
Mechanisms that Incentivize Data Sharing in Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Federated learning is typically considered a beneficial technology which allows
2 multiple agents to collaborate with each other, improve the accuracy of their models,
3 and solve problems which are otherwise too data-intensive / expensive to be solved
4 individually. However, under the expectation that other agents will share their
5 data, rational agents may be tempted to engage in detrimental behavior such as
6 *free-riding* where they contribute no data but still enjoy an improved model. In
7 this work, we propose a framework to analyze the behavior of such rational data
8 generators. We first show how a naive scheme leads to catastrophic levels of
9 free-riding where the benefits of data sharing are completely eroded. Then, using
10 ideas from contract theory, we introduce *accuracy shaping* based mechanisms to
11 maximize the amount of data generated by each agent. These provably prevent
12 free-riding without needing any payment mechanism.

13 1 Introduction

14 Data is a *non-rivalrous good*—once produced, it can be repeatedly used multiple times without
15 exhaustion. Thus, multiple firms can simultaneously use the data produced by any individual
16 firm, increasing societal utility/welfare [21]. To promote such multiple usage, data portability
17 requirements have been widely legislated, e.g., the GDPR in the EU, CCPA in California, etc [35].
18 As a consequence, services are required to enable a user to download any personal data collected and
19 potentially re-upload it to a different service. These desiderata form a solid economic and legal basis
20 for federated learning—a new paradigm in machine learning wherein multiple data-generating agents
21 collaborate with each other to train a model on their *combined* data so that all the agents end up with
22 a better model than they would have obtained on their own [22]. Such collaborative data sharing is
23 already common in genomics research [51], internet advertisement targeting [18], and is also gaining
24 traction between networks of hospitals [see, e.g., 45, 52, 41, 16].

25 It is clear that once a certain amount of data has been produced, privacy issues aside, societal welfare
26 is maximized by allowing free access to the data thereby making it a public good. However, under
27 such an expectation, a rational agent may be tempted to *free-ride*, i.e., consume the benefits of the
28 data production by others without contributing any data themselves. This may lead to a collapse in
29 the data generation with everyone wanting to free-ride. Such a problem inevitably arises with any
30 public good [5]. Further, even if no agent actually free-rides and everyone intends to contribute data
31 out of altruism, the mere perception that others may be free-riding reduces pro-social behavior and
32 willingness to contribute [11]. Thus, the long-term success of federated learning in particular and data
33 portability in general critically require overcoming free-riding. This motivates our main question:

34 *How do we design a system which incentivizes rational agents to contribute their*
35 *fair share of data, thereby maximizing the value of the resulting model and improv-*
36 *ing collaboration?*

37 **Contribution and summary of results.**

38 • We formulate a principal-agent model [31] where each agent has a cost associated with
39 generating a data point and wants to improve the value of a model while minimizing said
40 costs (Sec. 2). Our formulation borrows ideas from contract theory while introducing new
41 concepts that are specific to the federated learning setting.
42 • Using this framework we show how giving unconditional benefit of the combined data to all
43 agents (as is standard in federated learning) leads to catastrophic free-riding where almost
44 none of the agents contribute any data (Sec. 3) at their optimal responses.
45 • Accordingly, we propose to tune the value of the model received by an agent to their
46 contribution. In the full-information setting when the agent’s cost of data generation is
47 known, we derive an optimal mechanism which overcomes free-riding and leads to maximal
48 collaboration and data generation (Sec. 4).
49 • Finally, if the costs of an agent are unknown, we show (in App. D) how to design truth-
50 revealing value curves at some cost to the principal (i.e., information rent) to incentivize the
51 agents to report their true costs.
52 Our framework can capture free-riding and the need for collaboration when faced with challenging
53 learning problems. The latter is novel to our framework—we show that if the learning task is too
54 challenging, then it is not economically viable for any single agent to tackle the problem. However,
55 using incentivizing data-sharing mechanisms, it may be possible to share the costs among participants
56 and solve it collaboratively.

57 2 Modeling an Individual Agent

58 We begin with modelling the learning task and objective for an individual agent. We then provide
59 a characterization of the optimal data contribution for each single agent without participating in a
60 federated learning scheme.

61 2.1 Value of data

62 There are n agents all of whom want to solve a *common* learning problem. This is often true in
63 federated learning since coalitions form around solving some particular task. Concretely, we assume
64 that all agents want to maximize a value function, $v(\mathcal{D}) : 2^{\mathcal{D}} \rightarrow [0, 1]$, for a dataset \mathcal{D} . For simplicity,
65 we assume that every datapoint is *exchangeable* i.e. every datapoint has the same value as any
66 other datapoint. While this is a strong assumption, it holds true if the data is generated by manually
67 labelling a subset of an already public unlabelled dataset, as is common in machine learning; e.g.,
68 Cifar [30] and ImageNet [43]. This assumption is arguably also valid in our autonomous driving
69 example where each data point involves a random path taken under random external conditions. With
70 this, we can simplify the value function $v(\cdot)$ to depend only on the *size* of the dataset $m = |\mathcal{D}|$. For
71 convenience, we will treat dataset sizes as a continuous real. Thus, every agent wants to maximize

$$v(m) : \mathbb{R}_{\geq 0} \rightarrow [0, 1] = \max(0, b(m)), \text{ where } b \text{ is continuous, non-decreasing and concave. } (1)$$

72 We also assume w.l.o.g that $v(0) = 0$ and $\lim_{m \rightarrow \infty} v(m) > 0$. Perhaps the most natural way of
73 defining the value of data is via the test accuracy obtained by training a model on the data. For
74 example, each accurate product recommendation may lead to a sale or correct digital ad placement
75 may lead to a click and hence ad revenue. This is also true if each error represents costly consequences.
76 Each error by a medical diagnostic model, a loan application evaluation model, or an autonomous
77 driving model may lead to significant suffering. In all of these cases, the value of data comes by
78 directly improving generalization and guaranteeing test accuracy.

79 2.2 Agent’s objective and optimal solution

80 Each agent i has a marginal fixed cost $c_i > 0$ for producing a data point. Their cost for producing a
81 dataset \mathcal{D} with m number of data points is then:

$$cost_i(m) = c_i m. (2)$$

82 When manually labelling a dataset or when training an autonomous-driving model, this cost c_i may
83 represent the time spent by researchers/employees or an amount paid to crowd-sourced workers. The
84 cost $c_i m$ may also represent the risk associated with privacy loss for the agent for revealing m of
85 their data points. By incurring this cost, they can obtain a model with value $v(m)$. Thus, the net
86 utility of an agent is improve value for the least cost; i.e., to maximize

$$u_i(m) = v(m) - c_i m. (3)$$

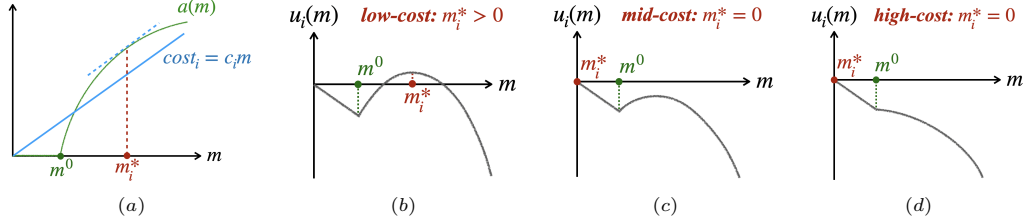


Figure 1: Illustration of the optimal amount of data for a single agent. (a): Value and cost versus the dataset size.(b)-(d): Utility function of a low/mid/high-cost agent versus the dataset size. Optimal amount for low-cost agent is positive but zero for mid and high cost agents.

87 **Theorem I** (Optimal individual generation). Consider an individual agent i with marginal cost per
 88 data point c_i and value function v satisfying (1) working on their own. Then, the optimal amount of
 89 data m_i^* is:
$$m_i^* = \begin{cases} 0 & \text{if } \max_{m_i \geq 0} u_i(m_i) \leq 0; \\ \alpha_i^*, \text{ such that } b'(\alpha_i^*) = c_i & \text{otherwise.} \end{cases} \quad (4)$$

90 Further, for agents i, j with costs $c_i \leq c_j$, their utility satisfies $u_i(m_i^*) \geq u_j(m_j^*)$ and $m_i^* \geq m_j^*$.

91 As Figure 1 shows, if the learning problem is too hard (m^0 is large) or
 92 if the marginal cost c_i is too high, the problem becomes infeasible for
 93 an individual agent to solve with $m_i^* = 0$. Such cases are especially
 94 important in federated learning where we want to enable agents to
 95 solve problems together which they cannot on their own. In other
 96 cases, the agent collects $m_i^* > 0$ data points.

97 We can simulate the value function arising from the generalization
 98 guarantees of an ERM problem ([38, Theorem 11.8]) with k mea-
 99 suring different difficulty levels (higher k means that the learning
 100 problem is harder). Figure 2.2 shows the optimal data contribution
 101 m_i^* versus the marginal cost c_i for different number of total agents on
 102 a log-log scale in such a setting. We see that the optimal contribution
 103 decreases with cost as $m_i^* \propto c_i^{-2/3}$, matching the theory. The vertical
 104 lines indicate the cutoff for minimum viability—beyond this, the cost is too high for the problem
 105 to be solvable by an individual. This minimum viability cost is smaller for more harder problems
 106 (larger k), but the optimum contribution increases with increasing k once this threshold is crossed.

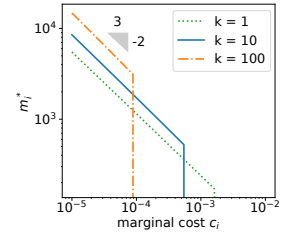


Figure 2: Optimal individual data contribution m_i^* versus the marginal cost c_i for different number of total agents.

107 3 Modeling Multiple Agents and Catastrophic Free-Riding

108 In this section we will study how agents behave when collaborating with each other as they do in
 109 federated learning. For this, we use a principal-*multi*-agent framework where the server who sets up
 110 the federated learning server is the principal.

111 3.1 Interaction between agents and server

112 The interaction between the federated learning server and the agents is formalized by a mechanism

$$\mathcal{M}(\mathbf{m}) : \mathbb{R}_{\geq 0}^n \rightarrow [0, 1]^n, \text{ which maps agents' contributions to values.} \quad (5)$$

113 We assume that each agent i generates and transmits m_i data points to the server. Based on these
 114 contributions, the mechanism assigns models to the clients with differing valuations ; i.e., if agent
 115 i contributes m_i data points it receives a model with value $v_i \in [0, 1]$, where $\mathcal{M}(m_1, \dots, m_n) =$
 116 (v_1, \dots, v_n) . The interaction proceeds in three steps: (i) first the server publishes a mechanism \mathcal{M} ,
 117 then (ii) each agent generates and transmits some data m_i to the server, and finally (iii) the server
 118 returns a trained model to each agent following the mechanism. Note that the agents decide how
 119 much data to generate adaptively *after* knowing the mechanism \mathcal{M} . However, they do not have any
 120 bargaining power—they cannot re-negotiate the mechanism—but can only decide if they join or not.
 121 We also disallow monetary compensation or exchanges between the parties since implementing them
 122 adds additional complexity. The only guarantee is that the server truthfully executes the protocol \mathcal{M} .

123 Given that the server necessarily needs to follow through on the mechanism, we need to make sure
 124 the mechanism is implementable.

125 **Definition A (Feasible mechanism).** A mechanism which returns value $[\mathcal{M}(\mathbf{m})]_i$ to agent i is said
 126 to be feasible if for any $i \in [n]$ and any $\mathbf{m} \in \mathbb{R}_{\geq 0}^n$, it satisfies $[\mathcal{M}(\mathbf{m})]_i \leq v(\sum_j m_j)$.

127 This is because we can pool together all the agent contributions \mathbf{m} and train a model to value
 128 $v(\sum_j m_j)$. Since $v(\cdot)$ is monotone, this is an upper bound on the value which can be obtained.
 129 However, it is always possible to use a subset of this data, or degrade the model in a controlled way
 130 using noisy perturbations. Thus, this captures mechanisms which are implementable in practice.

131 Faced with a potential feasible mechanism \mathcal{M} , an agent has to decide whether to join or simply train
 132 on their own. A mechanism which offers an especially bad value would discourage an agent and they
 133 would likely leave the server and train on their own. We will formalize this next.

134 **Definition B (Individual rationality (IR)).** Given data contributions \mathbf{m} by the n agents with costs
 135 \mathbf{c} , the mechanism provides a model with value $[\mathcal{M}(\mathbf{m})]_i$ to agent i . Such a mechanism \mathcal{M} is said to
 136 satisfy IR if for any agent $i \in [n]$ and any contribution \mathbf{m} ,

$$[\mathcal{M}(\mathbf{m})]_i - c_i m_i \geq v(m_i) - c_i m_i. \quad (6)$$

137 A mechanism which satisfies individual rationality guarantees that for any agent the value of the
 138 model received (and hence their utility) will be no worse than if they trained on their own. Since
 139 IR guarantees that all rational agents will participate in our mechanism, and participation is key to
 140 success of any platform, we will restrict our focus henceforth to mechanisms which satisfy IR.

141 Given any mechanism \mathcal{M} , we would like to argue about how rational agents would respond and how
 142 much data they would contribute. For this, we use the notion of an equilibrium.

