
No Bidding, No Regret: Pairwise-Feedback Mechanisms for Digital Goods and Data Auctions

Zachary Robertson
Department of Computer Science
Stanford
zroberts@stanford.edu

Oluwasanmi Koyejo
Department of Computer Science
Stanford
sanmi@stanford.edu

Abstract

The growing demand for data and AI-generated digital goods, such as personalized written content and artwork, necessitates effective pricing and feedback mechanisms that account for uncertain utility and costly production. Motivated by these developments, this study presents a novel mechanism design addressing a general repeated-auction setting where the utility derived from a sold good is revealed post-sale. The mechanism's novelty lies in using pairwise comparisons for eliciting information from the bidder, arguably easier for humans than assigning a numerical value. Our mechanism chooses allocations using an epsilon-greedy strategy and relies on pairwise comparisons between realized utility from allocated goods and an arbitrary value, avoiding the learning-to-bid problem explored in previous work. We prove this mechanism to be asymptotically truthful, individually rational, and welfare and revenue maximizing. The mechanism's relevance is broad, applying to any setting with made-to-order goods of variable quality. Experimental results on multi-label toxicity annotation data, an example of negative utilities, highlight how our proposed mechanism could enhance social welfare in data auctions. Overall, our focus on human factors contributes to the development of more human-aware and efficient mechanism design.

1 Introduction

Marketplaces generating digital goods, such as personalized written content and artwork based on user requests, have garnered significant attention in recent years due to their ability to scale and adapt to user preferences [Mor, 2023, Paul and Dang, 2023]. Such generative marketplaces possess immense potential to revolutionize the economy through applications such as online advertising [Paul and Dang, 2023], a market where spending is expected to exceed \$700 billion in 2023 [Sta, 2023]. However, they face challenges in collecting accurate and timely human feedback, as well as in managing compute costs for the most advanced models, which the CEO of OpenAI has described as "eye-watering" [Karpf, 2023]. This problem is particularly acute since the value each user derives from a fulfilled request is typically only known after allocation.

In this paper, we examine the general repeated-auction setting where the utility derived from sold digital goods is revealed to bidders post-sale. In this setting, digital goods are "made-to-order" based on user requests. A key challenge is that users could provide inaccurate or misleading feedback, which would harm revenue generation. To address these challenges, we propose an auction mechanism that is robust to strategic reporting on the user side and no-regret in revenue on the market side. Our pricing mechanism is based on a pairwise comparison model that asks the user to report if the value of their allocation is above an arbitrarily selected reference point, avoiding the learning-to-bid problem that has been a point of concern in previous works [Feng et al., 2018, Guo et al., 2022].

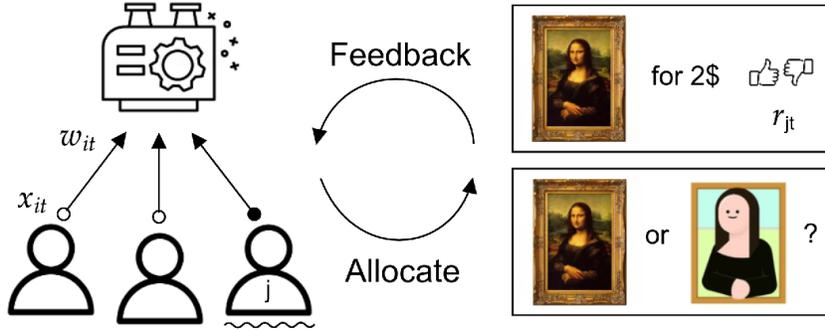


Figure 1: This is an illustration of our proposed mechanism introduced in Section 4. In the left panel we start period t , each agent $i \in [n]$ submits a request w_{it} (prompt) for a made-to-order digital good, and then the mechanism determines a boolean-valued allocation assignment x_{it} for each agent. If the j^{th} agent receives an allocation for their request ("Mona Lisa") then the agent receives a digital good made according to their request. After a digital good is produced, the agent self-reports a Boolean value, denoted by r_{jt} , indicating if the value of their allocation is above an arbitrarily selected reference point. Alternatively, the mechanism produces multiple digital goods and has the user rank them. A priori, it is not obvious how to design a mechanism robust to strategic manipulation. This paper devises a simple mechanism based on a second-price auction that is asymptotically truthful.

Our main contribution is a novel auction mechanism for selling digital goods that are costly to produce and whose utility to a particular user is uncertain. We outline our main contributions below:

1. **Feedback-based auction mechanisms:** This study introduces a feedback-driven, contextual, asymptotically truthful mechanism that eliminates the need for users to know their value for generating a digital good beforehand. By allocating goods to agents and subsequently collecting feedback on their satisfaction, the mechanism effectively sidesteps the learning-to-bid problem.
2. **Analysis of efficiency:** An in-depth analysis of the proposed mechanism is presented. We establish finite-time regret bounds for truthful reporting, participation, and welfare-revenue generation against the standard second-price auction. We also show the underlying expected utilities can be identified from pairwise comparisons without relying on distributional assumptions.
3. **Welfare maximizing data acquisition:** We also explore how to use our mechanism as a payment rule for toxicity annotation, a setting with negative utilities that are only realized after the mechanism purchases a label from a user. We discuss this setting further in Section 5.2.

We tackle two main technical challenges. First, akin to [Nazerzadeh et al., 2013], we utilize a learning algorithm for the expected utility function but diverge from their history-dependent rule, which requires $O(H^2)$ calls to the model for H allocation rounds to enhance computational efficiency. Our Lemma 5.2 outlines that providing inaccurate or misleading feedback isn't particularly profitable, achieved through a refined strategic reporting and regret analysis. Secondly, our mechanism features a simplified reporting rule, negating the need for agents to accurately learn to bid, as seen in works like [Feng et al., 2018, Guo et al., 2022]. Although common among all prior work we are aware of, we also dismiss the unrealistic assumption of precise value reporting as it essentially requires reporting a real number with infinite precision. In Theorem 5.5, we establish a mechanism that uses feedback reports free from distributional assumptions on underlying utilities.

The remainder of this paper is organized as follows. In Section 2, we review related work in auctions with incomplete information, mechanism design, and data pricing. In Section 3, we formalize our setting and describe the background for our approach. In Section 4 we introduce our proposed mechanism, and in Section 5 we present our main results for the proposed mechanism.

2 Related work

Our work builds on research in the fields of auctions with incomplete information, machine learning for mechanism design, and data pricing.

2.1 Mechanism design and machine learning

Our work intertwines mechanism design and machine learning to address challenges in pricing and feedback systems for digital goods and data [Balcan et al., 2005, Devanur and Kakade, 2009, Babaioff et al., 2009, 2015]. We also draw from the intersection of machine learning and mechanism design for allocating digital goods [Immorlica et al., 2005, Agarwal et al., 2009, Mahdian and Tomak, 2008, Nazerzadeh et al., 2013].

The proposed mechanism is distinguished by two key differences from existing works that also investigate online auction design [Devanur and Kakade, 2009, Nazerzadeh et al., 2013, Babaioff et al., 2009]. For exact truthfulness, strict characterizations on attainable regret rates have been established for deterministic payment rules [Devanur and Kakade, 2009, Babaioff et al., 2009]. Weaker asymptotic truthfulness has also been considered under the pay-per-action framework known in online advertising [Immorlica et al., 2005, Mahdian and Tomak, 2008, Nazerzadeh et al., 2013]. Most related to our work is [Nazerzadeh et al., 2013], which studies the pay-per-action setting, which allows reporting value after allocation, and proposes a history-dependent pricing rule. Our work builds on this progress by studying contextual auctions, removing the history-dependent pricing rule, and eliminating the requirement for an exact value offering a human-aware perspective and a more computationally-efficient mechanism. See Table 1 for a comparison between these approaches and our proposal.

