# Mechanism Design for Large Language Models[*]

Paul Dütting[†]     Vahab Mirrokni[†]     Renato Paes Leme[†]     Haifeng Xu[‡]

Song Zuo[†]

## Abstract

We investigate auction mechanisms for AI-generated content, focusing on applications like ad creative generation. In our model, agents' preferences over stochastically generated content are encoded as large language models (LLMs). We propose an auction format that operates on a token-by-token basis, and allows LLM agents to influence content creation through single dimensional bids. We formulate two desirable incentive properties and prove their equivalence to a monotonicity condition on output aggregation. This equivalence enables a second-price rule design, even absent explicit agent valuation functions. Our design is supported by demonstrations on a publicly available LLM.

## 1   Introduction

In the current web ecosystem, auctions are the primary mechanism used to decide which ads (and commercial content more broadly) are displayed to users (Edelman et al., 2007; Varian, 2007). In these auctions, advertisers bid for the opportunity to display their ad creatives alongside organic content. Many of the web formats such as text, banners, video, apps, ... have their own subtleties which led to the development of new auction tools to handle them. Our goal in this paper is to investigate auction mechanisms to support the emerging format of AI-generated content. More

specifically, we explore the use of auctions as a tool for influencing the output of large language models (LLMs) (e.g., Brown et al., 2020).

We consider a situation where a certain space in the web (which could be a part of a webpage, an UI element of an AI-chatbot, the dialog of a certain character in a video or a game, etc.) is designated for commercial content and different advertisers can bid to influence the content in that space. Each advertiser has an LLM that can generate content for that space, and is willing to pay a certain amount of money for the right to have their content displayed. A simple design is to collect bids from advertisers and let the highest bidder choose whatever content they wish to publish in that space. While simple, this design does not exploit the flexibility of LLMs which is to combine different concepts in a creative way.

Consider this example. First, we ask an LLM to produce different ads for the fictitious Stingray Resort and the equally fictitious Maui Airlines:

- *"Experience the magic of Hawaii at Stingray Resort, where stunning views, luxurious accommodations, and endless activities await. Book your stay today and create unforgettable memories in the heart of paradise."*

- *"Fly to Hawaii with Maui Airlines and experience the beauty of the Aloha State. We offer affordable flights to all the major islands, so you can start your Hawaiian vacation sooner. Book your flight today and let the island spirit take over!"*

For that use case, however, the LLM is flexible enough to produce a joint ad for both:

- *"Fly to paradise with Maui Airlines and experience the magic of Hawaii at Stingray Resort. Stunning views, luxurious accommodations, and endless activities await. Book your dream vacation today and create unforgettable memories."*

One can envision an auction mechanism that allows both Stingray Resort and Maui Airlines to submit their LLMs and bids, with these inputs determining their prominence in the final outcome.[1]

## 1.1 Unique Challenges

LLMs (Brown et al., 2020; Google et al., 2023; Thoppilan et al., 2022) are an emerging technology with new and unconventional aspects, many of which have direct implications to auction design (e.g., how preferences are represented/expressed). Our goal is to identify some of the key challenges and take a first step in designing mechanisms to address them:

---

[1]While this work's main focus is to create ad creatives that merge content from different advertisers, our designed auction mechanism for merging LLM outputs could also be used in other contexts.

- **Modelling and Expressing Preferences.** Auction theory typically models preferences via value functions that assign a value to each outcome. LLMs, however, as *generative models*, do not directly assign values. Instead, they succinctly encode preferences over outcomes within a stateless neural network model that predicts continuation probabilities.

- **Necessity of Randomization.** LLMs crucially rely on randomization. When forced to output tokens deterministically, LLMs often have a worse performance compared to situations that sample from a distribution (see, e.g., (Holtzman et al., 2019), for a performance comparison of different decoding strategies). Therefore, an auction that aggregates LLM outputs should preferably also output distributions.

- **Technical Compatibility.** Auction solutions should be compatible with current LLM technology, utilizing readily available information and integrating seamlessly. Ideally, the allocation and payments should be obtained from simple manipulations of the LLM outputs.

- **Computational Efficiency.** LLM models are expensive to query, so the auction computation should not add too much overhead. In particular, auctions should not increase the number of calls to inference the models beyond the minimum necessary.