143 **Theorem II (Existence of pure equilibrium).** Consider a feasible mechanism \mathcal{M} which can be
 144 expressed as:

$$[\mathcal{M}(m_i; \mathbf{m}_{-i})]_i = \max(0, \nu_i(m_i; \mathbf{m}_{-i})),$$

145 for a function $\nu_i(m_i; \mathbf{m}_{-i})$ which is continuous in \mathbf{m} and concave in m_i . For any such \mathcal{M} , there
 146 exists a pure Nash equilibrium in data contributions $\mathbf{m}^{eq}(\mathcal{M})$ which for any agent i satisfies,

$$[\mathcal{M}(\mathbf{m}^{eq}(\mathcal{M}))]_i - c_i m_i^{\mathcal{M}} \geq [\mathcal{M}(m_i, \mathbf{m}^{eq}(\mathcal{M})_{-i})]_i - c_i m_i, \text{ for all } m_i \geq 0. \quad (7)$$

147 Thus, under reasonable conditions on the mechanism \mathcal{M} which are satisfied for all the mechanisms
 148 we consider, an equilibrium always exists such that no agent can improve their utility by unilaterally
 149 changing their contribution. If all players are rational, then such an equilibrium point is a natural
 150 attractor with all the agents gravitating towards such contributions. Thus, it is reasonable to use the
 151 data contributions at this equilibrium to evaluate and compare different mechanisms.

152 Note that the mechanism is not concave because of the presence of a $\max(0, \cdot)$, and the resulting
 153 utilities of the agents are not even quasi-concave. Despite this, our proof uses the specific properties
 154 of our setting to prove existence. Our techniques may be more broadly applicable to study non-
 155 concavities arising from “minimum viability”.

156 3.2 Free-riding in the standard federated setting

157 We now examine the behavior of rational agents in the standard federated learning. Returning a
 158 model trained on the combined dataset to everyone corresponds to the mechanism

$$\mathcal{M}(\mathbf{m}) = \left(v(\sum_j m_j), \forall i \in [n] \right). \quad (8)$$

159 Clearly, this mechanism is feasible (Def. A) and also satisfies individual rationality (Def. B) since
 160 the value function $v(\cdot)$ is non-decreasing and $\sum_j m_j \geq m_i$ for any $i \in [n]$. In fact, given a data
 161 contribution \mathbf{m} , this mechanism may maximize the utility for all agents. This observation may at first
 162 seem like a strong argument in favor of this standard scheme. However, recall that the agents choose
 163 their contribution \mathbf{m} after the server publishes the mechanism \mathcal{M} . Thus, we need to first analyze
 164 how much data rational agents would contribute.

165 **Theorem III (Catastrophic free-riding).** Consider n agents with costs $\{c_i\}$ with a unique least cost
 166 agent $c_{\min} = \min_i c_i$. Let $\{m_i^*\}$ be the equilibrium contributions of agents when alone. The standard
 167 federated learning mechanism corresponding to $[\mathcal{M}(\mathbf{m})]_i = v(\sum_j m_j)$ for all clients i is feasible
 168 and IR, and has an unique equilibrium. At this equilibrium, only the lowest cost agent contributes:

$$m_i^{eq} = \begin{cases} m_i^* & \text{if } c_i = c_{\min} \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

169 The agent with the least cost $c_{\min} = \min_i c_i$ would have collected m_{\min}^* amount of data on their own.
 170 For any other agent i , $c_i \geq c_{\min}$ and so $m_i^* \leq m_{\min}^*$. Thus, agent i would already have access to
 171 data sufficient to satisfy them by the federated learning mechanism. The increase in value $v(\cdot)$ for
 172 collecting an additional data point beyond this is less than the marginal cost c_i incurred. This results
 173 in catastrophic free-riding, with only a single agent collecting data.

174 **Remark 1** (Collapse of collaboration). *Consider the case where $m_i^* = 0$ for all agents i , either*
 175 *because the learning problem is too hard or because the cost of data collection is too high for any*
 176 *individual agent. Theorem III implies that no data will be collected even with collaboration. Thus, if*
 177 *a problem is too costly to solve by an individual, it will remain insurmountable via standard federated*
 178 *learning. This is because everyone rationally assumes that everyone else will free-ride, defeating the*
 179 *main motivations of federated learning.*

180 4 Value Shaping under Verifiable Costs

181 How do we design mechanisms which prevent free-riding? In this section we will study this question
 182 assuming everyone (the server and the agents) know the costs \mathbf{c} involved in producing the data (we
 183 study the unknown costs setting in Section D), or that the costs can be verified; i.e., the agent cannot
 184 incur cost c and report a different cost \tilde{c} . This is justifiable in some cases—the cost of labelling a data
 185 point by a crowd-worker can be estimated by all parties. We formalize our goal of data maximization
 186 and give a simple optimal mechanism for it.

187 4.1 Value shaping mechanism

188 A mechanism \mathcal{M} is data-maximizing given costs \mathbf{c} if
 189 it maximizes the data collected at equilibrium.

190 **Definition C (Data Maximization).** *Suppose that*
 191 *given a mechanism \mathcal{M} , let $\mathbf{m}^{eq}(\mathcal{M})$ correspond to the*
 192 *amount of data generated by the agents at equilibrium.*
 193 *$\hat{\mathcal{M}}$ is data-maximizing if it maximizes the amount of*
 194 *data collected at equilibrium*

$$\hat{\mathcal{M}} \in \arg \max_{\mathcal{M}} \sum_j [m^{eq}(\mathcal{M})]_j, \quad (10)$$

195 *subject to \mathcal{M} being feasible and satisfying IR.*

196 **Mechanism description.** If we give Δm_i free data
 197 to agent i , then at equilibrium they will reduce the data
 198 they generate—they will only generate $(m_i^* - \Delta m_i)$ additional data. To prevent this, our key insight
 199 is to condition the amount of extra data on their actual contribution. For a given set of costs \mathbf{c} and
 200 some small $\varepsilon > 0$, consider the following mechanism:

$$[\mathcal{M}(\mathbf{m})]_i = \begin{cases} v(m_i) & \text{for } m_i \leq m_i^* \\ v(m_i^*) + (c_i + \varepsilon)(m_i - m_i^*) & \text{for } m_i \in [m_i^*, m_i^{\max}] \\ v(\sum_j m_j) & \text{for } m_i \geq m_i^{\max}, \end{cases} \quad (11)$$

201 where m_i^{\max} is defined such that $v(m_i^{\max} + \sum_{j \neq i} m_j) = v(m_i^*) + (c_i + \varepsilon)(m_i^{\max} - m_i^*)$. We
 202 illustrate the mechanism in Figure 3. Even without any external incentivization, agent i will compute
 203 m_i^* data points. Thus, for $m_i \leq m_i^*$ (10) returns a model trained on solely their own data. After m_i^* ,
 204 however, the marginal gain in value becomes smaller than the additional cost c_i . Hence, the agent
 205 requires active incentivization here and (10) ensures that for every additional data point computed,
 206 the marginal gain in value is strictly more than the cost c_i . However, the mechanism cannot provide
 207 unlimited value either and has to remain feasible, giving us our final constraint.

208 4.2 Analysis

209 **Theorem IV** (Data maximization with known costs). *The mechanism \mathcal{M} defined by (11) is data-*
 210 *maximizing for $\varepsilon \rightarrow 0^+$. At equilibrium, a rational agent i will contribute m_i^{\max} data points where*
 211 *$m_i^{\max} \geq m_i^*$, yielding a total of $\sum_j m_j^{\max}$ data points.*

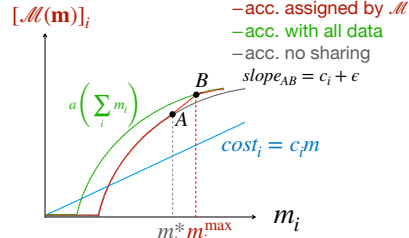


Figure 3: Illustration of value shaping. (red curve): model value returned to agent i by the mechanism; (grey curve): model value for agent i without participation; (green curve): model value if agent i receives all the data from the other agents.

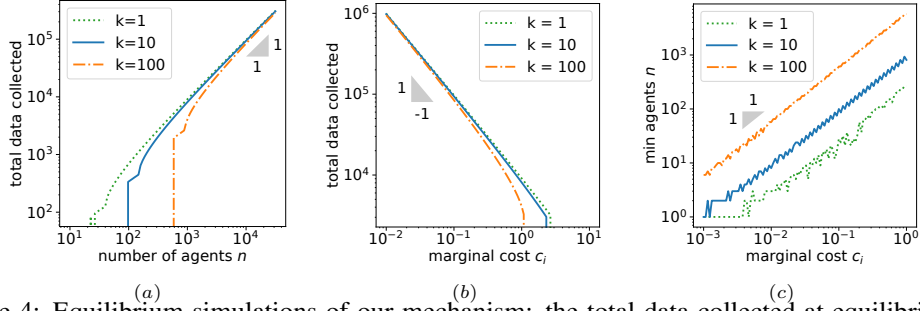


Figure 4: Equilibrium simulations of our mechanism: the total data collected at equilibrium (a) increases linearly with the number of agents n , (b) decreases as c_i^{-1} with increasing marginal cost per data point, and is relatively unaffected by the complexity k . (c) The number of agents required to cross the minimum viability threshold (i.e. smallest n for which $m_i^{\max} > 0$) increases linearly with both the marginal cost c_i and complexity k . Optimal individual contribution for all settings is $m_i^* = 0$; i.e., no data would be generated by standard federated learning.

212 Thus, to encourage collaboration we can set up a central repository of the data to which each agent
 213 is required to contribute. The cost for collecting a data point (say c) can easily be estimated and
 214 assumed to be the same for all agents. Using this estimate, we can compute a threshold. The agents
 215 don't receive any additional data for contributing up to this threshold. For each data point contributed
 216 beyond the threshold, the agents receive an increasing amount of additional data. By Theorem IV,
 217 this would prevent free-riding by the agents and ensure the best trained model reaches the consumers.

218 **Remark 2** (Deterrence). *At equilibrium, mechanism \mathcal{M} in (11) ensures all agents contribute $m_i^{\max} \geq$
 219 m_i^* ; i.e., they generate more data than they would on their own. Further, every agent receives a
 220 model trained on this combined dataset with value $v(\sum_j m_j^{\max})$. One can view our mechanism as
 221 using a deterrent which punishes free-riding, ensuring that all agents fully utilize the combined data.
 222 At equilibrium, such a deterrent is never actually invoked but just forms a credible threat.*

223 Suppose all agents have the same cost c . Theorem IV shows that the mechanism collects nm^{\max}
 224 data points in total. However, m^{\max} also depends on n . This is because with a larger pool of data
 225 contributions, the server can more strongly incentivize an individual and extract more data. There is
 226 a natural ceiling to this though—the value caps at 1. Thus, the absolute maximum data that can be
 227 extracted from an individual agent is m which satisfies $v(m^*) + c(m - m^*) = 1$. This gives us the
 228 range for the total data contributions to be $[nm^*, n(m^* + (1-v(m^*))/c)]$.

229 **Remark 3** (Collaboratively overcoming minimum viability). *When $m^* = 0$, i.e., the problem is
 230 not solvable by an individual agent, the net contribution from our mechanism nm^{\max} may still be
 231 positive. Suppose that the cost for all agents is the same c . Then, the total data collected is m^{tot}
 232 which satisfies*

$$c/n \cdot m^{\text{tot}} = v(m^{\text{tot}}).$$

233 *This implies that for sufficiently large n , the cost c is successfully shared and we obtain a positive
 234 data contribution. However, note that $m^{\text{tot}} = 0$ is also a valid solution and remains an equilibrium.
 235 If all other agents don't contribute, there is no extra data to share and so there is no incentive to
 236 compute extra data. In practice, this undesirable equilibrium is unlikely to be encountered since it has
 237 lower utility. It can also be prevented by the platform itself taking part as an agent and committing to
 238 non-zero data collection.*

239 Empirically, in Figure 4 we compute the equilibrium for value shaping assuming the value of the data
 240 is test accuracy. We assume all agents have the same cost, and observe the effect on the equilibrium
 241 data contribution as we vary the cost c and the total number of agents n . We used the following
 242 default parameters unless otherwise states: optimal value of $a_{\text{opt}} = 0.95$, marginal cost $c_i = 0.1$,
 243 participants $n = 10^4$. Under all parameter configurations of this experiment, the optimal individual
 244 contribution is $m_i^* = 0$, while the equilibrium data contributions are significantly larger as expected,
 245 validating our theory.

246 Finally, in App. B we discuss the incentive compatibility of our schemes, App. C performs additional
 247 simulations and detailed comparisons with prior work, and in Appendix D we extend our framework
 248 to the setting of unverifiable costs. Our conclusions translate to this unverifiable cost setting as well.

