus any feasible solution that satisfies the “small measure” constraints with some margin also satisfies these larger measure constraints.

Thus the main LP constraints reduce to requiring that

$$m_1 \log \frac{1}{1 - p_1} + m_2 \log \frac{1}{1 - q_1} \geq \log\left(\frac{H}{\delta}\right), \quad (8)$$

$$m_1 \log \frac{1}{1 - p_2} + m_2 \log \frac{1}{1 - q_2} \geq \log\left(\frac{H}{\delta}\right). \quad (9)$$

The solution  $\text{APPROX}(\mathcal{D}_1, \mathcal{D}_2)$  is the one that *minimizes*  $c_1 \cdot m_1 + c_2 \cdot m_2$  subject to [Equations \(8\) and \(9\)](#) (and additional constraints for any bigger disagreements, satisfied by slack).

If  $m_1 + m_2 \geq 2H \log H$ , we can show:

$$m_1 \log\left(\frac{1}{1 - p_i}\right) + m_2 \log\left(\frac{1}{1 - q_i}\right) > m_1 \log\left(\frac{1}{1 - p_1}\right) + m_2 \log\left(\frac{1}{1 - q_1}\right) + \log H$$

for each  $i \geq 3$ , hence those constraints are looser. Hence, satisfying [Equations \(8\) and \(9\)](#) by a small margin also satisfies the bigger-disagreement constraints.

**To satisfy the second bullet.** For any  $i = 3, 4, \dots$ , according to our construction, we have  $p_i \geq \frac{1}{n}$  and  $p_1, p_2 \leq \frac{1}{2n}$ . Thus, we have

$$\begin{aligned} 1 - p_1 &\geq 1 - \frac{1}{2n} = e^{-\frac{1}{2n}} - O\left(\frac{1}{n^2}\right) = e^{-\frac{1}{n}} e^{\frac{1}{2n}} - O\left(\frac{1}{n^2}\right) \geq \left(1 - \frac{1}{n}\right) e^{\frac{1}{2n}} - O\left(\frac{1}{n^2}\right) \\ &\geq (1 - p_i) e^{\frac{1}{2n}} - O\left(\frac{1}{n^2}\right). \end{aligned}$$

Thus, we have

$$\log\left(\frac{1}{1 - p_i}\right) \geq \log\left(\frac{1}{1 - p_1}\right) + \frac{1}{2n} - O\left(\frac{1}{n^2}\right) \geq \log\left(\frac{1}{1 - p_1}\right) + \frac{1}{3n}$$

for  $n$  big enough. Then for the constraint corresponding to the disagreement region  $\text{DIS}(h_0, h_i)$ , we have

$$\begin{aligned} m_1 \log \frac{1}{1-p_i} + m_2 \log \frac{1}{1-q_i} &\geq m_1 \log \frac{1}{1-p_1} + m_2 \log \frac{1}{1-q_1} + \frac{m_1 + m_2}{3n} \\ &\geq m_1 \log \frac{1}{1-p_1} + m_2 \log \frac{1}{1-q_1} + \log H, \end{aligned}$$

where the last inequality holds when  $m_1 + m_2 \geq 3H \log H$ . Similar results also hold for  $\mathcal{D}_2$ . That is to say, instead of satisfying Equations (8) and (9), it would be sufficient to satisfy

$$\begin{aligned} m_1 \log \frac{1}{1-p_1} + m_2 \log \frac{1}{1-q_1} &\geq \log \frac{3}{\delta}, \\ m_1 \log \frac{1}{1-p_2} + m_2 \log \frac{1}{1-q_2} &\geq \log \frac{3}{\delta}. \end{aligned}$$

Hence, for any  $(\mathcal{D}_1, \mathcal{D}_2) \in \mathcal{P}$ , we have  $\frac{\log 3/\delta}{\log H/\delta} \cdot \text{APPROX}(\mathcal{D}_1, \mathcal{D}_2)$  is feasible. For simplicity, let's fix  $\delta \leq 0.5$  and suppose  $H \geq 18$  from now. Then  $\frac{1}{2} \cdot \text{APPROX}(\mathcal{D}_1, \mathcal{D}_2)$  is sufficient and so is  $\text{APPROX}(\mathcal{D}_1, \mathcal{D}_2) - (1, 1)$ . Thus, we justify bullet 2.