2.2 Partially-informed auctions

Auction design with incomplete information has been a topic of interest in recent years [Bergemann and Pesendorfer, 2007, Feng et al., 2018, Epasto et al., 2021, Guo et al., 2022]. In particular, [Bergemann and Pesendorfer, 2007] considers single-item multi-bidder auctions where information is only partially revealed to bidders. [Feng et al., 2018] investigates the single-item setting where bidders learn to bid with partial feedback and obtain no-regret against the best fixed bid in hindsight. [Guo et al., 2022] extend this analysis by considering context and propose no-regret algorithms that are efficient from the buyer’s perspective with applications to privacy. However, all of this prior work requires that the agent learn-to-bid, which requires additional effort and is commonly understood to lower an agent’s welfare [Cai et al., 2015]. Our work takes a different approach to study partially-informed auctions by focusing on user compatibility to provide their value after allocation while still maintaining the connection to these privacy considerations [Epasto et al., 2021]. See Table 1 for a comparison between [Guo et al., 2022] and our proposal.

2.3 Welfare and truthful elicitation

In the realm of data pricing and acquisition, we draw from research on learning-based data pricing [Chen et al., 2023, Zhao and Ermon, 2021, Karimireddy et al., 2022], peer-prediction mechanisms [Prelec, 2004, Witkowski and Parkes, 2012, Cai et al., 2015]. In particular, [Cai et al., 2015] develops a model for constructing statistical estimators in the presence of costly information revelation, while [Prelec, 2004] and [Witkowski and Parkes, 2012] propose peer-prediction techniques for eliciting truthful information without the need for ground-truth data. These works have made significant contributions to understanding optimal mechanisms but are more concerned with obtaining truthful responses rather than the socially efficient allocation of goods.

In recent years, the question of worker welfare during data acquisition has become a central issue. In particular, recent work on OpenAI’s toxicity filter had to be halted because data annotation had a traumatic effect on workers [Perrigo, 2023]. Toxic data annotation, in general, is known to have traumatic effects on workers [Burns et al., 2008, Arsht and Etcovitch, 2018, Steiger et al., 2021, Perrigo, 2023]. In particular, [Steiger et al., 2021] proposes preventing or reducing exposure as a potential technological intervention strategy. One of our contributions is to formalize this problem as an instance of a reverse auction (negative utilities) in our setting and demonstrate theoretically and empirically that we can asymptotically maximize the social welfare of workers.

Table 1: Comparison of our mechanism with prior work

Auction Mechanism	Strategy Robust	Feedback Reporting	Local Payments	Efficient User-Strategy
[Devanur and Kakade, 2009]	✓	×	✓	×
[Nazerzadeh et al., 2013]	✓	✓	×	×
[Guo et al., 2022]	✓	×	✓	×
This Work	✓	✓	✓	✓

3 Preliminaries

In this section, we overview the problem setting under consideration and introduce our key definitions. Discussion of how to learn a payment rule is discussed in Section 4. As an example (Figure 1), agents could be competing for resources to generate artwork, where the space of prompts is \mathcal{W} , the space of digital goods is \mathcal{O} , and the mechanism determines the price and allocation of resources. We summarize our notation in Appendix A provided in the supplementary materials.

Problem Setting: We consider a scenario where a set of n agents compete for allocations across discrete periods $t = 1, 2, \dots, H$, up to a horizon H . A mechanism \mathcal{M} oversees pricing for allocations. At each period t the following happens:

1. Each agent $i \in [n]$ submits a request $w_{it} \in \mathcal{W}$ sampled independently from one another, and a Boolean-valued array of allocations x_{it} is generated for the agents.
2. If the j^{th} agent receives an allocation, then the agent receives a digital good $o_{jt} \in \mathcal{O}$ sampled some distribution conditioned on w_{jt} which is sold to the agent.
3. The non-negative value of the agent who receives an allocation during period t is denoted by a bounded random variable $u_{jt} : \mathcal{W} \times \mathcal{O} \rightarrow [0, 1]$.
4. The agent pays an amount p_{jt} determined by the mechanism and then reports a Boolean variable $r_{jt}(c)$ indicating if u_{jt} is above some value c which is to be chosen randomly.

We’ll emphasize that u_{it} has randomness from the requests w_{it} and the mechanism that generates outputs o_{it} . Since agent’s know their requests, we’ll use $u_{it}(w)$ to denote the utility random variable given a request w . Our main assumption is an independence assumption on the requests.

Assumption 3.1. For each agent $i \in [n]$, their request and output sequences w_{it} and o_{it} are independent of other agents and allocations.

Assumption 3.1 allows sequential generation, such as written content, but precludes collusion among agents or exploiting specific prompts for high utility. We also introduce notation for an ideal setting where we have perfect knowledge of the agents’ expected values.

$$\mu_i(w_{it}) := \mathbb{E}[u_{it}(w_{it}) | \{(w_{ik}, u_{ik})\}_{k < t}], \quad x_{it} := \mathbb{I}(\mu_i(w_{it}) \geq \mu_j(w_{jt}) \forall j) \quad (1)$$

This defines the agent’s expected utility and allocation. We will design \mathcal{M} (Section 4) so that agents have no-regret for participating and truthful reporting along with comparable revenue and social welfare to a standard auction format. To introduce the key definitions, we just need to know \mathcal{M} will determine allocations \hat{x}_{it} and prices p_{it} using value estimates $\hat{\mu}_{it} \sim \mu_i$ from the reports.

Definition 3.2. The i^{th} agent is considered truthful if $r_{it}(c) = \mathbb{I}(u_{it} \geq c)$ for all t and $c \in [0, 1]$.

This means they respond to the queries in Fig 1 accurately. In our setting, there are numerous reporting strategies $r_{it} : (\mathcal{W} \times \mathcal{O}) \times [0, 1] \rightarrow \{0, 1\}$ the agents could use. Our focus lies in designing a mechanism such that agent incentives are aligned with truth-telling. Another important criterion for each agent is that they have no-regret for participation.

Definition 3.3. \mathcal{M} is asymptotically ex ante individually rational if each agent $i \in [n]$ has no-regret for participation when they are truthful. Specifically, the long-term total utility of the agent is nonnegative:

$$\liminf_{H \rightarrow \infty} \mathbb{E} \left[\sum_{t=1}^H \hat{x}_{it} \mu_i - p_{it} \right] \geq 0$$

While definition 3.3 captures the rationality of each truthful agent’s participation, it does not guarantee no-regret against other reporting strategies.

Definition 3.4. Let $U_i(H)$ be the expected total utility of the i^{th} agent using a truthful reporting strategy and $\hat{U}_i(H)$ be the maximum expected profit whenever all other agents are truthful. We say that \mathcal{M} is asymptotically truthful if truthful reporting is no-regret against strategic reporting:

$$\hat{U}_i(H) - U_i(H) = o(H)$$

This definition is similar to previous definitions in that it ensures that deviating from truthful reporting is relatively unprofitable [Pavan et al., 2009, Nazerzadeh et al., 2013]. The main distinction is that it is regret-based which enables us to obtain rates in our analysis. While an asymptotic definition seems limiting, strong notions, such as dominant strategy incentive compatibility, are achievable only in limited settings [Pavan et al., 2009, Kakade et al., 2013].

We also desire \mathcal{M} to have no-regret against an idealized auction. We compare the welfare and revenue of our mechanism to a baseline given by the second-price auction known to be welfare and revenue maximizing [Myerson, 1981]. In this format, allocations go to the highest bidder, say, the i^{th} agent who will pay $\gamma_t = \max_{j \neq i} \mu_j(w_{jt})$ to \mathcal{M} . Otherwise, they pay nothing.