References

- 249
- 250 [1] Financials of openai. [https://web.archive.org/web/20220731013350/https://](https://web.archive.org/web/20220731013350/https://www.crunchbase.com/organization/openai/company_financials)
251 www.crunchbase.com/organization/openai/company_financials. Accessed: 2021-07-
252 30.
- 253 [2] Microsoft exclusively licenses openai’s groundbreaking gpt-3 text generation model. [http://](http://web.archive.org/web/20220731012339/https://www.theverge.com/2020/9/22/21451283/microsoft-openai-gpt-3-exclusive-license-ai-language-research)
254 [web.archive.org/web/20220731012339/https://www.theverge.com/2020/9/22/](http://web.archive.org/web/20220731012339/https://www.theverge.com/2020/9/22/21451283/microsoft-openai-gpt-3-exclusive-license-ai-language-research)
255 [21451283/microsoft-openai-gpt-3-exclusive-license-ai-language-research](http://web.archive.org/web/20220731012339/https://www.theverge.com/2020/9/22/21451283/microsoft-openai-gpt-3-exclusive-license-ai-language-research).
256 Accessed: 2021-07-30.
- 257 [3] Daron Acemoglu and Asu Ozdaglar. Lecture notes for course “6.207/14.15 networks”. [https://](https://economics.mit.edu/files/4711)
258 economics.mit.edu/files/4711. Accessed: 2022-01-30.
- 259 [4] Jeeyun Sophia Baik. Data privacy against innovation or against discrimination?: The case of
260 the california consumer privacy act (ccpa). *Telematics and Informatics*, 52, 2020.
- 261 [5] William J Baumol. Welfare economics and the theory of the state. In *The Encyclopedia of*
262 *Public Choice*, pages 937–940. Springer, 2004.
- 263 [6] Philippe Bich. Some fixed point theorems for discontinuous mappings. *Cahiers de la Maison*
264 *des Sciences Economiques*, 2006.
- 265 [7] Patrick Bolton and Mathias Dewatripont. *Contract Theory*. MIT Press, 2004.
- 266 [8] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir
267 Ivanov, Chloe Kiddon, Jakub Konečný, Stefano Mazzocchi, Brendan McMahan, et al. Towards
268 federated learning at scale: System design. *Proceedings of Machine Learning and Systems*, 1:
269 374–388, 2019.
- 270 [9] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
271 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
272 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 273 [10] Yong Cheng, Yang Liu, Tianjian Chen, and Qiang Yang. Federated learning for privacy-
274 preserving ai. *Communications of the ACM*, 63(12):33–36, 2020.
- 275 [11] Taehyon Choi and Peter J Robertson. Contributors and free-riders in collaborative governance: A
276 computational exploration of social motivation and its effects. *Journal of Public Administration*
277 *Research and Theory*, 29(3):394–413, 2019.
- 278 [12] Mingshu Cong, Han Yu, Xi Weng, and Siu Ming Yiu. A game-theoretic framework for incentive
279 mechanism design in federated learning. In *Federated Learning*, pages 205–222. Springer,
280 2020.
- 281 [13] Ningning Ding, Zhixuan Fang, and Jianwei Huang. Incentive mechanism design for federated
282 learning with multi-dimensional private information. In *2020 18th International Symposium on*
283 *Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks (WiOPT)*, pages 1–8.
284 IEEE, 2020.
- 285 [14] Kate Donahue and Jon Kleinberg. Model-sharing games: Analyzing federated learning under
286 voluntary participation. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
287 volume 35, pages 5303–5311, 2021.
- 288 [15] Kate Donahue and Jon Kleinberg. Optimality and stability in federated learning: A game-
289 theoretic approach. *Advances in Neural Information Processing Systems*, 34, 2021.
- 290 [16] Mona Flores, Ittai Dayan, Holger Roth, Aoxiao Zhong, Ahmed Harouni, Amilcare Gentili,
291 Anas Abidin, Andrew Liu, Anthony Costa, Bradford Wood, et al. Federated learning used for
292 predicting outcomes in SARS-COV-2 patients. *Research Square*, 2021.
- 293 [17] Yann Fraboni, Richard Vidal, and Marco Lorenzi. Free-rider attacks on model aggregation in
294 federated learning. In *International Conference on Artificial Intelligence and Statistics*, pages
295 1846–1854. PMLR, 2021.

- 296 [18] Google. Google ads data hub, 2022. URL [https://web.archive.org/web/](https://web.archive.org/web/20220423221048/https://developers.google.com/ads-data-hub/guides/intro)
 297 <https://developers.google.com/ads-data-hub/guides/intro>.
 298 Accessed on 2022.04.28.
- 299 [19] Jiyue Huang, Rania Talbi, Zilong Zhao, Sara Boucchenak, Lydia Y Chen, and Stefanie Roos.
 300 An exploratory analysis on users' contributions in federated learning. In *2020 Second IEEE*
 301 *International Conference on Trust, Privacy and Security in Intelligent Systems and Applications*
 302 *(TPS-ISA)*, pages 20–29. IEEE, 2020.
- 303 [20] Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nick Hynes, Nezihe Merve Gürel,
 304 Bo Li, Ce Zhang, Dawn Song, and Costas J Spanos. Towards efficient data valuation based on
 305 the Shapley value. In *The 22nd International Conference on Artificial Intelligence and Statistics*,
 306 pages 1167–1176. PMLR, 2019.
- 307 [21] Charles I Jones and Christopher Tonetti. Nonrivalry and the economics of data. *American*
 308 *Economic Review*, 110(9):2819–58, 2020.
- 309 [22] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Ar-
 310 jun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings,
 311 et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine*
 312 *Learning*, 14(1–2):1–210, 2021.
- 313 [23] Jiawen Kang, Zehui Xiong, Dusit Niyato, Shengli Xie, and Junshan Zhang. Incentive mechanism
 314 for reliable federated learning: A joint optimization approach to combining reputation and
 315 contract theory. *IEEE Internet of Things Journal*, 6(6):10700–10714, 2019.
- 316 [24] Jiawen Kang, Zehui Xiong, Dusit Niyato, Dongdong Ye, Dong In Kim, and Jun Zhao. Toward
 317 secure blockchain-enabled internet of vehicles: Optimizing consensus management using
 318 reputation and contract theory. *IEEE Transactions on Vehicular Technology*, 68(3):2906–2920,
 319 2019.
- 320 [25] Jiawen Kang, Zehui Xiong, Dusit Niyato, Han Yu, Ying-Chang Liang, and Dong In Kim.
 321 Incentive design for efficient federated learning in mobile networks: A contract theory approach.
 322 In *2019 IEEE VTS Asia Pacific Wireless Communications Symposium (APWCS)*, pages 1–5.
 323 IEEE, 2019.
- 324 [26] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,
 325 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
 326 models. *arXiv preprint arXiv:2001.08361*, 2020.
- 327 [27] Hyesung Kim, Jihong Park, Mehdi Bennis, and Seong-Lyun Kim. Blockchain-enabled on-device
 328 federated learning. *IEEE Communications Letters*, 24(6):1279–1283, 2019.
- 329 [28] Jakub Konečný, H Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated optimiza-
 330 tion: Distributed machine learning for on-device intelligence. *arXiv preprint arXiv:1610.02527*,
 331 2016.
- 332 [29] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh,
 333 and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv*
 334 *preprint arXiv:1610.05492*, 2016.
- 335 [30] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
 336 2009.
- 337 [31] Jean-Jacques Laffont and David Martimort. The theory of incentives. In *The Theory of*
 338 *Incentives*. Princeton University Press, 2009.
- 339 [32] Li Li, Yuxi Fan, Mike Tse, and Kuo-Yi Lin. A review of applications in federated learning.
 340 *Computers & Industrial Engineering*, 149:106854, 2020.
- 341 [33] Wei Yang Bryan Lim, Jianqiang Huang, Zehui Xiong, Jiawen Kang, Dusit Niyato, Xian-Sheng
 342 Hua, Cyril Leung, and Chunyan Miao. Towards federated learning in uav-enabled internet of
 343 vehicles: A multi-dimensional contract-matching approach. *IEEE Transactions on Intelligent*
 344 *Transportation Systems*, 22(8):5140–5154, 2021.

- 345 [34] Jierui Lin, Min Du, and Jian Liu. Free-riders in federated learning: Attacks and defenses. *arXiv*
346 *preprint arXiv:1911.12560*, 2019.
- 347 [35] James Mancini. Data portability, interoperability and digital platform competition: Oecd
348 background paper. *Interoperability and Digital Platform Competition: OECD Background*
349 *Paper (June 8, 2021)*, 2021.
- 350 [36] Eric Maskin. The existence of equilibrium. *The Review of Economic Studies*, 53(1):1–26, 1986.
- 351 [37] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
352 Communication-efficient learning of deep networks from decentralized data. In *Artificial*
353 *Intelligence and Statistics*, pages 1273–1282. PMLR, 2017.
- 354 [38] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*.
355 MIT press, 2018.
- 356 [39] Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In
357 *International Conference on Machine Learning*, pages 4615–4625. PMLR, 2019.
- 358 [40] Adam Richardson, Aris Filos-Ratsikas, and Boi Faltings. Budget-bounded incentives for
359 federated learning. In *Federated Learning*, pages 176–188. Springer, 2020.
- 360 [41] Nicola Rieke, Jonny Hancox, Wenqi Li, Fausto Milletari, Holger R Roth, Shadi Albarqouni,
361 Spyridon Bakas, Mathieu N Galtier, Bennett A Landman, Klaus Maier-Hein, et al. The future
362 of digital health with federated learning. *NPJ Digital Medicine*, 3(1):1–7, 2020.
- 363 [42] Ankit Rohatgi. Webplotdigitizer, 2017.
- 364 [43] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
365 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei.
366 ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*
367 (*IJCV*), 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- 368 [44] Yunus Sarikaya and Ozgur Ercetin. Motivating workers in federated learning: A Stackelberg
369 game perspective. *IEEE Networking Letters*, 2(1):23–27, 2019.
- 370 [45] Micah J Sheller, G Anthony Reina, Brandon Edwards, Jason Martin, and Spyridon Bakas.
371 Multi-institutional deep learning modeling without sharing patient data: A feasibility study
372 on brain tumor segmentation. In *International MICCAI Brainlesion Workshop*, pages 92–104.
373 Springer, 2018.
- 374 [46] Rachael Hwee Ling Sim, Yehong Zhang, Mun Choon Chan, and Bryan Kian Hsiang Low. Col-
375 laborative machine learning with incentive-aware model rewards. In *International Conference*
376 *on Machine Learning*, pages 8927–8936. PMLR, 2020.
- 377 [47] Stephen A Smith. *Contract Theory*. OUP Oxford, 2004.
- 378 [48] Mengmeng Tian, Yuxin Chen, Yuan Liu, Zehui Xiong, Cyril Leung, and Chunyan Miao. A con-
379 tract theory based incentive mechanism for federated learning. *arXiv preprint arXiv:2108.05568*,
380 2021.
- 381 [49] Paul Voigt and Axel Von dem Bussche. The EU general data protection regulation (gdpr). *A*
382 *Practical Guide, 1st Ed.*, Cham: Springer International Publishing, 10(3152676):10–5555,
383 2017.
- 384 [50] Tianhao Wang, Johannes Rausch, Ce Zhang, Ruoxi Jia, and Dawn Song. A principled approach
385 to data valuation for federated learning. In *Federated Learning*, pages 153–167. Springer, 2020.
- 386 [51] John N Weinstein, Eric A Collisson, Gordon B Mills, Kenna R Shaw, Brad A Ozenberger,
387 Kyle Ellrott, Ilya Shmulevich, Chris Sander, and Joshua M Stuart. The cancer genome atlas
388 pan-cancer analysis project. *Nature genetics*, 45(10):1113–1120, 2013.
- 389 [52] Y Wen, W Li, H Roth, and P Dogra. Federated learning powered by NVIDIA
390 Clara, 2019. URL [https://web.archive.org/web/20220221070237/https://](https://web.archive.org/web/20220221070237/https://developer.nvidia.com/blog/federated-learning-clara/)
391 developer.nvidia.com/blog/federated-learning-clara/. Accessed on 2022.04.28.

- 392 [53] Xinyi Xu, Lingjuan Lyu, Xingjun Ma, Chenglin Miao, Chuan Sheng Foo, and Bryan
393 Kian Hsiang Low. Gradient driven rewards to guarantee fairness in collaborative machine
394 learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- 395 [54] Yufeng Zhan, Peng Li, Zhihao Qu, Deze Zeng, and Song Guo. A learning-based incentive
396 mechanism for federated learning. *IEEE Internet of Things Journal*, 7(7):6360–6368, 2020.
- 397 [55] Ning Zhang, Qian Ma, and Xu Chen. Enabling long-term cooperation in cross-silo federated
398 learning: A repeated game perspective. *IEEE Transactions on Mobile Computing*, 2022.

399 Appendix

400 Contents

401	1 Introduction	1
402	2 Modeling an Individual Agent	2
403	2.1 Value of data	2
404	2.2 Agent’s objective and optimal solution	2
405	3 Modeling Multiple Agents and Catastrophic Free-Riding	3
406	3.1 Interaction between agents and server	3
407	3.2 Free-riding in the standard federated setting	4
408	4 Value Shaping under Verifiable Costs	5
409	4.1 Value shaping mechanism	5
410	4.2 Analysis	5
411	A Review on the Related Work and Contract Theory Background	12
412	B Incentive compatibility under Verifiable Costs	13
413	C Additional Simulations and Comparisons	13
414	C.1 Simulating collaborative training of GPT-3	13
415	C.2 Comparison with baselines	14
416	C.2.1 Many bad equilibria	14
417	C.2.2 Sensitivity to choice of p and under-performance	15
418	C.2.3 Discrimination against high-cost agents	15
419	D Data Maximization under Unverifiable Costs	17
420	D.1 Mechanism description	17
421	D.2 Analysis	18
422	E Proofs from Section 2 (Optimal Individual Contributions)	19
423	F Proofs from Section 3 (Modeling Multiple Agents and Catastrophic Free-riding)	19
424	G Proofs from Section 4 (value Shaping under Known Costs)	22
425	H Proofs from Appendix D (Data Maximization with Unverifiable Costs)	24

426 A Review on the Related Work and Contract Theory Background

427 The literature on mechanism design and federated learning is vast. We discussed the most closely
428 related work in three verticals in the main text; we include a detailed review of the broader literature
429 in this section.