**To satisfy the first bullet.** Given an  $\mathbf{m}$  and  $\mathbf{m}'$ , we pick  $\mathcal{D}_1$  specified by  $p_1, p_2$ ,  $\mathcal{D}_2$  specified by  $q_1, q_2$  and  $\mathcal{D}'_1$  specified by  $p'_1, p'_2$  so that Inequalities Equations (8) and (9) hold with equality and that  $\mathbf{m}$  and  $\mathbf{m}'$  is the only solution to these linear equalities w.r.t.  $(p_1, p_2, q_1, q_2)$  and w.r.t.  $(p'_1, p'_2, q_1, q_2)$ , respectively. Inequalities Equations (8) and (9) can be approximated using a first-order approximation as follows.

$$\begin{aligned} m_1 p_1 + m_2 q_1 &= \log \frac{H}{\delta} =: \alpha, \\ m_1 p_2 + m_2 q_2 &= \alpha, \\ m'_1 p'_1 + m'_2 q_1 &= \alpha, \\ m'_1 p'_2 + m'_2 q_2 &= \alpha. \end{aligned}$$

Let's pick

$$\begin{aligned} p_1 &= \frac{\alpha}{m_1 m_2}, & q_1 &= \frac{(1 - 1/m_2)\alpha}{m_2}, \\ p_2 &= \frac{(1 - 1/m_1)\alpha}{m_1}, & q_2 &= \frac{\alpha}{m_1 m_2}. \end{aligned}$$

We can justify  $(\mathcal{D}_1, \mathcal{D}_2) \in \mathcal{P}$  since  $m_1, m_2 \geq 2n \log n$ . Then if  $\mathbf{m}'$  differs from  $\mathbf{m}$  at  $m_1$ , by solving

$$\begin{aligned} m'_1 p'_1 + m_2 q_1 &= m_1 p_1 + m_2 q_1, \\ m'_1 p'_2 + m_2 q_2 &= m_1 p_2 + m_2 q_2, \end{aligned}$$

we have

$$p'_1 = \frac{m_1 p_1}{m'_1}, \quad p'_2 = \frac{m_1 p_2}{m'_1}.$$

It's easy for us to justify that  $(\mathcal{D}'_1, \mathcal{D}_2) \in \mathcal{P}$  since  $\mathcal{D}'_1$  is very close to  $\mathcal{D}_1$ .

If  $\mathbf{m}'$  differs from  $\mathbf{m}$  at  $m_2$ , by solving

$$\begin{aligned} m_1 p'_1 + m'_2 q_1 &= m_1 p_1 + m_2 q_1, \\ m_1 p'_2 + m'_2 q_2 &= m_1 p_2 + m_2 q_2, \end{aligned}$$

we have

$$p'_1 = \frac{(m_2 - m'_2)q_1}{m_1} + p_1, \quad p'_2 = \frac{(m_2 - m'_2)q_2}{m_1} + p_2.$$

By plugging in the values of  $p_1, p_2, q_1, q_2$ , we have

$$\begin{aligned} p'_1 &= \frac{(\pm 1)(1 - 1/m_2)\alpha}{m_1 m_2} + \frac{\alpha}{m_1 m_2}, \\ p'_2 &= \frac{(\pm 1)\alpha}{m_1^2 m_2} + \frac{(1 - 1/m_1)\alpha}{m_1}. \end{aligned}$$

It is easy to see that  $p'_2$  is very close to  $p_2$  and  $p'_1 < p'_2$ . For  $p'_1$ , we need to make it fall in  $(0, \frac{1}{2n})$ . When  $m_1, m_2 \geq 2n \log n$ , we have  $p'_1 \leq \frac{2}{m_1 m_2} < \frac{1}{2n}$ . Also,  $p'_1$  is always positive. So we are done with computing  $\mathcal{D}'_1$ . We can compute  $\mathcal{D}'_2$  in the same way.  $\square$

*Proof of Theorem 4.* When  $\mathbf{f}$  is strategyproof, it must hold that  $f_i(\mathbf{m}) - c_i \cdot m_i = f_i(\mathbf{m}') - c_i \cdot m'_i$  for any two neighboring  $\mathbf{m}, \mathbf{m}' \in M$ . If  $f_i(\mathbf{m}) - c_i \cdot m_i > f_i(\mathbf{m}') - c_i \cdot m'_i$ , agent  $i$  will misreport their distribution when the ground truth is  $(\mathcal{D}'_i, \mathcal{D}_{-i})$ ; else if  $f_i(\mathbf{m}) - c_i \cdot m_i < f_i(\mathbf{m}') - c_i \cdot m'_i$ , agent  $i$  will misreport when the ground truth is  $\mathcal{D}$ , which conflicts with truthfulness. Since  $M$  is connected, we have  $f_i(\mathbf{m}) - c_i \cdot m_i = C_i$  for all  $\mathbf{m} \in M$ .  $\square$