Definition 3.5. We say \mathcal{M} is asymptotically ex-ante welfare maximizing if it has no-regret against the welfare generated by a second-price auction:

$$\mathbb{E} \left[\sum_{t=1}^H \sum_{i=1}^n \hat{x}_{it} \mu_i(w_{it}) \right] - \mathbb{E} \left[\sum_{t=1}^H \max_i (\mu_i(w_{it})) \right] = o(H)$$

Definition 3.6. We say \mathcal{M} is asymptotically equivalent to the revenue of the second-price auction if it has no-regret against the revenue generated by a second-price auction:

$$\mathbb{E} \left[\sum_{t=1}^H \sum_{i=1}^n \hat{x}_{it} p_{it} \right] - \mathbb{E} \left[\sum_{t=1}^H \gamma_t \right] = o(H)$$

4 The proposed mechanism

Algorithm 1: Feedback-Driven Mechanism

Input : Exploration rate η_t , agent submissions w_{it}

Output : Tuple of context-report pairs

```

1 for  $t = 1, 2, \dots$  do
2   if explore with probability  $\eta_t$  then
3      $i = \text{sample}([1, \dots, n])$ ;
4      $x_{it} = 1$ ;
5      $c_{it} = \text{sample}([0, 1])$ ;
6      $p_{it} = 0$ ;
7      $r_{it}(c_{it}) = \text{agent-report}(w_{it}, c_{it})$ ;
8   else
9      $i = \text{argmax}_j \hat{\mu}_{jt}(w_{jt})$ ;
10     $x_{it}, y_{it} = 1$ ;
11     $p_{it}, c_{it} = \text{max}_{j \neq i} \hat{\mu}_{jt}(w_{jt})$ ;
12     $r_{it}(c_{it}) = \text{agent-report}(w_{it}, c_{it})$ ;
13  end
14 end

```

In our approach, as illustrated in Figure 1 and implemented in Algorithm 1, we aim to estimate the utility function of each agent using a learning algorithm \mathcal{L} , connecting with the high-level goals of the paper by designing an auction mechanism that improves welfare in digital goods and data auctions. Ideally, we would know μ_i for each agent $i \in [n]$ and allocate using a second-price auction. The challenge lies in the potential misreporting of observed utilities by agents to gain utility.

The basic mechanism is a second-price payment rule estimated with a learning algorithm \mathcal{L} . We fit $\hat{\mu}_{it}$ using \mathcal{L} and a data set of context and reporting tuples $\{(w_{ik}, r_{ik}, c_{ik})\}_{k \in S_{it}}$, where $c_{ik} = p_{ik}$ represents the price comparison, and S_{it} denotes periods of allocation to the i^{th} agent up to time t . Simultaneously, the proposed mechanism performs a variant of ϵ -greedy allocation, allocating to agents with the highest estimated value during exploitation rounds indicated by y_{it} or exploring by allocating for free to a randomly chosen agent with probability $\eta_t \in [0, 1]$. During exploration rounds, agents still compare to a price point c_{it} sampled from a distribution over $[0, 1]$. Finally, the exploitation round payments are $\hat{\gamma}_t = \max_{j \neq i} \hat{\mu}_{jt}(w_{jt})$, with i indicating the allocated agent. In general, the payment p_{it} equals $y_{it} \hat{\gamma}_t$.

To study allocation and payment, we define the best empirical estimate under \mathcal{L} of an agent's expected utility and allocation using the data set $\{(w_{ik}, c_{ik}, r_{ik})\}_{k \in S_{it}}$.

$$\hat{\mu}_{it}(w_{it}) := \mathbb{E}_{\mathcal{L}}[u_{it} | \{(w_{ik}, c_{ik}, r_{ik})\}_{k \in S_{it}}], \quad \hat{x}_{it} := \mathbb{I}(\hat{\mu}_{it}(w_{it}) > \hat{\mu}_{jt}(w_{jt}) \forall j) \quad (2)$$

where $\mathbb{E}_{\mathcal{L}}$ is an estimate under \mathcal{L} for the true expected value. Discussion of a concrete choice of a learning algorithm is delayed until Section 5.1. It is worth remarking that these definitions differ from equation 1 because the agent merely reports their relative utility against c_{ik} .

It is worth making a few remarks comparing our mechanism (see Table 1) to related works implementing online auctions with learned expected values [Devanur and Kakade, 2009, Nazerzadeh et al., 2013, Babaiouf et al., 2015, Guo et al., 2022]. Our design deviates significantly in key areas. We allow agents to self-report their satisfaction, circumventing the "learn-to-bid" assumption adopted by [Guo et al., 2022]. We introduce a simplified reporting rule that directly links reports to current payments via binary feedback, which simplifies value reporting for agents, a notable contrast from all of these works. This strategy also contrasts with the computationally demanding history-dependent payment rule used by [Nazerzadeh et al., 2013], which scales quadratically. This tailored approach marks our mechanism as both efficient and user-side oriented, as per our prior technical discussion.

5 Analysis

We develop sufficient conditions for the learning algorithm \mathcal{L} applied in Algorithm 1 to estimate the μ_i that results in a mechanism \mathcal{M} that meets our social efficiency criterion. We then explore how to implement the learning algorithm and examine a relevant data acquisition example. Our main condition involves the error from \mathcal{L} applied to truthful reporting data from exploration rounds:

$$\Delta_t := \mathbb{E}[\max_k |\mu_k(w_{kt}) - \hat{\mu}_{kt}(w_{kt} | x_{kt'} = 1, y_{kt'} = 0, 0 \leq t' < t)|] \quad (3)$$

As the number of exploitation allocations is dependent on the other agents' behavior, we focus on performance using just the randomly allocated exploration rounds. Our main result offers valuable insights into the performance of mechanism \mathcal{M} using intuitive conditions on the learning algorithm \mathcal{L} .

Theorem 5.1. *Suppose our mechanism \mathcal{M} estimates agent values bounded to the unit-interval using some learning algorithm \mathcal{L} . Suppose the expected error of this algorithm is monotone decreasing in the number of samples and that for all time,*

$$\sum_{t=1}^H \Delta_t = o\left(\sum_{t=1}^H \mathbb{E}[\eta(t)]\right) = o(H) \quad (4)$$

then the mechanism \mathcal{M} satisfies the following:

1. *Is asymptotically individually rational and asymptotically truthful*
2. *Is asymptotically welfare and revenue maximizing.*
3. *Compared to a second-price auction \mathcal{M} can obtain welfare and revenue regret $\tilde{O}(H^{2/3})$ if there is a learning algorithm \mathcal{L} for valuations with slow learning rate $\tilde{O}(H^{-1/2})$ and regret $\tilde{O}(H^{1/2})$ if there is an algorithm with a fast rate $\tilde{O}(H^{-1})$.*

The primary assumption is that the expected sum of errors for the learning algorithm decreases quickly relative to the sum of expected regret terms $\mathbb{E}[\eta(t)]$ over the entire horizon H . Although we can always increase the number of exploration rounds for a consistent algorithm to meet this condition, doing so may significantly reduce revenue. We discuss a concrete learning algorithm in Section 5.1.

The proof of this result relies on determining whether strategic reporting can be profitable for an agent, which we bound in terms of the error induced by the learning algorithm. We outline the main steps here and defer the proof to Appendix B.1 in the supplementary materials.

Lemma 5.2. *A single agent providing misleading feedback can increase their expected utility up to time H by no more than $6\sum_{t \in [H]} \Delta_t$.*

Proof Sketch. The proof of Lemma 5.2 relies on evaluating the expected utility of an agent deviating from a truthful strategy. We first ignore exploration allocations since they are strategy independent. We then consider an agent who deviates from a truthful reporting strategy T to another strategy L . We fix $\hat{\mu}_{it}$ as the utility estimated from data collected under strategy T . We then analyze the expected utility of both strategies. The expected profit of such a strategy is given by the difference between the expected utilities under the new strategy L and the truthful strategy T . We decompose this profit into three cases, corresponding to the different allocation times for the item to the agent: those in $S_L \setminus S_T$, those in $S_T \setminus S_L$, and those in $S_T \cap S_L$. For each case, we establish bounds on the differences between the expected utilities of the two strategies. These bounds involve the estimates of utilities for different agents under both strategies and the true utilities of the allocated items. We then combine the results of these cases and obtain an upper bound on the profit of deviating from the truthful strategy, which is $O(\mathbb{E}[\sum_t \Delta_t])$ over all time steps up to time H . \square

This bound is worse than [Nazerzadeh et al., 2013] by a constant factor since we do not use a history-dependent payment rule which can correct for its past estimation errors. Unlike this work, we use a regret-based framework to proceed with the rest of the analysis.

Lemma 5.3. *\mathcal{M} is asymptotically truthful and individually rational.*

Proof Sketch. Notice that the expected profit from exploration rounds is the same between strategies. Therefore, by Lemma 5.2 we have $\hat{U}_i(H) - U_i(H) \leq 6\sum_{t \in [H]} \Delta_t$. We also have asymptotic individual rationality because we can bound the overcharges to the agent,

$$\mathbb{E} \left[\sum_{t \in S_T} \hat{\gamma}_t - \mu_i(w_{it}) \right] \leq \mathbb{E} \left[\sum_{t \in S_T} \hat{\mu}_{it}(w_{it}) - \mu_t(w_{it}) \right] \leq \mathbb{E} \left[\sum_t \Delta_t \right]$$

As we send the horizon to infinity, we see that these overcharges are bounded by the error rate. By assumption, this term is dominated by the free allocations, so we have asymptotic ex-ante individual rationality. \square

Now we examine the question of regret concerning welfare and revenue objectives. Our mechanism may lose revenue during exploration and due to estimation error during exploitation.