430 Over the past decade, federated learning (FL) has emerged as an important paradigm in modern
431 large-scale machine learning [28, 27, 22, 32, 41]. Specifically, FL research has resulted in many
432 applications to overcome practical challenges such as data silos and data sensitivity: on one side,
433 since more training data often gives better model performance, data silos results in scarcity of
434 labeled training data and puts limit on the industrial performance; on the other side, in high-stakes
435 applications the data may contain private user information and thus the sharing of data is constrained
436 by regulations and laws [49, 4, 10, 35]. Given these challenges, FL provides a useful scheme for
437 different agents / parties to train collaboratively and leverage the benefit from other agents' data,
438 while the training data remains distributed over the agents. Such a framework has been shown to
439 be able to bring improved model performance to all the participants. Indeed, many prior works
440 have been devoted to develop more scalable and communication-efficient distributed optimization
441 algorithms for FL [29, 37, 8].

442 However, one cannot ignore an important aspect in the standard FL scheme, which is the incentives
443 aspect. The standard FL scheme may incentivize strategic agents to contribute less data in order
444 to minimize their data collection cost and maximize the gain from participating in the federated
445 learning mechanism. Although the participation and contribution of each agent is often legislated
446 by certain protocols, such free-riding behavior has been notoriously hard to regulate and prevent
447 in practice [17, 19]. Recently, a few works have started to explore such free-riding behavior in FL,
448 with various incentive models proposed [40, 44, 34, 17, 13, 55]. However, the majority this work
449 has focused on a taxonomy of free-rider attacks or the detection of attacks under the existing FL
450 scheme, instead of proposing mechanisms that incentivize maximal data contribution. In this work,
451 we strive for a mechanism for information sharing under the standard federated learning setting such
452 that rational agents are incentivized to contribute their maximal amount of data.

453 In this work, we focus on the free-riding behavior of FL agents in terms of data collection. In FL,
454 the data collection happen on the agents' side before they join the mechanism for training models.
455 Therefore, the cost of collecting data is often *private* information to each agent. Such an information
456 asymmetry brings difficulty to prevent free-riding, because the agents might simply report fake costs.
457 This brings the need to design *incentive* mechanisms for FL, under which the agents are incentivized
458 to behave truthfully, which is also guaranteed to lead to the best utility.

459 Indeed, designing incentive mechanisms under private costs is not new, and has been a main focus
460 of the contract theory literature Smith [47], Laffont and Martimort [31], Bolton and Dewatripont
461 [7]. Moreover, the existence of a central server (a "principal") in FL brings further convenience to
462 apply a principle-agent model. An emerging line of recent works have been exploring the application
463 of contract theory for federated learning [24, 25, 23, 33, 48, 12, 54]. In particular, Tian et al. [48]
464 proposed a contract-based aggregator under a multi-dimensional contract model over two possible
465 types of agents and showed improved model generalization accuracy under that contract. However,
466 their mechanism focused on eliciting the private type information instead of maximizing the data
467 contribution. To the best of our knowledge, our work is a first step to use contract theory for *data*
468 *maximization* in federated learning. Further, prior work has focused on how to design payments to
469 agents, rather than the value-shaping problem that we focus on here.

470 This work is related to the active line of research on mechanism design for collaborative machine
471 learning, which involves multiple parties each with their own data, jointly training a model or making
472 predictions in a common learning task [46, 53]. In collaborative machine learning, a major focus has
473 been the design of model rewards (i.e., data valuation) in order to ensure certain fairness or accuracy
474 objectives. Towards that goal, there has been model rewards proposed based on notions from the
475 cooperative game theory literature such as the Shapley value [20, 50]. However, the guarantees of
476 these model rewards depend on the assumption that the agents are already willing to contribute the
477 data they have. In this work, we study a different incentivization task for data maximization.

478 More broadly, apart from data maximization, there are other objectives which are of interest in
479 federated learning, such as fairness and welfare objectives, that have been under active study [15, 14,
480 39]. We defer a thorough analysis of the tradeoffs among various objectives to future research.

481 B Incentive compatibility under Verifiable Costs

482 One of our motivating reasons for preventing free-riding was to ensure that none of the participating
483 agents feel taken advantage of. That is, we wish to satisfy some notion of fairness. However, there
484 may potentially be new sources of unfairness in (11). In particular, consider two agents, $i, j \in [n]$,
485 with different costs: if $c_i \leq c_j$ then $m_i^* \geq m_j^*$. Here, an agent i with smaller cost c_i faces two
486 disadvantages under mechanism (11): (i) they have a larger threshold amount of data m_i^* they have
487 to contribute before receiving any benefit, and (ii) they receive a smaller increase in value ($c_i + \varepsilon$) for
488 each additional data point computed.

489 If the cost for generating each data point is inherently fixed (such as the cost of driving a vehicle) this
490 is arguably not an issue. However, in many other settings an agent may innovate and develop new
491 methods to reduce their cost of collecting a data point. In fact, the business model of large internet
492 advertising providers is based on systems which can cheaply capture consumer data in order to show
493 them better advertisements. Would our data-sharing mechanism (10) disincentivize agents from such
494 innovations? We show this is in fact not true.

495 **Theorem V** (Incentive compatibility). *Under given costs \mathbf{c} , consider our optimal mechanism (11)*
496 *with equilibrium contributions \mathbf{m}^{\max} , and agents working individually with equilibrium contributions*
497 *of \mathbf{m}^* . The utility of the every agent i remains unchanged:*

$$v(\sum_j m_j^{\max}) - c_i m_i^{\max} = v(m_i^*) - c_i m_i^*.$$

498 Thus, our mechanism does not induce any distortions in the incentive structure. Further, recall by
499 Theorem I, the utility $u_i(m_i^*) \geq u_j(m_j^*)$ if $c_i \leq c_j$. This implies that users with smaller costs
500 continue to receive a higher utility, encouraging them to innovate and reduce the costs; i.e., our
501 mechanism is incentive compatible. Of course this is assuming that the costs incurred by an agent is
502 verifiable. They cannot lie about the true cost, but may be able to choose between different collection
503 strategies.

504 **Remark 4** (Distribution of surplus). *One may ask where the additional surplus which is generated*
505 *by agents collaborating has disappeared, since the agents receive none of it. Our mechanism utilizes*
506 *this surplus in order to extract additional data, $m_i^{\max} - m_i^*$, from the agents. Thus, all the additional*
507 *surplus goes into improving the value of the model and hence to the end consumers of the model.*

508 C Additional Simulations and Comparisons

509 C.1 Simulating collaborative training of GPT-3

510 We illustrate through a pedagogical example how one may use our theory in practice. We first extract
511 the data of loss values obtained by training GPT-like language models on dataset of varying sizes
512 from [26, Fig. 1] using WebPlotDigitizer [42]. Then, we fit a simple linear regression model in the
513 log-log space to obtain a close fit $\text{loss}(m) = \left(\frac{5.4 \times 10^{13}}{m}\right)^{0.95}$. We use this to model how the loss
514 would decrease as data increases. Then, we can define accuracy as $(1 - \text{loss}(m))$. To define the value
515 function, we need to assign a dollar value to a perfectly trained model. Microsoft reportedly paid 1
516 billion (10^9) 2019 dollars for a license to GPT-3 [2] and hence this forms an estimate of the value of
517 a fully trained model to one company. Thus we have

$$v(m) = 10^9 \left(1 - \left(\frac{5.4 \times 10^{13}}{m}\right)^{0.95}\right).$$

518 GPT-3 was trained on 500 billion tokens [9]. OpenAI has raised an estimated 1 billion USD [1].
519 Suppose that accounting for salaries of all personnel involved etc., we allocate half of this money as
520 being spent on training GPT-3, this gives an estimate of the marginal cost per datapoint as $c_i = 10^{-3}$
521 USD. With these numbers, we see that we need at least 1000 companies (each for whom the trained
522 model is worth 1 billion USD) collaborating together to make it feasible. Note that with our estimated
523 costs and benefits, it seems like OpenAI is at a loss. This is true—it would likely need to license GPT-3
524 (or its successors) to many more companies before it breaks even. Finally, we emphasize that this was
525 more of a pedagogical exercise and not an actual prediction about outcomes. The stylized framework
526 here is meant to provide qualitative insights about the incentives at play.

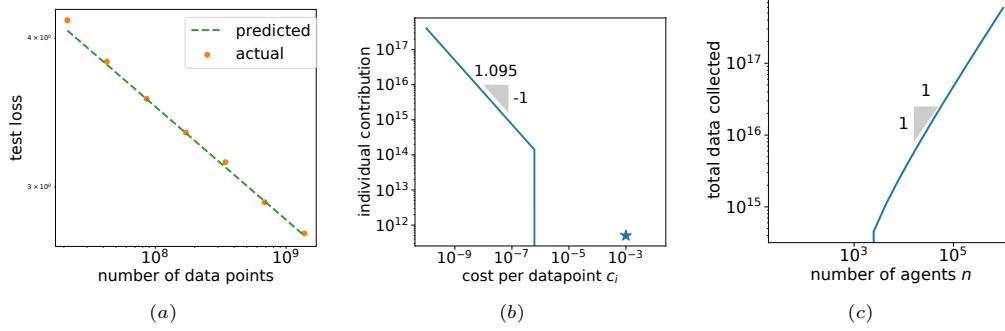


Figure 5: Simulating training of GPT-3. In (a), we take real world data (yellow dots) of how loss scales with data size and show that $(\frac{5.4 \times 10^{13}}{m})^{0.95}$ is a good fit (dashed green). We combine this with an estimate of the value of model and compute the optimal individual data contribution in (b) for different marginal data costs. The blue star shows the estimated marginal cost and total data collected by open-AI for training GPT-3. Finally, (c) shows the data collected by our data-maximizing mechanism with the estimated utility function. It shows that we need at least 1000 collaborative agents to train the GPT-3 model. The data collected initially grows super-linearly in n , but asymptotically becomes linear.

527 C.2 Comparison with baselines

528 We compare with some alternative fairness inspired deterrence mechanisms which do not rely on our
 529 contract theory framework. These can be summarized as if you do not contribute as much as your
 530 peers, then you may be punished. Suppose the agents submit \mathbf{m} number of data points. Then, the
 531 server chooses a feasible mechanism which returns a model of value less than $v(\sum_j m_j)$ to the client
 532 i . We consider the following mechanisms:

- 533 • **Proportionate data (PD)**. Penalize agent i if they submit less number of data points than
 534 their peers as

$$[\mathcal{M}(\mathbf{m})]_i = \left(\frac{m_i}{\max_j m_j} \right)^p v(\sum_j m_j) \quad \text{for } p \in [0, 1] \quad (12)$$

- 535 • **Proportionate value (PV)**. Penalize agent i for contributing less value than their peers:

$$[\mathcal{M}(\mathbf{m})]_i = \left(\frac{v(m_i)}{\max_j v(m_j)} \right)^p v(\sum_j m_j) \quad \text{for } p \in [0, 1] \quad (13)$$

- 536 • **Proportionate Shapley (PS) [46]**. Shapely values has a long history of being used a fair
 537 contribution measure. Thus, we can compute the Shapely value for each agent's contribution
 538 $\phi_i(\mathbf{m})$ and penalize as

$$[\mathcal{M}(\mathbf{m})]_i = \left(\frac{\phi_i(\mathbf{m})}{\max_j \phi_j(\mathbf{m})} \right)^p v(\sum_j m_j) \quad \text{for } p \in [0, 1]. \quad (14)$$

539 If all other agents are contributing m datapoints, the shapely value for agent i for contributing
 540 m_i datapoints simplifies as

$$\phi_i(m_i) = \frac{1}{n} \sum_{k \in [n]} v(km + m_i) - v(km).$$

541 In all cases, $p = 0$ returns to the standard federated learning scheme which, as we saw in Section 3,
 542 has catastrophic free-riding. These measures have numerous shortcomings which we explore in
 543 sequence.

544 C.2.1 Many bad equilibria

545 Because the mechanism only penalizes on relative performance among the different agents, there are
 546 multiple stable equilibrium. Consider n identical agents with same cost c here in a *low-cost* setting
 547 with positive optimal individual contributions $m^* > 0$. Our conclusions also hold in more general
 548 settings, but we focus on this setting for simplicity.

549 **Theorem VI.** Consider mechanisms (12)–(14) with n identical agents with marginal cost c . Let
550 $m^* > 0$ be the equilibrium individual contribution and m^{\max} is the equilibrium contribution by our
551 optimal mechanism. Then, there is a set of data contributions \mathcal{S} such that all agents contributing
552 $m \in \mathcal{S}$ constitutes an equilibrium. Further, $\frac{m^*}{n} \in \mathcal{S}$ is an equilibrium with the maximum utility.