## F RESULTS FOR EXPECTED ACCURACY OBJECTIVE

### F.1 APPROXIMATION ALGORITHM

Recall that the optimization problem with the expected accuracy objective given an error parameter  $\varepsilon$  is

$$\min_{\mathbf{m} \in \mathbb{N}^k} \mathbf{c}^\top \mathbf{m}$$

subject to

$$\max_{h^* \in \mathcal{H}} \mathbb{E}_{S \sim \mathcal{P}(\mathcal{D}, \mathbf{m}, h^*)} [\text{err}_{\mathcal{D}_i, h^*}^{\text{ERM}}(S)] \leq \varepsilon, \forall i \in \mathcal{A}. \quad (10)$$

Let's denote the optimal solution to the above problem as  $\mathbf{m}^{\star, \text{exp}}(\varepsilon)$ .

For any pair  $(h_j, h_t)$  write

$$E_{j,t} := \left\{ \text{no sample lies in } \text{DIS}(h_j, h_t) \right\}, \quad a_{j,t}^i := \mathcal{D}_i(\text{DIS}(h_j, h_t)),$$

the “no-sample” event and its probability mass under agent  $i$ 's distribution  $\mathcal{D}_i$ .

Without loss of generality, relabel the hypotheses so that the disagreement masses are non-increasing:

$$a_1^i \geq a_2^i \geq \dots \geq a_K^i.$$

Define the first “small” index by

$$n := \min \{ j \mid a_j^i \leq \frac{\varepsilon}{2} \} \quad (\text{set } n = K \text{ if no such } j \text{ exists}).$$

$$\Lambda_{<j} := \bigwedge_{k < j} \neg E_k, \quad \Delta_{\leq j} := \bigvee_{k \leq j} E_k, \quad a_{K+1}^i := 0.$$

These abbreviations mean, respectively, “no disagreement observed yet” and “some disagreement observed by step  $j$ .”

$$\begin{aligned} \mathbb{E}_{S \sim \mathcal{P}(\mathcal{D}, \mathbf{m}, h^*)} [\text{err}_{\mathcal{D}_i, h^*}^{\text{ERM}}(S)] &= \sum_{j=1}^{H-1} \Pr(E_j \wedge \Lambda_{<j}) a_j^i \\ &\leq \sum_{j=1}^n \Pr(E_j \wedge \Lambda_{<j}) a_j^i + \frac{\varepsilon}{2} \\ &= \sum_{j=1}^n \Pr(\Delta_{\leq j}) (a_j^i - a_{j+1}^i) + \frac{\varepsilon}{2}. \end{aligned}$$

For any index set  $[j] = \{1, \dots, j\}$  let

$$\Delta_{\leq j} := \bigvee_{t \leq j} E_t.$$

Then

$$\sup_{t \leq j} \Pr(E_t) \leq \Pr(\Delta_{\leq j}) \leq \sum_{t=1}^j \Pr(E_t) \leq j \sup_{t \leq j} \Pr(E_t). \quad (\text{UB})$$

Replacing  $\Pr(\Delta_{\leq j})$  by its upper bound  $\sum_{t=1}^j \Pr(E_t)$  in the telescoping sum from the previous step yields

$$\begin{aligned} & \sum_{j=1}^n \left( \sum_{t=1}^j \Pr(E_t) \right) (a_j^i - a_{j+1}^i) + \frac{\varepsilon}{2} \\ &= \sum_{t=1}^n \Pr(E_t) a_t^i + \frac{\varepsilon}{2} \leq \varepsilon, \end{aligned}$$

where we set  $a_{n+1}^i = 0$  for the telescoping identity. Each probability

$$\Pr(E_t) = \prod_{k=1}^K (1 - \mathcal{D}_k(\text{DIS}(h_0, h_t)))^{m_k}$$

is a product of log-affine functions of  $\mathbf{m} = (m_1, \dots, m_K)$  and is therefore convex; the entire left-hand side above is a non-negative weighted sum of convex functions, hence convex in  $\mathbf{m}$ .