Lemma 5.4. *\mathcal{M} is asymptotically welfare and revenue maximizing. Compared to a second-price auction \mathcal{M} can obtain welfare and revenue regret $\tilde{O}(H^{2/3})$ if there is a learning algorithm \mathcal{L} for valuations with slow learning rate $\tilde{O}(H^{-1/2})$ and regret $\tilde{O}(H^{1/2})$ if there is an algorithm with a fast rate $\tilde{O}(H^{-1})$.*

Remarks: Our results improve upon [Nazerzadeh et al., 2013], which does not provide finite-time regret rates for their algorithm. Furthermore, compared to [Guo et al., 2022], which only considers a single agent against a post price, our mechanism considers a multi-agent setting without assuming the agents' ability to bid. In this work, they establish this bound for applications to user privacy where the context is masked in order to preserve privacy. Since masking is deterministic in this work, we can use standard results from realizable learning theory, see Anthony et al. [1999], to conclude our algorithm also achieves $\tilde{O}(H^{-1/2})$ in the stochastic setting. In general, it is unclear if using an adaptive algorithm would help given $\Omega(H^{2/3})$ lower-bounds on regret in this setting [Devanur and Kakade, 2009, Babaioff et al., 2009]. Despite these limitations, our mechanism showcases the importance of feedback and welfare in shaping digital goods and data auction mechanisms that are more efficient and user-friendly.

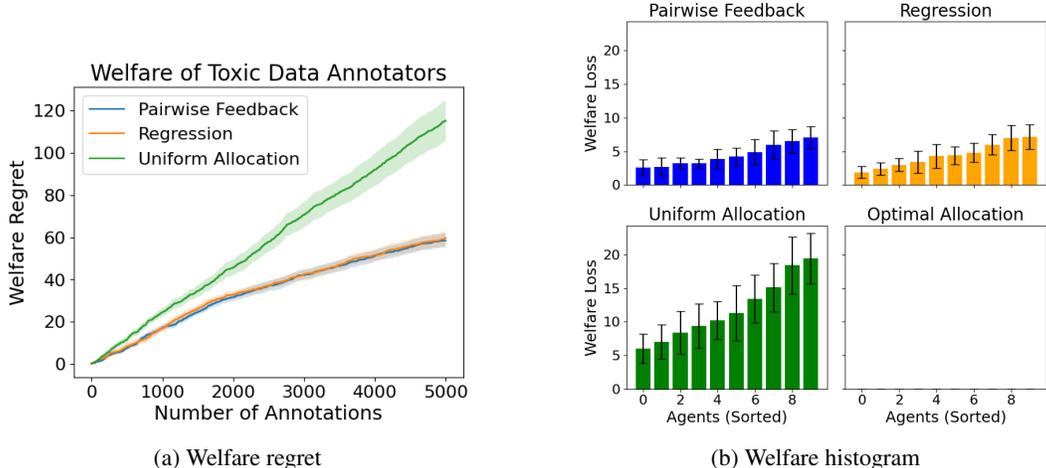


Figure 2: Welfare is the sum of the utilities of all agents across all the allocation periods. We compare the expected welfare regret (a) of different allocation mechanisms for data acquisition vs. an ideal allocation mechanism. Our relative feedback approach elicits relative utility information, regression fits values using utility reports from agents, and uniform allocation methods such as peer prediction make assignments at random. The welfare histogram (b) shows the distribution of welfare losses for each agent under different allocation mechanisms.

5.1 Estimating the value model with pairwise feedback

To make our theory concrete, we can consider μ_i that follows a linear model with d -features and with standard independent noise and realizability assumption. We will call the resulting mechanism using uniformly sampled comparison prices, linear regression algorithm \mathcal{L} , and exploration rate η by $\mathcal{M}_\eta(\text{linear})$.

Theorem 5.5. *For an exploration rate $\eta_t = t^{-1/3} \cdot (n \log(t))^{(1+2\epsilon)/3}$ we have that $\mathcal{M}_\eta(\text{linear})$ satisfies the conditions of Theorem 5.1 and so is asymptotically individually rational, incentive compatible, and no-regret in revenue and welfare.*

The full proof is provided in Appendix B.2 in the supplementary materials. The main technical step in the proof is to identify the underlying utility from pairwise reports. In our mechanism, we obtain reports r_{it} from agents regarding if the random utility u_{it} derived from their allocation for a request w_{it} is satisfactory. This is based on being above or below some reference payment c . Our model is that $r_{it}(c) = \mathbb{I}(u_{it} \geq c)$. In particular, we have the following,

Lemma 5.6. *If u is a nonnegative random variable, we have that,*

$$\mu = \int_0^1 \mathbb{P}(u \geq c) dc \tag{5}$$

interpreting the integral as a Lebesgue integral with respect to the Lebesgue measure.

An immediate corollary of Lemma 5.6 is that the least-squares estimator for the conditional expectation $\mathbb{E}[r_i|w]$ given by $\hat{r}_i(w)$ equals $\hat{\mu}_i(w)$ in the population setting under uniform random sampling of comparison prices. Another implication of this result is that we also obtain a data set of labeled comparisons expressing human feedback on the performance of the underlying generative process. One potential application is in the context of reinforcement learning from human feedback [Christiano et al., 2017]. In this setting, we would also allow comparisons between generated outputs and then train the value model using these pairwise comparisons as constraints.

5.2 Experiment with toxicity annotation

In some cases, allocations to users might result in negative utility, which would mean the mechanism pays users for reporting feedback. For instance, the employment of low-wage workers to enhance AI systems has brought about ethical concerns, such as the distressing impact on workers who review

harmful content [Steiger et al., 2021, Perrigo, 2023]. In particular, [Steiger et al., 2021] proposes preventing or reducing exposure as a potential technological intervention strategy. In this section, we present two experiments to evaluate our mechanism as an intervention strategy¹.

For the experiments, we use multi-label data from the Toxic Comment Classification Challenge on Kaggle, which contains a large number of Wikipedia comments labeled by human raters for toxic behavior, including categories such as toxic, severe-toxic, obscene, threat, insult, and identity-hate [Kag, 2018]. We assess the welfare regret of various allocation strategies by comparing them to the strategy that allocates resources to the agent with the highest expected welfare during each period. This concept is formally defined in definition 3.5. The expected utility is modeled with a linear model, as discussed in Section 5.1, and employs 30-dimensional PCA analysis to GloVe features for the data representation [Pennington et al., 2014]. Our experiment involves 10 agents and spans 5000 rounds of allocation. We assume that each agent possesses a fixed type sensitivity of which they are unaware, which determines their utility function. Every agent classifies examples sampled i.i.d from the dataset as toxic or non-toxic by reporting their relative utility from viewing the example. Moreover, when an agent encounters an example with a label matching their type sensitivity, they lose one unit of utility. For instance, an agent may be particularly sensitive to obscene examples.

We evaluate the welfare performance of three mechanisms: our method based on relative feedback, approaches that directly regress utility, and uniform assignment approaches. The direct utility regression methods, including [Devanur and Kakade, 2009, Babaioff et al., 2009, Nazerzadeh et al., 2013], vary in terms of payment structures, while learn-to-bid methods [Guo et al., 2022] assume workers estimate their own value and learn to bid accordingly. Uniform assignment approaches, on the other hand, make no attempt to intervene in the content allocation process and encompass most peer-prediction methods that pay based on conformity rather than utility, focusing on incentive compatibility issues [Prelec, 2004, Witkowski and Parkes, 2012, Cai et al., 2015].

Our experimental results in Figure 2a and 2b show that using an auction mechanism can significantly improve the welfare of the allocations given to agents over the peer-prediction method and performs favorably to the optimal allocation strategy. Moreover, our relative elicitation mechanism is competitive with the full information set but has the advantage of being simpler for users to report on. For example, while we make no further modifications to the calculation of welfare beyond what has already been discussed, some works assume there is a further cost to complicated elicitation strategies [Cai et al., 2015].