553 Note that this implies that every agent only contributing $\frac{m^*}{n}$ i.e. n times lesser than they would on
554 their own is also an equilibrium. With this only m^* datapoints would be collected by the server.
555 Further, this equilibrium corresponds to the maximum utility and so it is possible that all agents
556 will converge to this. In contrast, our optimal scheme has an unique equilibrium corresponding to
557 maximum data contribution.

558 *Proof.* Consider the generic mechanism $[\mathcal{M}(\mathbf{m})]_i = \left(\frac{\psi_i(\mathbf{m})}{\max_j \psi_j(\mathbf{m})} \right)^p v(\sum_j m_j)$ for any positive,
559 continuous, non-decreasing contribution measure ψ . Suppose that all agents submit m data-points.
560 This is an equilibrium if the following condition is satisfied:

$$-p \frac{\psi'_i(m)}{\psi_i(m)} v(nm) \geq v'(nm) - c \geq 0.$$

561 The right hand side is satisfied for $m = m^*/n$. Also note that ψ' is positive meaning the left hand
562 hand side is negative. This implies that as long as ψ' and ψ are continuous around m , there exist a set
563 of solutions all of which satisfy the above condition. Contributing $m = m^*/n$ is utility maximizing
564 for all the agents in general, and so corresponds to the maximum utility equilibrium as well. \square

565 One way out of this may be for the server to take part in the process as an agent and also contribute
566 data. This way by increasing its contribution, the server can force other agents choose an equilibrium
567 corresponding to a large equilibrium.

568 C.2.2 Sensitivity to choice of p and under-performance

569 Consider n agents each have identical *high-costs* c with optimal individual contribution is $m^* = 0$.
570 We use the same experimental setup as in Figure 4 with $a_{opt} = 0.95$, marginal cost $c = 0.1$, and
571 $n = 10^4$. In Figure 6, we numerically compute the equilibrium corresponding to the *maximum data*
572 *contribution* for each of alternative mechanisms, assuming the server may be able to intervene and
573 direct the agents towards the most beneficial equilibrium.

574 When we use $p = 0$ all the alternative schemes recover the standard FedAvg scheme. As we saw in
575 Sec. 3, this implies for $p = 0$ there is catastrophic free-riding and hence the equilibrium contribution
576 is 0. However, if we choose a value of p too large, it is possible that the mechanism no longer
577 satisfies *individual rationality*. This means that at equilibrium the agents drop out and effectively
578 contribute 0 data points again. The proportional value scheme significantly under performs, whereas
579 the proportional Shapley scheme suggested by Sim et al. [46] and the much simpler proportional data
580 scheme both perform reasonably well. However, note that even for the best value of p , both these
581 schemes do not match our much simpler data-maximizing value shaping mechanism.

582 C.2.3 Discrimination against high-cost agents

583 Consider $n = 10^4$ agents with $a_{opt} = 0.95$ and which are one of two types: either they have a low
584 cost of 1×10^{-4} , or they have a comparable but slightly higher cost of 2×10^{-4} . A fraction (say
585 $p_i \in [0, 1]$) are of the low cost and the rest have high cost. We assume the server is aware of (or
586 can verify) the cost of each of the agent. In Figure 7, we numerically compute the equilibrium data
587 contribution of the high and low cost agents for the data maximizing value shaping mechanism, and
588 the proportional data mechanism. For the latter PD, we compute the *maximum data contribution*
589 *equilibrium* as well as use the *optimal* p to compare value shaping against the best possible version of
590 the alternative mechanism.

591 We observe that the total amount of data collected by value shaping is much more (up to $50\times$) than
592 PD, especially when a large fraction of agents have high cost. The high cost agents continue to
593 contribute significant amount of data when using the value shaping mechanism. Recall that with
594 these contribution levels, they receive full value of the combined data from the mechanism. However,
595 very starkly, the high cost agent chooses to opt-out and contributed no data with the PD mechanism.
596 This means that with PD, the high cost agents *receive zero value*. This is a direct result of PD being

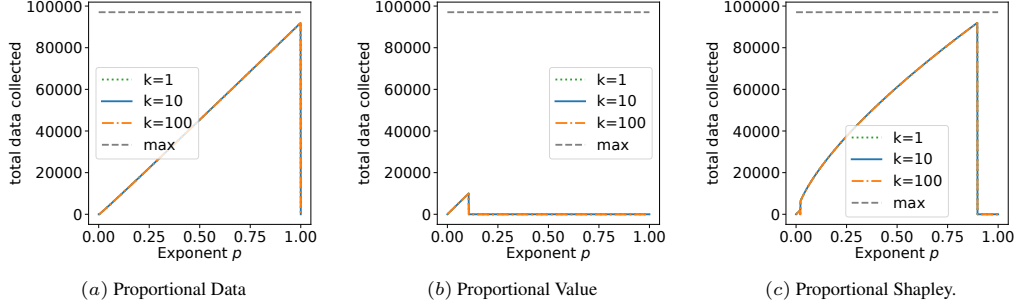


Figure 6: Total data collected as we vary the exponent p in the different mechanisms. The black dashed line represents the data collected by our data maximizing value-shaping mechanism outlined in Sec. 4. Using a small p is an insufficient deterrent and so each agent tends to free-ride yielding low overall data collection. Using too high value of p (say $p = 1$) makes the deterrent so strong that the agents rationally chooses to drop out and contribute 0.

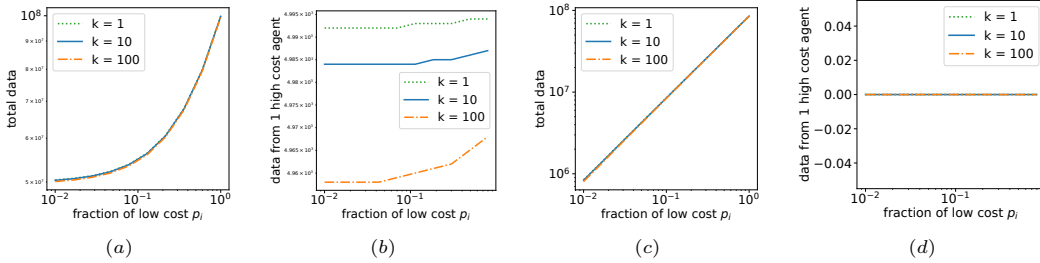


Figure 7: Data contributions in a mixture of low and high cost agents as we vary the fraction of low cost agents. Left shows the data maximizing value shaping mechanism and the (a)total data collected, and (b) data contributed by the high cost agents. Right (c) and (d) show the same with proportional data mechanism. Total data collected in (a) with value shaping is $50\times$ more than in (c) with proportional data. Further, the high cost agents continue to contribute and receive full value in (b) with value shaping, but they opt-out and receive zero value in (d) with proportional data. Thus, not accounting for costs may end up discriminating against high cost agents.

597 agnostic to the cost—it penalizes an agent for contributing lower amount of data without taking into
 598 consideration the difference in costs involved in such a process. In the face of such deterrent, the high
 599 cost agent rationally chooses to out-out. This is despite the costs differing by a mere $2\times$.

600 Our framework has thus thrown to light an important insight—‘fair’ mechanisms which do not account
 601 for difference in costs of the agents will end up being extremely unfair. Further, the high cost agents
 602 will rationally chose to drop out and derive no value from the system. Consider the setting where
 603 agents from resource-rich and resource-poor locations collaborate. The latter agents will have a
 604 higher data collection costs due to systemic barriers and will be unfairly penalized in a proportionate
 605 fairness approach.

606 D Data Maximization under Unverifiable Costs

607 Until now, we assumed that the cost of all agents is known to everyone involved, or is atleast verifiable.
 608 In some settings where the costs are universal and outside the control of the agent, this assumption
 609 may be justified. However, in numerous other cases, the exact process of the data generation may
 610 be a trade secret and so there is uncertainty about the cost incurred by an agent. In this section, we
 611 examine how to incorporate such uncertainties into our mechanism.

612 We focus on the simplest version of this uncertainty. Suppose that we know that the cost of each
 613 agent can either be low (\underline{c}) or high (\bar{c}). Further, suppose we have some prior knowledge where agent
 614 i has low cost \underline{c} with probability p_i and \bar{c} with $(1 - p_i)$. Note that there is an inherent *information*
 615 *asymmetry* in this setting. The agent knows the realization of their cost, $c_i \in \{\bar{c}, \underline{c}\}$, whereas the
 616 server only knows the distribution from which it was drawn. In particular, the server needs to present
 617 a mechanism \mathcal{M} to an agent without knowing their actual cost.

618 D.1 Mechanism description

619 Suppose each agent independently selects their cost to be low ($c_i = \underline{c}$) with probability p_i . Let an
 620 agent with low cost \underline{c} generate \bar{m}^* data points at equilibrium on their own (and correspondingly
 621 define \bar{m}^* for a high-cost agent). Then, for some small $\varepsilon > 0$ and $\bar{m}^* \leq m_i^\uparrow \leq m_i^\downarrow$, consider the
 622 following mechanism (illustrated in Figure 8)

$$[\mathcal{M}(\mathbf{m})]_i = \begin{cases} v(m_i) & \text{for } m_i \leq \bar{m}^* \\ v(\bar{m}^*) + (\bar{c} + \varepsilon)(m_i - \bar{m}^*) & \text{for } m_i \in [\bar{m}^*, m_i^\uparrow] \\ v(m_i^\downarrow + \sum_{j \neq i} m_j) - (\underline{c} + \varepsilon)(m_i^\downarrow - m_i) & \text{for } m_i \in [m_i^\uparrow, m_i^\downarrow] \\ v(\sum_j m_j) & \text{for } m_i \geq m_i^\downarrow. \end{cases} \quad (15)$$

623 Recall from Theorem I that $\bar{m}^* \geq \bar{m}^*$
 624 since $\underline{c} \leq \bar{c}$. Thus, agents with either
 625 costs do not need additional incentive to
 626 collect data up to \bar{m}^* . Now, consider a
 627 high-cost agent. After \bar{m}^* , they need
 628 a marginal gain in value of at least \bar{c}
 629 which they do not get on their own. Addi-
 630 tional supplementary data is provided by
 631 (15) until m_i^\uparrow to incentivize a high-cost
 632 agent. It is now in their best interest to
 633 contribute m_i^\uparrow . For the low-cost agent,
 634 the marginal gain in value is at least \underline{c}
 635 until m_i^\downarrow , making this their best contribu-
 636 tion. The specific values of m_i^\downarrow and m_i^\uparrow
 637 (points D and C) can then be chosen to
 638 maximize the expected data contribution
 639 $((1 - p_i)m_i^\uparrow + p_i m_i^\downarrow)$.

640 For $\Delta m_{-i} := \sum_{j \neq i} m_j$, let \bar{m}^{\max} be the
 641 maximum amount of data a high-cost
 642 agent can be incentivized to contribute
 643 as in (11) i.e. it is defined to be $v(\bar{m}^{\max} +$
 644 $\Delta m_{-i}) = v(\bar{m}^*) + \bar{c}(\bar{m}^{\max} - \bar{m}^*)$, and
 645 \bar{m}^{\max} defined correspondingly for the
 646 low-cost agent. Then, we define m_i^\downarrow (point D) to satisfy

$$v'(m_i^\downarrow + \Delta m_{-i}) = \min\left(\max\left(\underline{c} - \frac{p}{1-p}\bar{c}, v'(\bar{m}^{\max} + \Delta m_{-i})\right), v'(\bar{m}^{\max} + \Delta m_{-i})\right). \quad (16)$$

647 Then, we can define m_i^\uparrow (point C) as the intersection of the two linear curves (starting from A and D
 648 in Fig 8):

$$v(m_i^\downarrow + \sum_{j \neq i} m_j) - (\underline{c} + \varepsilon)(m_i^\downarrow - m_i^\uparrow) = v(\bar{m}^*) + (\bar{c} + \varepsilon)(m_i^\uparrow - \bar{m}^*). \quad (17)$$

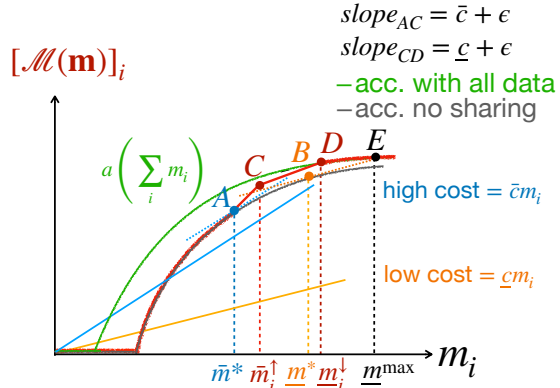


Figure 8: value shaping mechanism under unknown costs. (red curve): model value returned to agent i by the mechanism; (grey curve): model value for agent i without participation; (green curve): model value if agent i receives all the data from the other agents.

649 Note that our mechanism withholds some data from a high-cost agent resulting in a lower value
 650 model for them. This is necessary to prevent a contribution level targeted at high-cost agent from
 651 becoming attractive to a low-cost agent.