## 1.2 Our Contributions

**The Token Auction Model.** Our first contribution is a formalism ("The Token Auction Model") for studying this problem. *Tokens* are the units making up sentences and paragraphs.[2] Examples of tokens include (sub-)words, symbols, numbers, and special tokens indicating the beginning and ending of the text. In particular, any piece of text (potentially incomplete) can be represented as an array of tokens, and any array of tokens also encodes a piece of text.

One salient feature of the state-of-the-art LLMs is that they are stateless, i.e., they maintain no internal memory or state. Instead, they simply map a prefix string to a distribution over the next token. The output is then created in an autoregressive manner. Given an input prompt, the output is generated by repeatedly feeding the current sequence of tokens into the LLM, sampling a continuation token, and appending it to the sequence of tokens.

The proposed *token auction* operates on a token-by-token basis, and serves to aggregate several LLMs to generate a joint output. We assume the designer has access to algorithmic LLM agents represented by their respective text generation functions (the functions that map a sequence of

---

[2]More generally, one can consider tokens forming parts of images (Ramesh et al., 2021; Yu et al., 2022) and videos (Sun et al., 2019). For the purpose of this paper, we will stick with text generation.

tokens to a distribution over the next token). In addition, we allow each LLM agent to submit a single dimensional bid. The auction output will be an aggregated distribution together with a payment rule that defines payments for each agent.[3]

This approach may seem counterintuitive initially, as advertisers typically focus on the final generated text rather than individual word choices. This seems to suggest a dynamic planning of the generated token sequence. However, existing LLMs do not reason about full pieces of text, nor do they plan ahead; instead, their preferences are expressed as desired distributions over merely the next token. In other terms, we can think of an LLM as a succinct distillation of an agent's complex combinatorial preferences over sequences of tokens into a generative token-by-token model.[4]

The problem of aggregating LLMs forces the designer to understand the preferences of the agents away from the distilled LLM. This appears to be a very difficult problem. Specifically, we believe it is implausible or at least impractical to assume an individual agent can meaningfully manipulate the distribution over tokens at any given stage, to direct the produced text to a more preferred one. Our auction formulation seeks to strike a balance: By truthfully revealing the LLM to the designer, the agent gives the auction mechanism a hint as to what their preferred distribution is. The bids, in turn, can be used to tradeoff between agents, and in particular help the designer determine their relative weights.

**Simple and Robust Token Auctions.**   Motivated by the challenges in modeling agents' preferences over generated distributions, we take a robust design approach aiming for token auctions that provide desirable incentive properties, while imposing minimal assumptions on the agents' preferences over distributions.

Specifically, we model agents' preferences as entailing partial orders over distributions. Based on this partial preference order[5], we formulate two desirable incentive properties, which we consider minimal requirements:

- *Payment monotonicity*: Given two different bids by the same agent, a final distribution is closer to the desired distribution if and only if the payment is higher.

- *Consistent aggregation:* If for two different bids of the same agent, the final distribution is closer to the preferred distribution for some bids of the other agents, then it should be so for all bids of the other agents.

---

[3]See our discussion later this section on the rationale of the indirect mechanism formulation.

[4]See our discussion in Section 4, and Propositions 4.1 and 4.3 for additional support for the stateless approach.

[5]Partial orders are more general than total orders, and hence our key results (such as Lemma 3.7, Lemma 3.6, and Theorem 3.5) apply to any complete preference order model.

We show that any mechanism with these two properties is *strategically equivalent* to a mechanism that satisfies a monotonicity requirement on the distribution aggregation function.

We then investigate whether it is possible to equip such distribution aggregation functions with payment rules that satisfy additional incentive properties. Specifically, we investigate whether such aggregation rules admit an analogue of the *second-price payment rule*. In the single-item second-price (or Vickrey) auction (Vickrey, 1961), the payment corresponds to the critical bid where an agent transitions from losing to winning. To port this notion to our setting, we show that under robust preferences (see Definition 2.1) any monotone aggregation rule can be written as a distribution over deterministic allocations from bids to tokens such that there is a critical bid where the allocation transitions from a less preferred to a more preferred token. Such a critical bid then serves as a natural candidate for a payment rule. This hence leads to an analogue of the second-price auction for our token auction model that only requires ordinal preferences. The resulting class of auctions is applicable whenever the agent valuations are compatible with the partial order, and provides robust incentives for all of these.