6 Limitations and future work

Our proposed approach has some limitations that warrant further exploration. While we propose an intervention method to improve the social welfare of workers doing toxic annotation, other aspects, such as negative psychological impacts and systemic issues around who does such work, are left unaddressed. While we provide a mechanism for the dynamic setting, we assume value evolves independently of other agents and the mechanism allocations, which prohibits us from studying collusion or adversarial scenarios. In particular, agents who make alias accounts could game the mechanism for free exploration allocations. Also, it is possible a randomized approach could improve upon the rates presented. Finally, we think further exploration as auction mechanism for selling digital goods in real-world settings is an important direction.

7 Conclusions and societal impact

In this paper, we have presented a novel approach to auctioning AI services that emphasizes user-friendly bidding, extends to multi-agent and contextual settings, offers simpler mechanisms, and improved bounds. Our approach has the potential for significant societal impact by facilitating a more efficient allocation of AI services and enabling a wider range of users to access these services without requiring them to have a deep understanding of their value. At the same time, we recognize that the collection of feedback or toxic annotation data may have negative externalities on the privacy and welfare of workers. Addressing the limitations and ethical concerns identified, we can move towards a future where AI services are more accessible, efficient, and ethically responsible, ultimately leading to a positive impact on society.

¹We provide the code for our experiments in the supplementary materials

Acknowledgments and Disclosure of Funding

I'd like to thank Ellen Vitercik for advising during the initial development of theory for the mechanism. I'd also like to thank Neil Band for suggesting I pursue simplifying the reporting method.

References

- Mar 2023. URL <https://www.jpmorgan.com/insights/research/generative-ai>.
- Katie Paul and Sheila Dang. Facebook owner meta announces tests of generative ai ads tool, May 2023. URL <https://www.reuters.com/technology/facebook-owner-meta-announces-tests-generative-ai-ads-tool-2023-05-11/>.
- Feb 2023. URL <https://www.statista.com/outlook/dmo/digital-advertising/worldwide>.
- David Karpf. Money will kill chatgpt's magic, Jan 2023. URL <https://www.theatlantic.com/technology/archive/2022/12/chatgpt-ai-chatbots-openai-cost-regulations/672539/>.
- Zhe Feng, Chara Podimata, and Vasilis Syrgkanis. Learning to bid without knowing your value. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 505–522, 2018.
- Wenshuo Guo, Michael Jordan, and Ellen Vitercik. No-regret learning in partially-informed auctions. In *International Conference on Machine Learning*, pages 8039–8055. PMLR, 2022.
- Hamid Nazerzadeh, Amin Saberi, and Rakesh Vohra. Dynamic pay-per-action mechanisms and applications to online advertising. *Operations Research*, 61(1):98–111, 2013.
- M-F Balcan, Avrim Blum, Jason D Hartline, and Yishay Mansour. Mechanism design via machine learning. In *46th Annual IEEE Symposium on Foundations of Computer Science (FOCS'05)*, pages 605–614. IEEE, 2005.
- Nikhil R Devanur and Sham M Kakade. The price of truthfulness for pay-per-click auctions. In *Proceedings of the 10th ACM conference on Electronic commerce*, pages 99–106, 2009.
- Moshe Babaioff, Yogeshwer Sharma, and Aleksandrs Slivkins. Characterizing truthful multi-armed bandit mechanisms. In *Proceedings of the 10th ACM conference on Electronic commerce*, pages 79–88, 2009.
- Moshe Babaioff, Robert D Kleinberg, and Aleksandrs Slivkins. Truthful mechanisms with implicit payment computation. *Journal of the ACM (JACM)*, 62(2):1–37, 2015.
- Nicole Immorlica, Kamal Jain, Mohammad Mahdian, and Kunal Talwar. Click fraud resistant methods for learning click-through rates. In *Internet and Network Economics: First International Workshop, WINE 2005, Hong Kong, China, December 15-17, 2005. Proceedings 1*, pages 34–45. Springer, 2005.
- Nikhil Agarwal, Susan Athey, and David Yang. Skewed bidding in pay-per-action auctions for online advertising. *American Economic Review*, 99(2):441–447, 2009.
- Mohammad Mahdian and Kerem Tomak. Pay-per-action model for on-line advertising. *International Journal of Electronic Commerce*, 13(2):113–128, 2008.
- Dirk Bergemann and Martin Pesendorfer. Information structures in optimal auctions. *Journal of economic theory*, 137(1):580–609, 2007.
- Alessandro Epasto, Andrés Muñoz Medina, Steven Avery, Yijian Bai, Robert Busa-Fekete, CJ Carey, Ya Gao, David Guthrie, Subham Ghosh, James Ioannidis, et al. Clustering for private interest-based advertising. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 2802–2810, 2021.
- Yang Cai, Constantinos Daskalakis, and Christos Papadimitriou. Optimum statistical estimation with strategic data sources. In *Conference on Learning Theory*, pages 280–296. PMLR, 2015.

- Siyu Chen, Jibang Wu, Yifan Wu, and Zhuoran Yang. Learning to incentivize information acquisition: Proper scoring rules meet principal-agent model. *arXiv preprint arXiv:2303.08613*, 2023.
- Shengjia Zhao and Stefano Ermon. Right decisions from wrong predictions: A mechanism design alternative to individual calibration. In *International Conference on Artificial Intelligence and Statistics*, pages 2683–2691. PMLR, 2021.
- Sai Praneeth Karimireddy, Wenshuo Guo, and Michael I Jordan. Mechanisms that incentivize data sharing in federated learning. *arXiv preprint arXiv:2207.04557*, 2022.
- Drazen Prelec. A bayesian truth serum for subjective data. *science*, 306(5695):462–466, 2004.
- Jens Witkowski and David Parkes. A robust bayesian truth serum for small populations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 26, pages 1492–1498, 2012.
- Billy Perrigo. Exclusive: Openai used kenyan workers on less than \$2 per hour to make chatgpt less toxic. *Last accessed*, 19, 2023.
- Carolyn M Burns, Jeff Morley, Richard Bradshaw, and José Domene. The emotional impact on and coping strategies employed by police teams investigating internet child exploitation. *Traumatology*, 14(2):20–31, 2008.
- Andrew Arsht and Daniel Etcovitch. The human cost of online content moderation. *Harvard Journal of Law and Technology*, 2018.
- Miriah Steiger, Timir J Bharucha, Sukrit Venkatagiri, Martin J Riedl, and Matthew Lease. The psychological well-being of content moderators: the emotional labor of commercial moderation and avenues for improving support. In *Proceedings of the 2021 CHI conference on human factors in computing systems*, pages 1–14, 2021.
- Alessandro Pavan, Ilya R Segal, and Juuso Toikka. Dynamic mechanism design: Incentive compatibility, profit maximization and information disclosure. *Profit Maximization and Information Disclosure (May 1, 2009)*, 2009.
- Sham M Kakade, Ilan Lobel, and Hamid Nazerzadeh. Optimal dynamic mechanism design and the virtual-pivot mechanism. *Operations Research*, 61(4):837–854, 2013.
- Roger B Myerson. Optimal auction design. *Mathematics of operations research*, 6(1):58–73, 1981.
- Martin Anthony, Peter L Bartlett, Peter L Bartlett, et al. *Neural network learning: Theoretical foundations*, volume 9. cambridge university press Cambridge, 1999.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
2018. URL <https://tinyurl.com/jigsaw-toxic-comment-challenge>.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- László Györfi, Michael Köhler, Adam Krzyżak, and Harro Walk. *A distribution-free theory of nonparametric regression*, volume 1. Springer, 2002.