652 D.2 Analysis

653 We now analyze the properties of our expected data-maximization algorithm.

654 **Theorem VII** (Expected data maximization). *Mechanism (15) is feasible, satisfies IR, and has a*
 655 *unique Nash equilibrium: $m_i^{eq} = m_i^\uparrow$ if $c_i = \bar{c}$ and otherwise $m_i^{eq} = m_i^\downarrow$. Further, for $\varepsilon \rightarrow 0^+$,*
 656 *the mechanism (15) maximizes the expected (over the sampling of the true costs) amount of data*
 657 *collected with*

$$\sum_j (1 - p_j) m_j^\uparrow + p_j m_j^\downarrow = \max_{\mathcal{M}} \left\{ \sum_j \mathbb{E}_{\mathbf{c}} [m_j^{\mathcal{M}}], \text{ subject to } \mathcal{M} \text{ being feasible and IR} \right\}.$$

658 **Remark 5** (Decreased data collection). *By construction of our mechanism, the contribution of a*
 659 *high-cost agent would be $m_i^\uparrow \in [\bar{m}^*, \bar{m}^{\max}]$ i.e. they contribute more than they would on their own,*
 660 *but lesser than the max possible under known costs. Further, our assumption that $v(\cdot)$ is concave*
 661 *means $v'(\cdot)$ is non-increasing. Hence, (16) implies that the data contributed by a low-cost agent is*
 662 *$m_i^\downarrow \in [\bar{m}^{\max}, \underline{m}^{\max}]$. However, if $p_i \geq \frac{c}{c+\bar{c}}$, (15) always implies that $m_i^\downarrow = \underline{m}^{\max}$.*

663 We extract lesser data than if we knew the agent's true cost i.e. $m_i^\uparrow \leq \bar{m}^{\max}$. However, they also
 664 receive a model which has worse value with $v(\bar{m}^*) + \bar{c}(m_i^\uparrow - \bar{m}^*) \leq v(\sum_j m_j)$ i.e. it is not trained
 665 on the combined data. This is because if we offered a full value model to a high-cost agent at
 666 $\bar{m}^{\max} \leq m_i^\downarrow$ contribution, the low-cost agent can claim they are actually high-cost and cheat our
 667 system. Instead, now the low-cost agent will contribute $m_i^\downarrow \geq \bar{m}^{\max}$ and will receive a model trained
 668 on the combined data with value $v(\sum_j m_j)$.

669 **Theorem VIII** (Information rent). *Consider our optimal mechanism (15) with equilibrium contri-*
 670 *butions $m_i^{eq} = m_i^\uparrow$ for a high-cost agent and $m_i^{eq} = m_i^\downarrow$ for the low-cost agent. Further, let \bar{m}^**
 671 *and \underline{m}^* be the equilibrium individual contributions. Then, the utility of the high-cost agent remains*
 672 *unchanged with $v(\bar{m}^*) + \bar{c}(m_i^\uparrow - \bar{m}^*) - \bar{c}m_i^\uparrow = v(\bar{m}^*) - \bar{c}\bar{m}^*$. The utility of a low-cost agent,*
 673 *however, improves by $\left(\underline{c}(\underline{m}^{\max} - m_i^\downarrow) - v(\underline{m}^{\max} + \Delta m_{-i}) + v(m_i^\downarrow + \Delta m_{-i}) \right) \geq 0$.*

674 Because a low-cost agent can always lie and pretend to be high cost, they hold some power over the
 675 server when $m_i^\downarrow < \underline{m}^{\max}$. This is reflected in the extra utility they manage to extract and is called
 676 information rent. The utility of the high-cost agent remains unchanged since they hold no such power.

677 **E Proofs from Section 2 (Optimal Individual Contributions)**

678 **Theorem I** (Optimal individual generation). *Consider an individual agent i with marginal cost per*
 679 *data point c_i and value function v satisfying (1) working on their own. Then, the optimal amount of*
 680 *data m_i^* is:*

$$m_i^* = \begin{cases} 0 & \text{if } \max_{m_i \geq 0} u_i(m_i) \leq 0; \\ \alpha_i^*, \text{ such that } b'(\alpha_i^*) = c_i & \text{otherwise.} \end{cases} \quad (4)$$

681 *Further, for agents i, j with costs $c_i \leq c_j$, their utility satisfies $u_i(m_i^*) \geq u_j(m_j^*)$ and $m_i^* \geq m_j^*$.*

682 *Proof.* Recall that the utility function (see Eq. 3) of a single agent is:

$$u_i(m_i) = v(m_i) - c_i m_i.$$

683 Thus we have,

$$u_i'(m_i) = v'(m_i) - c_i.$$

684 Denote $\arg \max_m v(m) = 0$ as m^0 . By definition, $\forall m_i > m^0, v(m_i) = b(m_i) > 0$. Given that $b(\cdot)$
 685 is concave, $b'(m_i), m_i \geq m^0$ (or $u_i'(m_i)$) is maximized when $m_i = m^0$.

686 **Case 1 (high-cost agent):** $u_i'(m^0) \leq 0$. Then, for $\forall m_i \geq m^0, u_i'(m_i) \leq u_i'(m^0) \leq 0$. On the
 687 other hand, $\forall 0 \leq m_i \leq m^0, u_i'(m_i) = -c_i \leq 0$. Thus $u_i(m_i)$ is non-increasing, and $m_i^* = 0$. The
 688 utility function of an agent in this case is illustrated in Figure 1 (d).

689 **Case 2 (mid-cost agent):** $u_i'(m^0) > 0$ and $\max_{m_i} u_i(m_i) \leq 0$. When $u_i'(m^0) > 0$, that implies
 690 that at $m^0, b'(m^0) > c_i$. Moreover, for $m_i \geq m^0$, we have that

$$u_i'(m_i) = b'(m_i) - c_i < b'(m^0) - c_i.$$

691 Therefore, since $b(m_i)$ is concave, it is possible that $u_i(m_i)$ increases first after m^0 . However, as
 692 long as $\max_{m_i} u_i(m_i) \leq 0$, we still have that $m_i^* = 0$. The utility function of an agent in this case is
 693 illustrated in Figure 1 (c).

694 **Case 3 (low-cost agent):** $u_i'(m^0) > 0$ and $\max_{m_i} u_i(m_i) > 0$. Recall that for a mid-cost agent, it
 695 is possible that $u_i(m_i)$ increases first after m^0 . Moreover, given that $v(m_i) \leq 1$, as $m_i \rightarrow \infty$,

$$u_i'(m_i) = b'(m_i) - c_i \leq 0.$$

696 Therefore, there exists $\alpha_i^* > m^0 > 0$ such that $b'(\alpha_i^*) = c_i$. The utility function of an agent in this
 697 case is illustrated in Figure 1 (b).

698 Combining the three cases above completes the first part of the proof.

699 Next, consider two agents with costs $c_i \leq c_j$. Note that for any fixed $m, u_i(m) \geq u_j(m)$. Hence,
 700 the inequality also holds after minimizing both sides. Finally, note that if j is not a low-cost agent,
 701 it is clear that $m_i^* \geq m_j^* = 0$. If both i and j are low-cost agents, note that $m_i^* = b'^{-1}(c_i)$ and
 702 $m_j^* = b'^{-1}(c_j)$. Since $b(\cdot)$ is concave and positive, b' (and hence b'^{-1}) is non-increasing. This
 703 implies that $m_i^* \geq m_j^*$ finishing the theorem. \square

704 **F Proofs from Section 3 (Modeling Multiple Agents and Catastrophic**
 705 **Free-riding)**

706 **Theorem II** (Existence of pure equilibrium). *Consider a feasible mechanism \mathcal{M} which can be*
 707 *expressed as:*

$$[\mathcal{M}(m_i; \mathbf{m}_{-i})]_i = \max(0, \nu_i(m_i; \mathbf{m}_{-i})),$$

708 *for a function $\nu_i(m_i; \mathbf{m}_{-i})$ which is continuous in \mathbf{m} and concave in m_i . For any such \mathcal{M} , there*
 709 *exists a pure Nash equilibrium in data contributions $\mathbf{m}^{eq}(\mathcal{M})$ which for any agent i satisfies,*

$$[\mathcal{M}(\mathbf{m}^{eq}(\mathcal{M}))]_i - c_i m_i^M \geq [\mathcal{M}(m_i, \mathbf{m}^{eq}(\mathcal{M})_{-i})]_i - c_i m_i, \text{ for all } m_i \geq 0. \quad (7)$$

710 *Proof.* For a set of contributions \mathbf{m} , define the following best response mapping:

$$[B(\mathbf{m})]_i := \arg \max_{\tilde{m}_i \geq 0} \{u_i(\tilde{m}_i, \mathbf{m}_{-i}) := [\mathcal{M}(\tilde{m}_i, \mathbf{m}_{-i})]_i - c_i \tilde{m}_i\}, \quad (18)$$

711 where recall $[\mathcal{M}(\tilde{m}_i, \mathbf{m}_{-i})]_i$ is the value returned by the mechanism upon agent i submitting \tilde{m}_i
 712 data points and the rest contributing \mathbf{m}_{-i} . Note that the mapping defined above is a multi-valued
 713 function i.e. $B : \mathbb{R}^n \rightarrow 2^{\mathbb{R}^n}$. This is because the $\arg \max$ defined above may potentially return
 714 multiple values. Nevertheless, suppose that there existed a fixed point to the mapping B i.e. there
 715 existed $\tilde{\mathbf{m}}$ such that $\tilde{\mathbf{m}} \in B(\tilde{\mathbf{m}})$. Then, $\tilde{\mathbf{m}}$ is the required equilibrium contribution since by definition
 716 of the arg-max we have for any $m_i \geq 0$,

$$[\mathcal{M}(\tilde{m}_i, \tilde{\mathbf{m}}_{-i})]_i - c_i \tilde{m}_i \geq [\mathcal{M}(m_i, \tilde{\mathbf{m}}_{-i})]_i - c_i m_i.$$

717 So, we only have to prove that the mapping B has a fixed point. Since the mechanism \mathcal{M} is feasible,
 718 by Definition A and equation (1) we have

$$[\mathcal{M}(\mathbf{m})]_i \leq v(\sum_j m_j) \leq \lim_{m \rightarrow \infty} v(m) \leq 1.$$

719 This implies that

$$0 \geq u_i(\tilde{\mathbf{m}}) \leq 1 - c_i \tilde{m}_i \Rightarrow \tilde{m}_i \leq 1/c_i.$$

720 Thus, we can restrict our search space to a *convex* and *compact* product set $\mathcal{C} := \prod_j [0, 1/c_i] \subset \mathbb{R}^n$
 721 and our mapping is then over $B : \mathcal{C} \rightarrow 2^{\mathcal{C}}$. Next by assumption on the mechanism \mathcal{M} , our utility
 722 function can be written as

$$u_i(m_i, \mathbf{m}_{-i}) = \max(-c_i m_i, \nu(m_i, \mathbf{m}_{-i}) - c_i m_i),$$

723 where $\nu(m_i, \mathbf{m}_{-i}) - c_i m_i$ is concave in m_i . Unfortunately, u_i may not be quasi-concave in m_i
 724 because of the max. If it was quasi-concave, the mapping $B(\mathbf{m})$ would be continuous in \mathbf{m} and
 725 applying Kakutani's theorem would yield the existence of the required fixed point (see Maskin [36]
 726 or Acemoglu and Ozdaglar [3, Lecture 11], for details).

727 **Lemma 6** (Kakutani's fixed point theorem). *Consider a multi-valued function $F : \mathcal{C} \rightarrow 2^{\mathcal{C}}$ over*
 728 *convex and compact domain \mathcal{C} for which the output set $F(\mathbf{m})$ i) is convex and closed for any fixed \mathbf{m} ,*
 729 *and ii) changes continuously as we change \mathbf{m} . For any such F , there exists a fixed point \mathbf{m} such that*
 730 *$\mathbf{m} \in F(\mathbf{m})$.*

731 However, our utility function is not quasi-concave and the mapping B may be discontinuous. While
 732 there have been recent extensions of Kakutani's fixed point theorem to half-continuous functions (e.g.
 733 Bich [6, Theorem 3.2]), the mapping B does not satisfy this either. We next study the exact nature of
 734 discontinuity.

735 **Lemma 7.** *Consider the best-response mapping B in (18) over convex and compact domain \mathcal{C} . For*
 736 *any \mathbf{m} , either the mapping $[B(\mathbf{m})]_i$ is convex, closed, and continuous in \mathbf{m} , or $0 \in [B(\mathbf{m})]_i$.*

737 *Proof.* Figure 9 looks at the best response mapping B_i depending on the utility curve $u_i(\cdot, \mathbf{m}_{-i})$.
 738 Even if the utility itself is smoothly varying with the parameters \mathbf{m} , the best response may be
 739 discontinuous. In Fig. 9, for a small change in the utility curve between $u_2(m_i)$ to $u_3(m_i)$, the best
 740 response drastically changes from $\tilde{m}_i = 0$ (A) to $\tilde{m}_i > 125$ (B). However, this is the only source of
 741 discontinuity.