Impose the term-wise bound  $\Pr(E_t) a_t^i \leq \frac{\varepsilon}{2H}$  for every disagreement index  $t$ . Define the set of “large-mass” pairs

$$\mathcal{P}_{\varepsilon/2} := \left\{ (h_1, h_2) \in \mathcal{H}^2 : \max_{i \leq k} \mathcal{D}_i(\text{DIS}(h_1, h_2)) > \frac{\varepsilon}{2} \right\}.$$

For each such pair write

$$p_i(h_1, h_2) := \mathcal{D}_i(\text{DIS}(h_1, h_2))$$

and

$$a(h_1, h_2) := \min_{i: p_i(h_1, h_2) > \varepsilon/2} p_i(h_1, h_2).$$

The resulting LP is

$$\begin{aligned} & \min_{\mathbf{m} \in \mathbb{N}^k} \mathbf{c}^\top \mathbf{m} \\ & \text{s.t.} \quad \sum_{i=1}^k m_i \log(1 - p_i(h_1, h_2)) \leq \log\left(\frac{\varepsilon}{2H a(h_1, h_2)}\right), \quad \forall (h_1, h_2) \in \mathcal{P}_{\varepsilon/2}. \end{aligned} \tag{11}$$

**Theorem 7.** *The solution to Equation 11 is a feasible solution to Equation 10. Given the optimal solution  $\mathbf{m}^{\star, \text{exp}}(\frac{\varepsilon}{4})$  to Equation 10 with error parameter  $\frac{\varepsilon}{4}$ , then  $\log(2H)\mathbf{m}^{\star, \text{exp}}(\frac{\varepsilon}{4})$  is a feasible solution to Equation 11.*

*Proof of Theorem 7.* It is direct to see the solution to Equation 11 is a feasible solution to Equation 10. Given contribution  $\mathbf{m} = \mathbf{m}^{\star, \text{exp}}(\frac{\varepsilon}{4})$ , we have  $\Pr(E_{j,t}) a_{j,t}^i \leq \frac{\varepsilon}{4}$  for all  $h_j, h_t \in \mathcal{H}$ . That is to say,

$$\prod_{i'=1}^k (1 - a_{j,t}^{i'})^{m_{i'}} \leq \frac{\varepsilon}{4a_{j,t}^i}.$$

For  $a_{j,t}^i > \frac{\varepsilon}{2}$ , since  $\log(2H) = \log H + 1 \geq \frac{\log(a_{j,t}^i H / \varepsilon)}{\log(4a_{j,t}^i / \varepsilon)} = \frac{\log H + \log(a_{j,t}^i / \varepsilon)}{\log(4a_{j,t}^i / \varepsilon)}$ , we have

$$\prod_{i'=1}^k (1 - a_{j,t}^{i'})^{m_{i'} \cdot \log(2H)} \leq \prod_{i'=1}^k (1 - a_{j,t}^{i'})^{m_{i'} \cdot \frac{\log(a_{j,t}^i H / \varepsilon)}{\log(4a_{j,t}^i / \varepsilon)}} \leq \frac{\varepsilon}{a_{j,t}^i H}.$$

□

**The performance of running ERM over the cover  $\bar{\mathcal{H}}$  for the expected accuracy objective.** For any  $\mathcal{H}$ , let  $\bar{\mathcal{H}}$  be a  $\gamma$ -cover of  $\mathcal{H}$ . Then for any contribution  $\mathbf{m}$  satisfying this constraint, any  $h^* \in \mathcal{H}$  and any agent  $i$ , given the labeled data  $S^{h^*} = \cup_{j \in [k]} \{(x, h^*(x)) | x \in U_j\}$ , let  $\bar{\text{err}} = \max\{\text{err}_{\mathcal{D}_i}(\bar{h}) | \bar{h} \in \bar{\mathcal{H}} : \text{err}_{S^{h^*}}(\bar{h}) = \min_{h' \in \bar{\mathcal{H}}} \text{err}_{S^{h^*}}(h')\}$  denote the worst case error when running ERM over  $\bar{\mathcal{H}}$ . We can show that  $\bar{\text{err}}$  is small in expectation. More specifically, when we can achieve  $\varepsilon/2$ -expected accuracy guarantee for  $\bar{\mathcal{H}}$ , we have