**Designing Aggregation Functions.** We then move to designing concrete aggregation functions. Our approach considers aggregated loss functions inspired by state-of-the-art LLM training, and derives optimal distribution aggregation functions that minimizes such aggregated loss functions.

We focus on specific forms of aggregated loss functions based on KL-divergence, a commonly used loss function in LLMs. We consider two natural formulations inspired by current LLM training, and show that the corresponding optimal aggregation rules are the weighted (log-space) convex combination of the target distributions from all participants.

The linear and log-linear aggregation rules we identify have different pros and cons. Both share the advantage that they are optimal for the respective aggregated loss functions. The linear rule turns out to be monotone with respect to robust preferences, and is therefore compatible with the robust incentives approach. However, the log-linear rule is not.

**Demonstration.** We conclude with demonstrations to support our token auction formulation, obtained by prompt-tuning of a publicly available LLM. A two-advertiser demonstrative example is considered, under both the linear and log-linear aggregation rules. We show how the combined output varies as a function of $\lambda = b_1/(b_1+b_2)$, where $b_1$ and $b_2$ are the advertisers' bids. Both approaches lead to meaningful and interpretable texts that smoothly transition from favoring one to favoring another advertiser, with a joint ad produced for intermediate values of $\lambda$.

**Discussion/Design Choices.** An alternative to our approach of designing an *indirect mechanism* would be to aim for a *direct mechanism*. Such a mechanism, instead of asking agents for a scalar bid along with query access to the agents' LLMs, would elicit the agents' full preferences directly. However, this appears unrealistic in our new domain due to multiple reasons: (1) Allocation outcomes in our setting are a high-dimensional distribution, whereas a classic mechanism's allocation is typically a subset of items, and often a single item in tractable setups. (2) While it is reasonable in the classic setup to elicit a valuation for an item or a subset of items, it does not appear realistic to elicit a high-dimensional utility function over all possible token distributions. (3) Eliciting full preferences over any token distribution would require solving a problem that is strictly harder than what current LLMs are trained to do (namely, merely output the most preferred distribution). This level of complexity might go beyond current technological capabilities and would likely be computationally inefficient.

## 1.3 Additional Related Work

To the best of our knowledge, the exact research question and our approaches in this work have not been previously studied. However, our work is indeed connected to a few lines of research.

**Related LLM Research.** Our work shares some similarities with the literature on fine-tuning LLMs, with reinforcement learning from human feedback (RLHF) as a representative approach (Bai et al., 2022; Bakker et al., 2022; Ouyang et al., 2022; Wei et al., 2021). At a high level, fine-tuning and RLHF seek to align a generally pre-trained LLM with certain desirable behaviors. This is in spirit analogous to our goal of designing LLMs to better align with a group of agents' overall preferences. However, our research challenges and methods are both different from those in the fine-tuning literature. Specifically, fine-tuning refines the underlying model's parameters whereas our approach is one-layer up and directly aggregates the token distributions from multiple models. The main challenge we address is the potential incentive misalignment while eliciting LLM agents' preferences, whereas human labelers or other models that generate reward feedback for RLHF are assumed to be genuine and do not misrepresent their own preferences.

The literature on in-context learning (Brown et al., 2020; Wei et al., 2022, 2023) is similar to us in the sense that this approach also does not change the model parameters. A main difference to our work is that this literature seeks to influence token distributions by conditioning on better-generated prefix contexts, whereas we directly aggregate distributions from multiple LLM agents.

**Connections in Mechanism Design.** Our work is related to the literature on (combinatorial) public projects (Dughmi, 2011; Papadimitriou et al., 2008). The connection is that one can view the output of the aggregated LLM in our situation as a public project that benefits the agents to different degrees. Similar to these earlier studies, a core challenge in our problem is to elicit preferences about the public project from unknown agents. However, the design problem in our case is fundamentally different — we choose a high-dimensional distribution from an $\mathbb{R}^T$ space with only partial knowledge about agents' preferences, whereas previous work has focused on the problem of choosing from a discrete (often exponentially large) set with clear agent valuation functions (Dughmi, 2011; Papadimitriou et al., 2008).