742 Recall that our utility function u_i is a max of a decreasing linear function and a concave function.
 743 Thus it has at most two local maxima: either 0, or the maxima of the concave function $f(m_i) =$
 744 $\nu_i(m_i; \mathbf{m}_{-i}) - c_i m_i$. The set of maxima of a continuous concave function is continuous, closed and
 745 convex. Hence, either 0 is part of the best response, or $[B(\mathbf{m})]_i$ is continuous, closed and convex. \square

746 Armed with Kakutani's fixed point theorem Lemma 6 and a description of the discontinuities in the
 747 best response mapping Lemma 7, we can continue with the proof of existence of a fixed point for B .
 748 Given any index set $\mathcal{I} \subseteq [n]$, we can define the following sub-domain $\mathcal{C}_{\mathcal{I}} := \prod_{i \in \mathcal{I}} [0, 1/c_i]$. Given
 749 any vector $\mathbf{p} \in \mathcal{C}_{\mathcal{I}}$, we can construct its extension $m(\mathbf{p}; \mathcal{I}) \in \mathcal{C}$ as

$$[m(\mathbf{p}; \mathcal{I})]_i := \mathbf{p}_i \text{ if } i \in \mathcal{I}, \text{ and } 0 \text{ otherwise.}$$

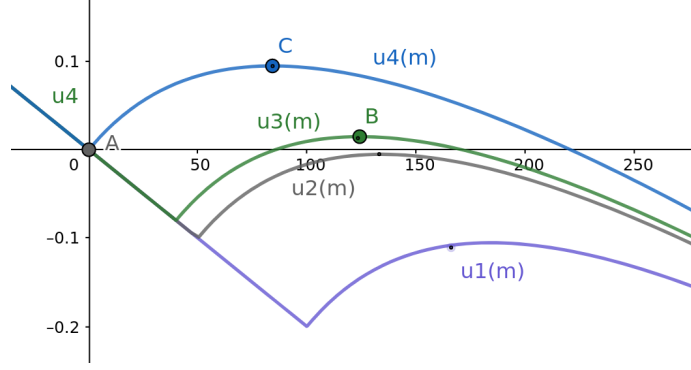


Figure 9: Utility curves $u(m_i)$ of some agent i , and the corresponding discontinuous best responses ($\tilde{m}_i \geq 0$ which maximizes $u(m_i)$). For both $u1$ and $u2$, the best response of the agent is $\tilde{m}_i = 0$ (point A), and points B and C are the best responses for $u3$ and $u4$. A small change in the utility curves (from $u2$ to $u3$) can result in a large change in the best response (from A to B).

750 We will omit the \mathcal{I} dependence and use $m(\mathbf{p})$ when clear from context. Given this mapping between
 751 sub-domain $\mathcal{C}_{\mathcal{I}}$ and the full domain \mathcal{C} , we can define a mapping:

$$B_{\mathcal{I}}(\mathbf{p}) : \mathcal{C}_{\mathcal{I}} \rightarrow 2^{\mathcal{C}_{\mathcal{I}}} := ([B(m(\mathbf{p}))]_i \text{ for } i \in \mathcal{I})$$

752 Finally, for any $\mathbf{m} \in \mathcal{C}$, define the set of indices $\mathcal{I}(\mathbf{m}) \subseteq [n]$ as

$$\mathcal{I}(\mathbf{m}) := \{i \text{ for which } 0 \notin [B(\mathbf{m})]_i . \}$$

753 Let us start from $\mathbf{m} = \mathbf{0}$. If $\mathcal{I}(\mathbf{0}) = \emptyset$, we are done since this implies $\mathbf{0} \in B(\mathbf{0})$. Otherwise,
 754 Lemma 7 states that the mapping $B_{\mathcal{I}(\mathbf{0})}(\mathbf{p})$ over the compact convex domain $\mathcal{C}_{\mathcal{I}(\mathbf{0})}$ is convex,
 755 compact and continuous. Hence, by Lemma 6, it has a fixed point such that $\mathbf{p}^1 \in B_{\mathcal{I}(\mathbf{p}^1)}$. We can
 756 inductively continue applying the same argument. If $m(\mathbf{p}^1)$ is a fixed point of the full mapping
 757 B with $m(\mathbf{p}^1) \in B(m(\mathbf{p}^1))$, we are done. Otherwise, $\mathcal{I}(m(\mathbf{p}^1)) \supset \mathcal{I}(\mathbf{0})$ and we can continue
 758 repeating the same argument inductively. Since the size of \mathcal{I} is at most n , the recursion will stop and
 759 yield a fixed point $\tilde{\mathbf{m}} \in \mathcal{C}$ such that $\tilde{\mathbf{m}} \in B(\tilde{\mathbf{m}})$. As we initially proved, this fixed point $\tilde{\mathbf{m}}$ to the
 760 best response dynamics is also the equilibrium of our mechanism.

761

□

762 **Theorem III** (Catastrophic free-riding). Consider n agents with costs $\{c_i\}$ with a unique least cost
 763 agent $c_{\min} = \min_i c_i$. Let $\{m_i^*\}$ be the equilibrium contributions of agents when alone. The standard
 764 federated learning mechanism corresponding to $[M(\mathbf{m})]_i = v(\sum_j m_j)$ for all clients i is feasible
 765 and IR, and has an unique equilibrium. At this equilibrium, only the lowest cost agent contributes:

$$m_i^{eq} = \begin{cases} m_i^* & \text{if } c_i = c_{\min} \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

766 *Proof.* Let \tilde{i} be the agent with the least cost. By Theorem I, it follows that $m_j^* < m_{\tilde{i}}^*$ for all agents
 767 $j \in [n]$.

768 First, suppose that $m_{\tilde{i}}^* > 0$. In this setting, all other agents agent j will have access to data contributed
 769 by \tilde{i} which is $m_{\tilde{i}}^*$. Now given access to this, the marginal gain in value for an additional data-point
 770 for any agent j is less than their cost i.e. $b'(m_{\tilde{i}}^*) = c_{\tilde{i}} < c_j$. Hence, it is optimal for agent j to just
 771 use $m_{\tilde{i}}^*$ data-points, and not generate any additional data-points. Thus, the equilibrium is all other
 772 agents contribute no data, and agent \tilde{i} computes $m_{\tilde{i}}^*$ datapoints as if on its own.

773 Next consider the case where $m_{\tilde{i}}^* = 0$ and hence all $m_i^* = 0$. Suppose all the other agents in total
 774 contribute Δm datapoints, which is given to all agents unconditionally. With this extra free data, it
 775 is possible that there exists some agent for whom $m_i^* > 0$. However, the incentives of the agents
 776 remain identical and so the agent with the least cost collects the most data. Given access to this
 777 $\Delta m + m_{\tilde{i}}^*$ amount of data, agent j with higher cost $c_j \geq c_{\tilde{i}}$ has no incentive to collect any additional

778 data. Hence, only agent \tilde{i} would collect any data and so $\Delta m = 0$. However, if $\Delta m = 0$, agent \tilde{i} also
 779 has no incentive to collect any data. This implies that all agents contributing 0 datapoints is the only
 780 Nash equilibrium possible. \square

781 G Proofs from Section 4 (value Shaping under Known Costs)

782 **Theorem IV** (Data maximization with known costs). *The mechanism \mathcal{M} defined by (11) is data-*
 783 *maximizing for $\varepsilon \rightarrow 0^+$. At equilibrium, a rational agent i will contribute m_i^{\max} data points where*
 784 *$m_i^{\max} \geq m_i^*$, yielding a total of $\sum_j m_j^{\max}$ data points.*

785 **Proof.** We will do the proof in two steps. Consider the best response $B_{\mathcal{M}}$ of the agents to a
 786 mechanism \mathcal{M} similar to the definition in the proof of Theorem II:

$$[B_{\mathcal{M}}(\mathbf{m})]_i := \arg \max_{m_i \geq 0} [\mathcal{M}(m_i, \mathbf{m}_{-i})]_i - c_i m_i.$$

787 We will first prove that given a fixed contribution \mathbf{m} from other users, the best response for our
 788 mechanism m_i^{\max} is higher than that of any other feasible and IR mechanism. We will then show that
 789 this necessarily implies that the equilibrium contribution of the agent is also data maximizing.

790 **Lemma 8.** *For a given data contribution \mathbf{m} and any feasible and IR mechanism $\tilde{\mathcal{M}}$, define best*
 791 *responses $B_{\mathcal{M}}(\mathbf{m})$ and $B_{\tilde{\mathcal{M}}}(\mathbf{m})$ for our mechanism \mathcal{M} (defined in (11)) and the other mechanism*
 792 *$\tilde{\mathcal{M}}$. Then, for any agent i and contribution \mathbf{m} ,*

$$[B_{\mathcal{M}}(\mathbf{m})]_i \geq [B_{\tilde{\mathcal{M}}}(\mathbf{m})]_i.$$

793 *Further, the best response $[B_{\mathcal{M}}(\mathbf{m})]_i$ is non-decreasing in the net contribution from other agents*
 794 *($\sum_{j \neq i} \mathbf{m}_j$).*

795 For now, we will assume that the above lemma and continue with our proof. As shown in the proof of
 796 Theorem II, the equilibrium of all feasible mechanisms (if they exist) lie in the range $\mathcal{C} := \prod_i [0, 1/c_i]$.
 797 Suppose that $\tilde{\mathbf{m}} \in \mathcal{M}$ is the equilibrium of mechanism $\tilde{\mathcal{M}} \in \mathcal{M}$. Note that $\tilde{\mathbf{m}}$ is also the fixed point
 798 of the best response with $\tilde{\mathbf{m}} \in B_{\tilde{\mathcal{M}}}(\tilde{\mathbf{m}})$. Now, define the following subspace

$$\mathcal{C}_{\geq \tilde{\mathbf{m}}} := \prod_j [\tilde{\mathbf{m}}_j, 1/c_j].$$

799 The set $\mathcal{C}_{\geq \tilde{\mathbf{m}}}$ is compact and convex. Thus, we can apply Theorem II to our optimal mechanism \mathcal{M}
 800 to prove that there exists an equilibrium point $\mathbf{m} \in \mathcal{C}_{\geq \tilde{\mathbf{m}}}$ such that

$$[\mathbf{m}]_i \in \arg \max_{m_i \geq \tilde{m}_i} [\mathcal{M}(m_i, \mathbf{m}_{-i})]_i - c_i m_i.$$

801 We will next show that the above point \mathbf{m} is in fact a fixed of $B_{\mathcal{M}}(\mathbf{m})$ and satisfies:

$$[\mathbf{m}]_i \in \arg \max_{m_i \geq 0} [\mathcal{M}(m_i, \mathbf{m}_{-i})]_i - c_i m_i.$$

802 Note that the only difference between the two claims is that in the latter the arg max is taken over
 803 ≥ 0 where as it was more constrained in the former. For the sake of contradiction, suppose this is
 804 not true i.e. there exists an agent i such that $\mathbf{m}_i \notin [B_{\mathcal{M}}(\mathbf{m})]_i$ and $[B_{\mathcal{M}}(\mathbf{m})]_i < \tilde{\mathbf{m}}_i$. However, this
 805 leads to a contradiction:

$$\begin{aligned} \sum_{j \neq i} \mathbf{m}_j &\geq \sum_{j \neq i} \tilde{\mathbf{m}}_j \\ \Rightarrow [B_{\mathcal{M}}(\mathbf{m})]_i &\geq [B_{\mathcal{M}}(\tilde{\mathbf{m}})]_i \geq [B_{\tilde{\mathcal{M}}}(\tilde{\mathbf{m}})]_i = \tilde{\mathbf{m}}_i. \end{aligned}$$

806 The first inequality is because $\mathbf{m} \in \mathcal{C}_{\geq \tilde{\mathbf{m}}}$. The first inequality in the second step follows from the
 807 latter part of Lemma 8 while the next inequality is from the first part. Finally, the last equality
 808 follows because $\tilde{\mathbf{m}}$ is a fixed point of $B_{\tilde{\mathcal{M}}}$. Hence, we have proven that there exists a fixed point
 809 $\mathbf{m} \in B_{\mathcal{M}}(\mathbf{m})$ such that $\mathbf{m} \in \mathcal{C}_{\geq \tilde{\mathbf{m}}}$ i.e. the equilibrium contribution of every agent under \mathcal{M} is at
 810 least as much as $\tilde{\mathcal{M}}$. \square

811 **Proof of Lemma 8.** Recall the optimal mechanism defined in (11) restated below:

$$[\mathcal{M}(\mathbf{m})]_i = \begin{cases} v(m_i) & \text{for } m_i \leq m_i^* \\ v(m_i^*) + (c_i + \varepsilon)(m_i - m_i^*) & \text{for } m_i \in [m_i^*, m_i^{\max}] \\ v(\sum_j m_j) & \text{for } m_i \geq m_i^{\max}. \end{cases} \quad (19)$$

812 For now, suppose that $m_i^* > 0$. Recall, from [case 3, Theorem I], that this implies $v'(m_i^*) =$
813 $b'(m_i^*) = c_i$.