A Summary of notation

Table 2: Summary of notation

Notation	Definition
$i \in [n]$	Agent i out of n
$t \in [H]$	Periods t in horizon H
$w_{it} \in \mathcal{W}$	Agent request
$o_{it} \in \mathcal{O}$	Produced good from request
$u_{it} \in [0, 1]$	Agent i utility at time t
μ_i	Expected value function of agent i
\mathcal{L}	Learning algorithm for μ_i
Δ_t	\mathcal{L} 's error - exploration rounds
η_t	Exploration frequency
$\hat{\mu}_{it}$	Estimate for μ_i with \mathcal{L}
$\hat{\gamma}_t$	Estimated second-price
S_{it}	Allocations of i up to time t
$x_{it} \in \{0, 1\}$	Allocation indicator
\hat{x}_{it}	Allocations with \mathcal{L}
$y_{it} \in \{0, 1\}$	Exploitation indicator
$r_{it} \in \{0, 1\}$	Report of agent i at time t
p_{it}	Payment of agent i at time t
c_{it}	Comparison price

B Omitted proofs

B.1 Proof of theorem 5.1

Theorem 5.1. *Suppose our mechanism \mathcal{M} estimates agent values bounded to the unit-interval using some learning algorithm \mathcal{L} . Suppose the expected error of this algorithm is monotone decreasing in the number of samples and that for all time,*

$$\sum_{t=1}^H \Delta_t = o\left(\sum_{t=1}^H \mathbb{E}[\eta(t)]\right) = o(H) \quad (4)$$

then the mechanism \mathcal{M} satisfies the following:

1. *Is asymptotically individually rational and asymptotically truthful*
2. *Is asymptotically welfare and revenue maximizing.*
3. *Compared to a second-price auction \mathcal{M} can obtain welfare and revenue regret $\tilde{O}(H^{2/3})$ if there is a learning algorithm \mathcal{L} for valuations with slow learning rate $\tilde{O}(H^{-1/2})$ and regret $\tilde{O}(H^{1/2})$ if there is an algorithm with a fast rate $\tilde{O}(H^{-1})$.*

We proceed by first showing that strategic reporting isn't particularly profitable compared to the profit from free allocations. Following this, we establish asymptotic truthfulness and individual rationality properties. Finally, we analyze the welfare and revenue regret of our mechanism.

Lemma 5.2. *A single agent providing misleading feedback can increase their expected utility up to time H by no more than $6\sum_{t \in [H]} \Delta_t$.*

Proof. Recall that Δ_t is the expected maximum error of the learning algorithm over all the agents $i \in [n]$ with respect to exploration round data. We denote the set of allocation times with $S := \{t \in [H] : \hat{x}_{it} = 1\}$. Recall there are numerous reporting strategies $r_{it} : (\mathcal{W} \times \mathcal{O}) \times [0, 1] \rightarrow \{0, 1\}$ each agent could use. Accordingly, when it is important to specify the reporting strategies used by the agents, we will adorn these definitions in some fashion with the relevant strategy L e.g., S_L , $\hat{x}_{it}(L)$, $\hat{\mu}_{it}(w_{it}|L)$.

We ignore exploration allocations since they are strategy independent and therefore cancel each other out. Still, our agent may try to deviate from being truthful under some reporting strategy T to some other reporting strategy L that induces a new set of allocation times S_L for the item to the agent. We fix $\hat{\mu}_{it}$ to represent the utility estimated from data collected under strategy T . By the tower property the expected utility for a given strategy satisfies,

$$\begin{aligned} U_i(L, H) &= \mathbb{E} \left[\sum_{t \in S_L} u_{it}(w_{it}) - \hat{\gamma}_t(L) \right] \\ &= \mathbb{E} \left[\sum_{t \in S_L} \mathbb{E} [u_{it}(w_{it}) - \hat{\gamma}_t(L) | \{(w_{ik}, c_{ik}, r_{ik})\}_{k \in [t]}] \right] \\ &= \mathbb{E} \left[\sum_{t \in S_L} \mu_i(w_{it}) - \hat{\gamma}_t(L) \right] \end{aligned}$$

Accordingly, the expected profit of such a strategy is given by,

$$\begin{aligned} \text{Profit}(L, H) &= U_i(L, H) - U_i(T, H) \\ &= \mathbb{E} \left[\sum_{t \in S_L} \mu_i(w_{it}) - \hat{\gamma}_t(L) \right] - \mathbb{E} \left[\sum_{t \in S_T} \mu_i(w_{it}) - \hat{\gamma}_t(T) \right] \\ &= \mathbb{E} \left[\sum_{t \in S_L \setminus S_T} \mu_i(w_{it}) - \hat{\gamma}_t(L) \right] - \mathbb{E} \left[\sum_{t \in S_T \setminus S_L} \mu_i(w_{it}) - \hat{\gamma}_t(T) \right] + \mathbb{E} \left[\sum_{t \in S_T \cap S_L} \hat{\gamma}_t(T) - \hat{\gamma}_t(L) \right] \end{aligned}$$

In the last line, observe that whenever $t \in S_T \setminus S_L$ we have that $\hat{\mu}_{it}(w_{it}|T) > \hat{\gamma}_t(T)$ and so,

$$\begin{aligned} \mu_i(w_{it}) - \hat{\gamma}_t(T) &= \mu_i(w_{it}) - \hat{\mu}_{it}(w_{it}|T) + \underbrace{\hat{\mu}_{it}(w_{it}|T) - \hat{\gamma}_t(T)}_{>0} \\ &\Rightarrow \mu_i(w_{it}) - \hat{\gamma}_t(T) \geq \mu_i(w_{it}) - \hat{\mu}_{it}(w_{it}|T), \quad t \in S_T \setminus S_L \end{aligned}$$

For $t \in S_L \setminus S_T$ let j be the agent who would receive the item if agent i were truthful. We have $\hat{\mu}_{jt}(w_{jt}|T) \geq \hat{\mu}_{it}(w_{it}|T)$. So we have,

$$\mu_i(w_{it}) - \hat{\gamma}_t(L) \leq \mu_i(w_{it}) - \hat{\mu}_{it}(w_{it}|T) + \hat{\mu}_{jt}(w_{jt}|T) - \hat{\gamma}_t(L) \quad t \in S_L \setminus S_T$$

We also have that $\hat{\mu}_{jt}(w_{jt}|L) \leq \hat{\gamma}_t(L)$ so we have,

$$\mu_i(w_{it}) - \hat{\gamma}_t(L) \leq (\mu_i(w_{it}) - \hat{\mu}_{it}(w_{it}|T)) + \max_{k \neq i} \{\hat{\mu}_{kt}(w_{kt}|T) - \hat{\mu}_{kt}(w_{kt}|L)\} \quad t \in S_L \setminus S_T$$

Finally, for $t \in S_T \cap S_L$ we let $j \neq i$ be the agent with the highest $\hat{\mu}_{jt}(w_{jt}|T)$. So we have,

$$\hat{\gamma}_t(T) - \hat{\gamma}_t(L) \leq \hat{\mu}_{jt}(w_{jt}|T) - \hat{\mu}_{jt}(w_{jt}|L) \leq \max_{k \neq i} \{\hat{\mu}_{kt}(w_{kt}|T) - \hat{\mu}_{kt}(w_{kt}|L)\}$$

After substitution, we end up with a simplified expression for the i^{th} agent's profit:

$$\begin{aligned} \text{Profit}(L, H) &\leq \mathbb{E} \left[\sum_{t \in S_L \setminus S_T} (\mu_i(w_{it}) - \hat{\mu}_{it}(w_{it}|T)) + \max_{k \neq i} \{\hat{\mu}_{kt}(w_{kt}|T) - \hat{\mu}_{kt}(w_{kt}|L)\} \right] \quad (6) \\ &\quad - \mathbb{E} \left[\sum_{t \in S_T \setminus S_L} \mu_i(w_{it}) - \hat{\mu}_{it}(w_{it}|T) \right] + \mathbb{E} \left[\sum_{t \in S_T \cap S_L} \max_{k \neq i} \{\hat{\mu}_{kt}(w_{kt}|T) - \hat{\mu}_{kt}(w_{kt}|L)\} \right] \end{aligned}$$

To proceed, we observe that:

$$\max_{k \neq i} \{\hat{\mu}_{kt}(w_k|T) - \hat{\mu}_{kt}(w_{kt}|L)\}$$

$$\begin{aligned}
&\leq \max_{k \neq i} |\mu_k(w_k) - \hat{\mu}_{kt}(w_{kt}|T)| + |\mu_k(w_{kt}) - \hat{\mu}_{kt}(w_{kt}|L)| \\
&\leq \max_{k \neq i} |\mu_k(w_{kt}) - \hat{\mu}_{kt}(w_{kt}|T)| + \max_{k \neq i} |\mu_k(w_{kt}) - \hat{\mu}_{kt}(w_{kt}|L)|
\end{aligned}$$