Another related stream of work includes (Freeman et al., 2019; Goel et al., 2019), which studies the problem of truthfully aggregating budget proposals. Their mechanisms output a distribution over budgets that best serves the population, just like our mechanisms output distributions over tokens. However, the objectives and techniques between our work and theirs are both different. First, their problem is mechanism design without money, whereas our problem has monetary transfers involved. A direct consequence of this first difference is that their mechanisms will treat every participant with equal weight, whereas the weights of our participants are determined by their bids. Second, the research on truthful budget proposal aggregations typically assumes explicit valuation functions (e.g., $l_1$ distance between preferred and output distributions), under which the VCG mechanism is truthful. Their main research question hence is to study additional properties of the mechanisms such as Pareto-efficiency and certain fairness properties (Freeman et al., 2019). Assuming such an explicit valuation function does not appear realistic in our problem, so our core research question is to design robust mechanisms that enjoy good incentive properties simultaneously for a broad range of valuation functions.

From this perspective, our work also bears some similarity to the rich literature on robust mechanism design. Most of this literature still assume existence of value functions with uncertainty modeled by Bayesian beliefs or in a max-min sense (Bergemann and Morris, 2005, 2012; Carroll, 2015; Dütting et al., 2019; Roughgarden and Talgam-Cohen, 2016). However, assuming such a valuation function over tokens or their distributions does not appear realistic in creatives generation, thus our model is more similar to a worst-case style consideration during which we only assume partial ("obvious") preferences.

**Follow-Up Work.** Several papers follow-up on our work, by studying mechanism design problems for LLMs. Dubey et al. (2024) consider bidders that bid for placement of their content within a summary generated by a large language model. Soumalias et al. (2024) design a truthful mech-

anism that generates several samples from a reference LLM, and incentivizes bidders to truthfully reveal their preferences. Mordo et al. (2024) consider sponsored question answering, in which an organic answer to a search query is fused with an ad to create a sponsored answer, and advertisers bid on the sponsored answers.

# 2  Preliminaries

In this section, we first provide an abstraction of typical generative models and then introduce the basic formalism of the mechanism design problem we study. For concreteness we adopt a terminology that suits the important LLM use case where the creative is text.

## 2.1  Abstraction of Large Language Models

*Large language models* (LLMs) (Brown et al., 2020; Google et al., 2023; Thoppilan et al., 2022) can be abstracted as functions mapping from a partial sentence to the distribution of the next *token* that extends the partial sentence.

Formally, let $T$ be the set of tokens and $\Delta(T)$ be the set of distributions over $T$. Let $T^* = T \cup T^2 \cup \cdots \cup T^K$ denote the set of sequences of tokens, where $K$ is the maximum sequence length that the LLM can handle. Each LLM is modeled as a function $f : T^* \to \Delta(T)$ that maps any sequence of tokens to a distribution over the next token.

**Autoregressive Text Generation.**  A *prompt* is an initial set of tokens $s_0 \in T^*$ provided with instructions of what text to generate. An LLM produces a text in response to the prompt by sampling a token $\tau_1 \sim f(s_0)$ and constructing $s_1 = s_0 \oplus \tau_1$ (where $\oplus$ is the operation to append a token to an array). We then repeat the process of $\tau_k \sim f(s_k)$ and $s_k = s_{k-1} \oplus \tau_k$ until a special *end-of-sentence* token is sampled. If at some point the sequence of tokens becomes too long (larger than $K$) we trim $s_k$ to its length-$K$ suffix.

Note that LLMs are stateless: They lack internal memory beyond the generated token sequence, and each token is sampled independently.