814 First we show that m_i^{\max} is the unique equilibrium contribution for an agent i . The slope of the utility
815 of agent i is

$$u'_i(m_i; \mathcal{M}) = \frac{\partial[\mathcal{M}(\mathbf{m})]_i}{\partial m_i} - c_i.$$

816 By construction, this slope is $u'_i(m_i; \mathcal{M}) > 0$ for any $m_i < m_i^{\max}$. Suppose the contribution of
817 all other agents is fixed to $\Delta m_{-i} = \sum_{j \neq i} m_j$. The slope of the utility at m_i^{\max} is $u'_i(m_i^{\max}; \mathcal{M}) =$
818 $v'(m_i^{\max} + \Delta m_{-i}) - c_i$. Again, by construction, $m_i^{\max} + \Delta m_{-i} \geq m_i^*$. Since b is concave and b' is
819 non-increasing,

$$u'_i(m_i^{\max}; \mathcal{M}) = v'(m_i^{\max} + \Delta m_{-i}) - c_i = b'(m_i^{\max} + \Delta m_{-i}) - c_i \leq b'(m_i^*) - c_i = 0.$$

820 Thus, m_i^{\max} is the unique equilibrium contribution of agent i . Next, we have to demonstrate the
821 data-maximizing property. For the sake of contradiction, suppose there existed some other mechanism
822 $\tilde{\mathcal{M}}$ such that

$$\arg \max_{m_i} [\tilde{\mathcal{M}}(\mathbf{m})]_i - c_i m_i =: \tilde{m}_i > m_i^{\max}.$$

823 This implies that $u'_i(m_i; \tilde{\mathcal{M}}) > 0$ for any $m_i \leq \tilde{m}_i$, i.e. $\frac{\partial[\tilde{\mathcal{M}}(\mathbf{m})]_i}{\partial m_i} > c_i$. In particular, this implies
824 that

$$\frac{\partial[\tilde{\mathcal{M}}(\mathbf{m})]_i}{\partial m_i} > \frac{\partial[\mathcal{M}(\mathbf{m})]_i}{\partial m_i} \text{ for all } m_i \in [m_i^*, m_i^{\max}].$$

825 Further, $\tilde{\mathcal{M}}$ satisfies individual rationality and so at $m_i = m_i^*$ we have

$$[\tilde{\mathcal{M}}(m_i^*, \mathbf{m}_{-i})]_i \geq v(m_i^*) = [\mathcal{M}(m_i^*, \mathbf{m}_{-i})]_i.$$

826 Together, these two conditions imply that for all $m_i \in [m_i^*, m_i^{\max}]$, we have $[\tilde{\mathcal{M}}(\mathbf{m})]_i > [\mathcal{M}(\mathbf{m})]_i$.
827 In particular at $m_i = m_i^{\max}$, we have

$$[\tilde{\mathcal{M}}(m_i^{\max}, \mathbf{m}_{-i})]_i > v(\sum_j m_j).$$

828 This gives us a contradiction since it violates feasibility. Thus, m_i^{\max} is the maximum data which can
829 be extracted from agent i .

830 The proofs for the low and medium cost agents are similar, while noting that $m_i^* = 0$. This finishes
831 the proof of the first part. The second part of the lemma follows directly from the definition of \mathcal{M} and
832 the fact that the value function $v(m_i + \sum_{j \neq i} m_j)$ is non-decreasing in the contributions $\sum_{j \neq i} m_j$.
833 \square

834 **Theorem V (Incentive compatibility).** *Under given costs \mathbf{c} , consider our optimal mechanism (11)*
835 *with equilibrium contributions \mathbf{m}^{\max} , and agents working individually with equilibrium contributions*
836 *of \mathbf{m}^* . The utility of the every agent i remains unchanged:*

$$v(\sum_j m_j^{\max}) - c_i m_i^{\max} = v(m_i^*) - c_i m_i^*.$$

837 *Proof.* This statement is true by construction of our mechanism. When $\varepsilon \rightarrow 0$, the slope of the utility
838 becomes

$$u'_i(m_i; \tilde{\mathcal{M}}) = \frac{\partial[\mathcal{M}(\mathbf{m})]_i}{\partial m_i} - c_i = c_i + \varepsilon - c_i = 0 \text{ for all } m_i \in [m_i^*, m_i^{\max}].$$

839 Further, note that at $m_i = m_i^*$, we have $[\mathcal{M}(m_i^*, \mathbf{m}_{-i})]_i = v(m_i^*)$. Thus, for all $m_i \in [m_i^*, m_i^{\max}]$,
840 the utility of agent i with our mechanism \mathcal{M} remains constant and equal to the optimal individual
841 utility $u_i(m_i^*)$. \square

842 **H Proofs from Appendix D (Data Maximization with Unverifiable Costs)**

843 **Theorem VII** (Expected data maximization). *Mechanism (15) is feasible, satisfies IR, and has a*
 844 *unique Nash equilibrium: $m_i^{eq} = m_i^\uparrow$ if $c_i = \bar{c}$ and otherwise $m_i^{eq} = m_i^\downarrow$. Further, for $\varepsilon \rightarrow 0^+$,*
 845 *the mechanism (15) maximizes the expected (over the sampling of the true costs) amount of data*
 846 *collected with*

$$\sum_j (1 - p_j) m_j^\uparrow + p_j m_j^\downarrow = \max_{\mathcal{M}} \left\{ \sum_j \mathbb{E}_{\mathbf{c}} [m_j^{\mathcal{M}}], \text{ subject to } \mathcal{M} \text{ being feasible and IR} \right\}.$$

847 *Proof.* Recall that we had defined the mechanism (15) as

$$[\mathcal{M}(\mathbf{m})]_i = \begin{cases} v(m_i) & \text{for } m_i \leq \bar{m}^* \\ v(\bar{m}^*) + (\bar{c} + \varepsilon)(m_i - \bar{m}^*) & \text{for } m_i \in [\bar{m}^*, m_i^\uparrow] \\ v(m_i^\downarrow + \sum_{j \neq i} m_j) - (\underline{c} + \varepsilon)(m_i^\downarrow - m_i) & \text{for } m_i \in [m_i^\uparrow, m_i^\downarrow] \\ v(\sum_j m_j) & \text{for } m_i \geq m_i^\downarrow. \end{cases} \quad (20)$$

848 First, we have to show that m_i^\uparrow and m_i^\downarrow are equilibrium for the high and low cost players \bar{c} and \underline{c}
 849 respectively. For the sake of simplicity, we first assume that $\bar{m}^* > 0$ and $\underline{m}^* > 0$. The proofs directly
 850 extend to the other cases. Now, note that $v'(\bar{m}^*) = b'(\bar{m}^*) = \bar{c}$. Thus, by constructions, we have
 851 that for a high cost agent,

$$u'_i(m_i; \mathcal{M}) = \frac{\partial [\mathcal{M}(\mathbf{m})]_i}{\partial m_i} - \bar{c} > 0 \text{ for all } m_i \leq m_i^\uparrow.$$

852 where as for $m_i > m_i^\uparrow$, the slope $u'_i(m_i; \mathcal{M}) = \underline{c} + \varepsilon - \bar{c} < 0$. Assuming ε is small enough, a high
 853 cost agent obtains optimal utility at m_i^\uparrow . Similarly, for the low cost agent, $u'_i(m_i; \mathcal{M}) > 0$ for all
 854 $m_i < m_i^\downarrow$ and is negative after (similar to Theorem IV). Thus, the optimum contribution of the low
 855 cost player is m_i^\downarrow .

856 Next, recall that we had defined in (16) that m_i^\downarrow satisfies

$$m_i^\downarrow = \min \left(\max \left(b'^{-1} \left(\underline{c} - \frac{p_i}{1-p_i} \bar{c} \right) - \Delta m_{-i}, \bar{m}^{\max} \right), \underline{m}^{\max} \right). \quad (21)$$

857 We will show that m_i^\downarrow defined this way maximizes the expected data for agent i :

$$\max_{m_i^\downarrow} \left\{ (1 - p_i) m_i^\uparrow + p_i m_i^\downarrow \right\} \text{ subject to } m_i^\uparrow, m_i^\downarrow \text{ are feasible for } \mathcal{M}. \quad (22)$$

858 This involves some variational calculus (see Fig. 10). As shown in Fig. 10, reducing the value of m_i^\downarrow
 859 results in an increase in m_i^\uparrow . Suppose we push the blue bar vertically by a small value dx . Because the
 860 slope of AC is \bar{c} , this results in increase of $\frac{dx}{\bar{c}}$ in m_i^\uparrow . Correspondingly, we can show that the decrease
 861 in m_i^\downarrow will be $\frac{dx}{\underline{c} - v'(m_i^\downarrow + \Delta m_{-i})}$. Putting these together, the net expected change in data contribution is

$$(1 - p_i) \frac{dx}{\bar{c}} - \frac{dx}{\underline{c} - b'(m_i^\downarrow + \Delta m_{-i})}.$$

862 The local unconstrained maxima can then be derived by setting the above to 0 i.e when

$$b'(m_i^\downarrow + \Delta m_{-i}) = \underline{c} - \frac{p_i}{1-p_i} \bar{c}.$$

863 Of course, we have to respect the constraints that $m_i^\downarrow \in [\bar{m}^{\max}, \underline{m}^{\max}]$ giving us our final result.
 864 Thus, the value of m_i^\downarrow as chosen by (22) is optimal for these class of mechanisms.

865 Now, we have to show that any data-maximizing mechanism corresponds to \mathcal{M} with some choice of
 866 m_i^\downarrow . Consider a mechanism $\tilde{\mathcal{M}}$ whose equilibrium contributions are (\tilde{m}, \tilde{m}) for a high and low-cost
 867 agent respectively (see points I and J in Fig. 10). Now, from the optimality of \underline{m}^{\max} , we know that
 868 $\tilde{m} \leq \underline{m}^{\max}$. Let us connect \bar{m}^* (point A) to \tilde{m} (point I) and then to \underline{m} (point J). Recall that we
 869 assumed that $\tilde{\mathcal{M}}$ is different from \mathcal{M} in (15). This means that the slope AI $\neq \bar{c}$ or IJ $\neq \underline{c}$. Consider

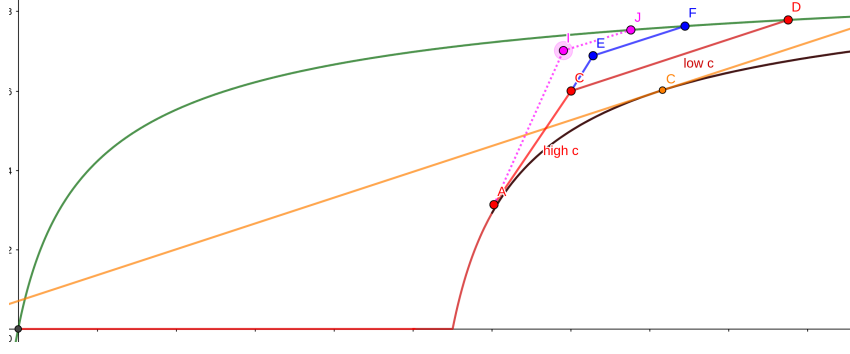


Figure 10: Value shaping mechanism under unknown costs. (*red curve*): model value returned to agent i by the mechanism; (*grey curve*): model value for agent i without participation; (*green curve*): model value if agent i receives all the data from the other agents. Point C and D (in red) and points E and F (in blue) represent two different choices for $(m_i^\uparrow$ and $m_i^\downarrow)$ respectively. If we choose a smaller value of m_i^\downarrow (shown in blue by point F), we would see an increase in m_i^\uparrow to point E. Thus, the optimum value balances these two depending on the probability p_i . Finally, points I and J (in magenta) represent potential other mechanisms \mathcal{M} .

870 the latter. Combined with I and J corresponding to equilibria, we have slope of $IJ > \underline{c}$. This implies
 871 that starting from point I, we could have instead drawn a line segment of slope \underline{c} and increased the
 872 data contribution by the low cost agent, while keeping the contribution of the high-cost agent fixed.
 873 Similarly, we can show that the optimal slope for AI is \bar{c} . Together, this implies that any optimal
 874 mechanism \mathcal{M} must be of the form (15), finishing our proof. \square

875 **Theorem VIII** (Information rent). *Consider our optimal mechanism (15) with equilibrium contri-*
 876 *butions $m_i^{eq} = m_i^\uparrow$ for a high-cost agent and $m_i^{eq} = m_i^\downarrow$ for the low-cost agent. Further, let \bar{m}^**
 877 *and \underline{m}^* be the equilibrium individual contributions. Then, the utility of the high-cost agent remains*
 878 *unchanged with $v(\bar{m}^*) + \bar{c}(m_i^\uparrow - \bar{m}^*) - \bar{c}m_i^\uparrow = v(\bar{m}^*) - \bar{c}\bar{m}^*$. The utility of a low-cost agent,*
 879 *however, improves by $(\underline{c}(m_i^{\max} - m_i^\downarrow) - v(m_i^{\max} + \Delta m_{-i}) + v(m_i^\downarrow + \Delta m_{-i})) \geq 0$.*

880 *Proof.* For a high cost player, the statement easily follows since $\frac{\partial[\mathcal{M}(\mathbf{m})]_i}{\partial m_i} - \bar{c} = 0$ for all $m_i \in [\bar{c}^*, c_i^\uparrow]$.
 881 Thus, a high cost player's utility remains constant during this period and is equal to utility at $m_i = \bar{m}^*$
 882 which is $v(\bar{m}^*) - \bar{c}\bar{m}^*$.

883 For a low cost agent, $\frac{\partial[\mathcal{M}(\mathbf{m})]_i}{\partial m_i} - \underline{c} = 0$ for all $m_i \in [m_i^\uparrow, m_i^\downarrow]$, and hence their utility is constant in
 884 this region. In particular, the difference in utility with mechanism \mathcal{M} and alone is

$$v(m_i^\downarrow + \Delta m_{-i}) - \underline{c}m_i^\downarrow - v(m_i^{\max} + \Delta m_{-i}) + \underline{c}m_i^{\max}.$$

885 The above quantity is always non-negative since $v'(m_i + \Delta m_{-i}) \leq \underline{c}$ for all $m_i \in [m_i^\downarrow, m_i^{\max}]$. \square