Recall the assumption that only the agent i is misreporting. Therefore, all the agents except i are truthful so each term can be bounded by Δ_t and so we end up with,

$$\begin{aligned}
\max_{k \neq i} \{ \hat{\mu}_{kt}(w_{kt}|T) - \hat{\mu}_{kt}(w_{kt}|L) \} &\leq 2 \max_k |\mu_k(w_{kt}) - \hat{\mu}_{kt}(w_{kt}|T)| \\
\Rightarrow \mathbb{E}[\max_{k \neq i} \{ \hat{\mu}_{kt}(w_{kt}|T) - \hat{\mu}_{kt}(w_{kt}|L) \}] &\leq 2\Delta_t
\end{aligned}$$

Here we are using the assumption that the learning algorithm is proper since Δ_t is the expected error over *only* the exploration rounds. Substituting the result into Eq. 6 we obtain:

$$\text{Profit}(L, H) \leq 6 \sum_{t=1}^H \Delta_t$$

as claimed, which completes the proof. □

Lemma 5.3. \mathcal{M} is asymptotically truthful and individually rational.

Proof. Notice that we have,

$$\begin{aligned}
U_i(H) &= \mathbb{E} \left[\sum_{t \in S_T} (\mu_i(w_{it}) - \hat{\gamma}_t) \right] + \frac{1}{n} \mathbb{E} \left[\sum_{t=1}^H \eta_t \mu_i(w_{it}) \right] \\
\hat{U}_i(H) - U_i(H) &\leq 6 \sum_{t \in [H]} \Delta_t = o \left(\mathbb{E} \left[\sum_t \eta_t \right] \right) = o(H)
\end{aligned}$$

The last step follows from Lemma 5.2 and since, by assumption, the free allocation rate strictly dominates the error rate of the estimation procedure.

We also have asymptotic individual rationality because we can bound the overcharges to the agent,

$$\mathbb{E} \left[\sum_{t \in S_T} \hat{\gamma}_t - \mu_i(w_{it}) \right] \leq \mathbb{E} \left[\sum_{t \in S_T} \hat{\mu}_{it}(w_{it}) - \mu_i(w_{it}) \right] \leq \mathbb{E} \left[\sum_t \Delta_t \right]$$

As we send the horizon to infinity, we see that these overcharges are bounded by the error rate and, by assumption, this term is dominated by the free allocations, so we have,

$$\mathbb{E} \left[\sum_{t \in S_T} \hat{\gamma}_t - \mu_i(w_{it}) \right] \leq \mathbb{E} \left[\sum_t \Delta_t \right] = o \left(\mathbb{E} \left[\sum_t \eta_t \right] \right) = o(H)$$

which establishes asymptotic ex-ante individual rationality. □

Lemma 5.4. \mathcal{M} is asymptotically welfare and revenue maximizing. Compared to a second-price auction \mathcal{M} can obtain welfare and revenue regret $\tilde{O}(H^{2/3})$ if there is a learning algorithm \mathcal{L} for valuations with slow learning rate $\tilde{O}(H^{-1/2})$ and regret $\tilde{O}(H^{1/2})$ if there is an algorithm with a fast rate $\tilde{O}(H^{-1})$.

Proof. We present the argument for revenue regret here. The argument for welfare regret goes similarly. The revenue minus free allocations satisfies,

$$\begin{aligned} \text{Revenue}(H) &= \mathbb{E} \left[\sum_{t=1}^H (1 - \eta_t) \hat{\gamma}_t \right] \geq \mathbb{E} \left[\sum_{t=1}^H (1 - \eta_t) (\gamma_t - \Delta_t) \right] \\ &\geq \mathbb{E} \left[\sum_{t=1}^H \gamma_t \right] - \mathbb{E} \left[\sum_{t=1}^H \eta_t \right] - \mathbb{E} \left[\sum_{t=1}^H \Delta_t \right] \end{aligned}$$

We conclude that,

$$\text{Regret}(H) \leq \mathbb{E} \left[\sum_{t=1}^H \eta_t \right] + \mathbb{E} \left[\sum_{t=1}^H \Delta_t \right]$$

We call $N_t = \sum_{k=1}^t \eta_k$ the expected number of exploration rounds and $N_t(i)$ the number of exploration rounds given to the i^{th} agent. Now we consider events when they are near their expectation. We use the following multiplicative Chernoff bound for sums of random variables bounded to the unit-interval²:

$$\mathbb{P}(N_t \leq (1 - \delta)\mu) \leq e^{-\delta^2\mu/2}$$

We define the event $\mathcal{A} = \{N_t \geq \frac{1}{2}\mathbb{E}[N_t]\}$ to capture the situation where the empirically observed number of exploration rounds is not too far below the expectation. Using the Chernoff bound we see,

$$\mathbb{P}(\neg\mathcal{A}) \leq e^{-\mathbb{E}[N_H]/8}$$

We also the event $\mathcal{B} = \{N_t(i) \geq \frac{1}{2n}N_t\}$ that the number of exploration rounds for a particular agent i is not too far from the expectation. Using the Chernoff bound once again we have,

$$\mathbb{P}(\neg\mathcal{B}) \leq e^{-N_H/8n} \Rightarrow \mathbb{P}(\neg\mathcal{B}|\mathcal{A}) \leq e^{-\mathbb{E}[N_H]/16n}$$

Conditioned on both of these events we have,

$$\mathbb{E}[\Delta_t] \leq \mathbb{E}[\Delta_t|\mathcal{A} \wedge \mathcal{B}] \cdot \mathbb{P}(\mathcal{A} \wedge \mathcal{B}) + 1 \cdot \mathbb{P}(\mathcal{A}) \cdot \mathbb{P}(\neg\mathcal{B}|\mathcal{A}) + 1 \cdot \mathbb{P}(\neg\mathcal{A}) \quad (7)$$

We will use the fact that conditional on $\mathcal{A} \wedge \mathcal{B}$ we have $\mathbb{E}[N_t(i)] \geq \frac{t\eta_t}{4n}$. There will be two cases to consider. The first corresponds to a slow learning algorithm, and the second to a fast learning algorithm.

Case 1: We have $\mathbb{E}[\Delta_t|N_t(i) = t] = O(\sqrt{\log(t)/t})$.

Without loss of generality, we can scale the bound and shift t by a constant so that we can assume $\mathbb{E}[\Delta_t|N_t(i) = t] \leq \sqrt{\log(t)/t}$. From this, we deduce that:

$$\begin{aligned} \mathbb{E}[\Delta_t|\mathcal{A} \wedge \mathcal{B}] &\leq \sqrt{\frac{4n \log\left(\frac{t\eta_t}{4n}\right)}{t\eta_t}} \leq \sqrt{\frac{4n \log(t/4)}{t\eta_t}} \\ \Rightarrow \mathbb{E}[\Delta_t] &\leq \sqrt{\frac{4n \log(t/4)}{t\eta_t}} + e^{-t\eta_t/8} + e^{-t\eta_t/16n} \end{aligned}$$

where in the last step we substitute into Eq. 7. For $t \geq 4e$ we have,

$$\sqrt{\frac{4n \log(t/4)}{t\eta_t}} \geq \sqrt{\frac{4n}{t\eta_t}}$$

Therefore,

$$\sup_{t \geq 4e} \frac{e^{-t\eta_t/16n}}{\sqrt{\frac{4n \log(t/4)}{t\eta_t}}} \leq \sup_{t \geq 4e} \frac{e^{-t\eta_t/16n}}{\sqrt{\frac{4n}{t\eta_t}}} = \sqrt{\frac{2}{e}}$$

²https://en.wikipedia.org/wiki/Chernoff_bound

$$\sup_{t \geq 4e} \frac{e^{-t\eta_t/8}}{\sqrt{\frac{4n \log(t/4)}{t\eta_t}}} \leq \sup_{t \geq 4e} \frac{e^{-t\eta_t/8}}{\sqrt{\frac{4n}{t\eta_t}}} = \sqrt{\frac{1}{n \cdot e}}$$