**Training of LLMs.**  An LLM $f$ is parameterized by a neural network structure $M$ and a set of weights $W$. The weights are often obtained by three stages of optimization (see first three rows in Table 1). The initial stage is very computationally intensive but task independent. Subsequent stages are less costly and their goal is to adapt the general purposed model obtained in the first stage to more specific tasks. Each of the stages usually minimizes a different loss function over a

| Training/Learning stages | Data | Cost | Goal |
|---|---|---|---|
| Pre-training[*] | General texts from web, books, etc | Very high | A common baseline shared across downstream tasks |
| Instruction fine-tuning[†] | Task specific data | Medium | Optimize the behavior for specific tasks |
| RLHF[‡] | Human evaluations | Medium | Security control, reducing harmful behavior, etc |
| In-context few-shot learning | Carefully designed prompts as inputs | Very low | Effectively influence the behavior in real-time |

[*] Brown et al. (2020); Google et al. (2023)  [†] Wei et al. (2021)  [‡] Ouyang et al. (2022)

Table 1: Common Training Stages of LLMs.

different dataset. The details of the training process are not particularly relevant to our discussion, though a more detailed discussion can be found in Section 4.1. We will note, however, that some of the mechanisms we discuss for combining the inputs of different LLMs resemble the functional forms used in the reinforcement learning and fine-tuning steps.

## 2.2   Token Auctions for LLMs

We now formalize the mechanism design problem of combining the output of different LLM-represented algorithmic agents. As discussed in the introduction, we will design an auction to act on the token-by-token generation stage. Our goal is to keep the auction technically aligned with the state-of-the-art LLM systems.

**Robust Modeling of LLM Agents' Preferences.**    A key challenge in designing mechanisms for LLM agents is comparing their "utilities" over different output distributions. To illustrate, suppose an LLM agent's preferred distribution over two tokens is $p = (0.6, 0.4)$, and consider two possible generated distribution outcomes: $q_1 = (0.5, 0.5)$ and $q_2 = (0.8, 0.2)$. Between $q_1$ and $q_2$, it is unclear which one this LLM agent would prefer since while $q_2$ appears more distant from $p$ than $q_1$, it has a higher probability on the first token which appears more preferably by the LLM's initial distribution $p$.

Despite this incomparability between $q_1$ and $q_2$, it does appear clear that $q_2$ would be less preferred by the LLM than $q_3 = (0.7, 0.3)$. This is because $q_3$ deviates from $p$ along the same directions as $q_2$ for each entry (i.e., both increase or both decrease), but deviates less in terms of the absolute value of deviation.

The above observation illustrates that while it is difficult to model LLM agents' complete preferences over all the generated distributions, it seems natural to assume certain *partial order* over the distributions. This motivates us to consider a *robust* modeling of LLM agents' preferences through general partial orders, and sometimes, the following specific notion of *robust preferences*.

**Definition 2.1** (Robust Preferences over Distributions)**.** Consider any LLM agent $i$ with preferred distribution $p_i$, and any two aggregation distribution $q, q' \in \Delta(T)$. We say $q$ is (weakly) *robustly*

9

*preferred* over $q'$ by agent $i$, or formally, $q \succeq_i q'$, if

$$\forall t \in T, \qquad |q(t) - p_i(t)| \le |q'(t) - p_i(t)| \qquad (1)$$

$$\text{and} \quad (q(t) - p_i(t))(q'(t) - p_i(t)) \ge 0. \qquad (2)$$

Moreover, if $q \ne q'$, then $q$ is strictly preferred over $q'$ by $i$, i.e., $q \succ_i q'$.

In other words, $q$ is robustly preferred by $i$ over $q'$ when (1) the deviation of $q$ from $p_i$ is smaller than the deviation of $q'$ from $p_i$ for every entry; and (2) these deviations are along the same direction for every entry. Note that Definition 2.1 only specifies a partial ordering among aggregated distributions. Thus it is possible that two distributions are not comparable, i.e., $q \nsucceq_i q'$ and $q' \nsucceq_i q$.

**Token Auctions.** Our goal is to design simple, practical auction mechanisms that work well under minimal assumptions about the agents' private preferences. Specifically, we seek to design *token auction mechanisms* $\mathcal{M} = \langle q, z \rangle$, where $q$ is a distribution aggregation function and $z$ is a payment function. A token auction mechanism operates on a token-by-token basis, and lets $n$ algorithmic LLM agents influence the output distribution and payments through scalar bids. The bid profile of all agents is denoted as $\boldsymbol{b} = (b_1, \ldots, b_n) \in \mathbb{Q}_+^n$ where any bid $b_i \in \mathbb{Q}_+$ is a non-negative rational number. We assume that the initial prompt $s_0 \in T^*$, and the text aggregation functions $f_1, \ldots, f_n$ of the $n$ LLM agents are publicly known.