When we can achieve  $\varepsilon$  expected accuracy guarantee for the cover  $\bar{\mathcal{H}}$ , we have

$$\begin{aligned} \mathbb{E}[\bar{\text{err}}] &\leq \mathbb{E}\left[\bar{\text{err}} \mid \text{err}_{S^{h^*}}(\bar{h}^*) = 0\right] \mathbb{P}\left(\text{err}_{S^{h^*}}(\bar{h}^*) = 0\right) + \mathbb{P}\left(\text{err}_{S^{h^*}}(\bar{h}^*) \neq 0\right) \\ &\leq \mathbb{E}\left[\bar{\text{err}} \mid \text{err}_{S^{h^*}}(\bar{h}^*) = 0\right] + \mathbb{P}\left(\text{err}_{S^{h^*}}(\bar{h}^*) \neq 0\right) \\ &\leq \varepsilon/2 + (1 - (1 - \gamma)^{\frac{k \text{VCdim}(\mathcal{H})}{\varepsilon}}). \end{aligned}$$

By setting  $\gamma = O(\frac{\varepsilon^2}{k \cdot \text{VCdim}(\mathcal{H})})$ , we have  $\mathbb{E}[\bar{\text{err}}] \leq \varepsilon$ .

## F.2 LOCAL OBLIVIOUSNESS OF THE APPROXIMATION ALGORITHM

Note that the image space of this approximation algorithm is  $[0, \frac{2 \log(H)}{\varepsilon}]^k$  according to Equation 11. Then we show that most area of  $[0, \frac{2 \log(H)}{\varepsilon}]^k$  is oblivious.

**Lemma 5.** *For any  $k < H \in \mathbb{N}$  and  $\varepsilon < 1 - \frac{1}{2H}$ , there exists an instance of expected accuracy learning instance of  $(\mathcal{H}, \varepsilon)$  with  $|\mathcal{H}| = H$  in the  $k$ -agent setting such that the approximation algorithm introduced in Theorem 7 is oblivious at all  $\mathbf{m} \in [2 + \frac{\log(1/\varepsilon)}{\log(2H)}, \frac{2 \log(H)}{\varepsilon}]^k$ .*

*Proof.* Let the input space be the  $H$  points  $\mathcal{X} = \{x_1, \dots, x_H\}$ . The hypothesis class  $\mathcal{H}$  consists of all singletons over  $\{x_1, \dots, x_{H-1}\}$ —denote them  $h_1, \dots, h_{H-1}$ —together with the all-negative hypothesis  $h_0$ . Fix a contribution vector  $\mathbf{m} = (m_1, \dots, m_k)$ . For each agent  $i \in [k]$  choose  $c_i \in [\varepsilon/2, 1 - \frac{1}{2H}]$  that solves

$$(1 - c_i)^{m_i} c_i = \frac{\varepsilon}{2H}. \quad (1)$$

Such a solution exists whenever

$$m_i \in \left[1 + \frac{\log(1/\varepsilon)}{\log(2H)}, \frac{2 \log H}{\varepsilon}\right].$$

Define  $\mathcal{D}_i$  by setting  $\mathcal{D}_i(x_i) = c_i$  and  $\mathcal{D}_i(x_H) = 1 - c_i$ .

Fix an agent  $i$  and a hypothesis  $h_j$  with  $j \in [H - 1]$ . Under  $h_j$  the expected ERM error of agent  $i$  equals

$$\mathbb{E}_{S \sim \mathcal{P}(\mathcal{D}, \mathbf{m}, h_j)}[\text{err}_{\mathcal{D}_i, h_j}^{\text{ERM}}(S)] = \Pr[x_i \text{ and } x_j \text{ both unseen in } S] c_i \leq \Pr[x_i \text{ unseen in } S] c_i.$$

For the all-negative hypothesis  $h_0$  the same bound holds:

$$\mathbb{E}_{S \sim \mathcal{P}(\mathcal{D}, \mathbf{m}, h_0)}[\text{err}_{\mathcal{D}_i, h_0}^{\text{ERM}}(S)] = \Pr[x_i \text{ unseen in } S] c_i.$$

Because  $c_i \leq 1 - \frac{1}{2H}$ , if agent  $i$  contributes only  $m_i - 1$  samples then

$$\Pr[x_i \text{ unseen in } S] c_i = (1 - c_i)^{m_i - 1} c_i \leq \frac{\varepsilon}{2H(1 - c_i)} \leq \varepsilon.$$

Hence  $\mathbf{m} - \mathbf{1}$  is also feasible for the profile  $\mathcal{D}$ . More generally, any Hamming neighbour  $\mathbf{m}'$  of  $\mathbf{m}$  is feasible for some modified profile  $(\mathcal{D}'_i, \mathcal{D}_{-i})$ , so the approximation algorithm remains oblivious throughout the specified hyper-rectangle.  $\square$