Subsequently, we have for $t \geq 4e$ that,

$$\mathbb{E}[\Delta_t] \leq c \cdot \sqrt{\frac{n \log(t)}{t\eta_t}}, \quad c = 2 \cdot (1 + \sqrt{2/e} + \sqrt{1/e})$$

Now we take $\eta_t = t^{-1/3} \cdot (n \log(t))^{(1+2\epsilon)/3}$ and then for $t \geq 4e$ we arrive at,

$$\begin{aligned} \mathbb{E}[\eta_t] &= t^{-1/3} \cdot (n \log(t))^{(1+2\epsilon)/3} \\ \mathbb{E}[\Delta_t] &\leq c \cdot \sqrt{\frac{n \log(t)}{t^{2/3} \cdot (n \log(t))^{(1+2\epsilon)/3}}} = c \cdot t^{-1/3} (n \log(t))^{(1-\epsilon)/3} \end{aligned}$$

We see from an integration that $\text{Regret}(H) = O(H^{2/3} \cdot \log(H)^{(1+2\epsilon)/3})$ which establishes the result. Since $\sum_t \mathbb{E}[\Delta_t] = o(\sum_t \mathbb{E}[\eta_t])$ this regret bound is consistent with the assumption in Theorem 5.1.

Case 2: We have $\mathbb{E}[\Delta_t | N_t(i) = t] = O(\log(t)/t)$.

Without loss of generality, we can scale the bound and shift t by a constant so that we can assume $\mathbb{E}[\Delta_t | N_t(i) = t] \leq \log(t)/t$. From this we deduce that,

$$\begin{aligned} \mathbb{E}[\Delta_t | \mathcal{A} \wedge \mathcal{B}] &\leq \frac{4n \log\left(\frac{t\eta_t}{4n}\right)}{t\eta_t} \leq \frac{4n \log(t/4)}{t\eta_t} \\ \Rightarrow \mathbb{E}[\Delta_t] &\leq \frac{4n \log(t/4)}{t\eta_t} + e^{-t\eta_t/8} + e^{-t\eta_t/16n} \end{aligned}$$

where in the last step we substitute into Eq. 7. For $t \geq 4e$ we have,

$$\frac{4n \log(t/4)}{t\eta_t} \geq \frac{4n}{t\eta_t}$$

Therefore,

$$\begin{aligned} \sup_{t \geq 4e} \frac{e^{-t\eta_t/16n}}{\frac{4n \log(t/4)}{t\eta_t}} &\leq \sup_{t \geq 4e} \frac{e^{-t\eta_t/16n}}{\frac{4n}{t\eta_t}} = \frac{4}{e} \\ \sup_{t \geq 4e} \frac{e^{-t\eta_t/8}}{\frac{4n \log(t/4)}{t\eta_t}} &\leq \sup_{t \geq 4e} \frac{e^{-t\eta_t/8}}{\frac{4n}{t\eta_t}} = \frac{2}{n \cdot e} \end{aligned}$$

Subsequently, we have for $t \geq 4e$ that:

$$\mathbb{E}[\Delta_t] \leq c \cdot \frac{n \log(t)}{t\eta_t}, \quad c = 4 \cdot (1 + 6/e)$$

Now we take $\eta_t = t^{-1/2} \cdot (n \log(t))^{(1+\epsilon)/2}$ and for $t \geq 4e$ we arrive at,

$$\begin{aligned} \mathbb{E}[\eta_t] &= t^{-1/2} \cdot (n \log(t))^{(1+\epsilon)/2} \\ \mathbb{E}[\Delta_t] &\leq c \cdot \frac{n \log(t)}{t^{1/2} \cdot (n \log(t))^{(1+\epsilon)/2}} = c \cdot t^{-1/2} (n \log(t))^{(1-\epsilon)/2} \end{aligned}$$

We see from an integration that $\text{Regret}(H) = O(H^{1/2} \cdot \log(H)^{(1+\epsilon)/2})$ which establishes the result. Since $\sum_t \mathbb{E}[\Delta_t] = o(\sum_t \mathbb{E}[\eta_t])$ this regret bound is consistent with the assumption in Theorem 5.1. \square

B.2 Proof of theorem 5.5

Theorem 5.5. *For an exploration rate $\eta_t = t^{-1/3} \cdot (n \log(t))^{(1+2\epsilon)/3}$ we have that $\mathcal{M}_\eta(\text{linear})$ satisfies the conditions of Theorem 5.1 and so is asymptotically individually rational, incentive compatible, and no-regret in revenue and welfare.*

To get the result, we need to show linear regression converges and that it converges to expected utility. Here we'll establish that the expectation of the reports does in fact equal expected utility. Later we will establish convergence of the regression to expected utility.

Lemma 5.6. *If u is a nonnegative random variable, we have that,*

$$\mu = \int_0^1 \mathbb{P}(u \geq c) dc \quad (5)$$

interpreting the integral as a Lebesgue integral with respect to the Lebesgue measure.

Proof. Given any measurable function $u : \Omega \rightarrow [0, \infty]$ from a sample space Ω to the nonnegative real numbers, there is a sequence of nonnegative simple functions u_t increasing pointwise to u . This means we can construct a sequence of nonnegative simple random variables u_t , increasing to u . To proceed, we look at the range set R_t of the u_t , which will be of finite cardinality. Moreover, each element in the range will be separated with some margin $\epsilon > 0$. Index these sets with $R_t(k)$ in order taking $R_t(0)$ to be a greatest lower bound. We have,

$$\begin{aligned} \sum_{k=R_t(0)}^{|R_t|} P(u_t \geq R_t(k)) &= \sum_{k=R_t(0)}^{|R_t|} \sum_{m=k}^{|R_t|} P(m + \epsilon > u_t \geq m) \\ &= \sum_{m=R_t(0)}^{|R_t|} \sum_{k=1}^m P(m + \epsilon > u_t \geq m) \\ &= \sum_{m=R_t(0)}^{|R_t|} m P(m + \epsilon > u_t \geq m) = \sum_{m=R_t(0)}^{|R_t|} \int_{[m, m+\epsilon)} m d\nu \\ &= \sum_{m=R_t(0)}^{|R_t|} \int_{[m, m+\epsilon)} u_t d\nu = \int u_t d\nu = \mathbb{E}[u_t] \end{aligned}$$

Our construction also implies $\mathbb{P}(u_t \geq z)$ increases to $\mathbb{P}(u \geq z)$ monotonically with respect to the sequence. It is not hard to see that this sequence of functions is measurable. Therefore, by the monotone convergence theorem,

$$\mathbb{E}[u] = \lim_{n \rightarrow \infty} \int_0^\infty \mathbb{P}(u_t \geq z) dz = \int_0^\infty \mathbb{P}(u \geq z) dz$$

which concludes the proof. □

For the main result. It is relatively well-known that well-specified linear regression converges in its predictions.

Lemma B.1. *([Györfi et al., 2002]) Let $\hat{\mu}_t$ be an empirical least squares estimator for a linear valuation function μ , which may depend on the input data w_1, \dots, w_t . Denote the $L_2(\nu)$ norm with respect to the probability measure ν of the data as $\|\cdot\|$. Then,*

$$\mathbb{E}[\|\hat{\mu}_t - \mu\|^2] \leq c \cdot \frac{(\log(t) + 1) \cdot d}{t}$$

for some constant c .

Since ν is a probability measure, by the monotone property for $L_p(\nu)$ norms we have that,

$$\mathbb{E}[\|\hat{\mu}_t - \mu\|] \leq \sqrt{\mathbb{E}[\|\hat{\mu}_t - \mu\|^2]} \leq \sqrt{c \cdot \frac{(\log(t) + 1) \cdot d}{t}}$$

Recall that Δ_t is the expected maximum error of the learning algorithm over all the agents $i \in [n]$ with respect to exploration round data. A simple bound works for our purposes,

$$\Delta_t \leq n \cdot \mathbb{E}[\|\hat{\mu}_t - \mu\|] = O\left(\sqrt{\frac{\log(t)}{t}}\right)$$

so with $\eta_t = t^{-1/3} \cdot (n \log(t))^{(1+2\epsilon)/3}$ we satisfy the necessary condition to invoke Theorem 5.1.