*Distribution Aggregation Function.* This is the first ingredient to a token auction mechanism. A distribution aggregation function $q$ takes as input a vector of bids $\boldsymbol{b} \in \mathbb{Q}_+^n$ and $n$ distributions $\boldsymbol{p} \in \Delta(T)^n$ and maps these to a distribution over tokens:

$$\text{aggregation function:} \quad q : \mathbb{Q}_+^n \times \Delta(T)^n \to \Delta(T).$$

For fixed bids, a distribution aggregation function can be used in the same way as a text aggregation function. Namely, starting from the initial prompt $s_0 \in T^*$, we can repeatedly sample $\tau_k$ from distribution $q_k = q((b_1, \ldots, b_n), (f_1(s_{k-1}), \ldots, f_n(s_{k-1})))$ for each $k \ge 1$ to generate $s_k = s_{k-1} \oplus \tau_k$. Note the alignment with LLMs, which already produce the distributions $f_i(s_{k-1})$ for $i \in [n]$. No additional calls to the LLMs are needed.

*Payment Function.* In addition to the distribution aggregation function, we seek to design payment functions. Here, we want to operate on a token-by-token basis and seek a stage-independent design akin to the stage independence of state-of-the-art LLMs' token generation. Formally, for

each agent $i$, we aim to define a

$$\text{pricing function:} \quad \zeta_i : \mathbb{Q}_+^n \times \Delta(T)^n \times T \to \mathbb{R},$$

with the interpretation that for bids $\boldsymbol{b} \in \mathbb{Q}_+^n$, distributions $\boldsymbol{p} \in \Delta(T)^n$, and token $t \sim q(\boldsymbol{b}, \boldsymbol{p})$, the payment from agent $i$ is $\zeta_i(\boldsymbol{b}, \boldsymbol{p}, t)$. These pricing functions naturally lead to expected payments by taking expectations over tokens. Namely, for each agent $i$, we define

$$\text{payment function:} \quad z_i : \mathbb{Q}_+^n \times \Delta(T)^n \to \mathbb{R},$$

as the function that takes as input a vector of bids $\boldsymbol{b} \in \mathbb{Q}_+^n$ and distributions $\boldsymbol{p} \in \Delta(T)^n$, and maps these to $z_i(\boldsymbol{b}, \boldsymbol{p}) = \mathbb{E}_{t \sim q(\boldsymbol{b}, \boldsymbol{p})}[\zeta_i(\boldsymbol{b}, \boldsymbol{p}, t)]$.

**Discussion.**    Token auctions provide a suitable abstraction for analyzing the strategic aspects of LLM aggregation. Fully expressing agent preferences over all generated content is impractical. Representing agents as LLMs is a plausible approach, as LLMs distill preferences into token distributions. Thus, auctioning tokens based on LLM-expressed preferences is a natural mechanism.

At the same time, the detailed functioning of LLMs remains rather opaque, and it seems implausible that agents could meaningfully misreport the outcome distributions of their LLMs in order to achieve a more desirable aggregated output. Our auction formulation offers a middle ground. We assume the designer has access to the LLMs, but let the agents influence the aggregation process through a single dimensional bid.

## 3   Incentives in Token Auctions

In this section, we examine the strategic properties of token auctions. Our goal is robust incentive properties that rely on as few assumptions about the agents' preferences as possible. We first formulate two natural properties that any reasonable mechanism should satisfy, and show that they are essentially equivalent to a monotonicity requirement on the distribution aggregation function. We then show that, when agents have robust preferences (see Definition 2.1), for any such monotone distribution aggregation function, it is possible to define a natural second-price payment rule.

## 3.1 Desirable Incentive Properties

We begin by formulating two conditions that reasonable token auction mechanisms should satisfy with respect to general partial orders $\succeq_i$. The first is a *monotonicity* condition on the payment function. It requires that agents' pay increases if and only if they obtain better distributions.

**Definition 3.1** (Payment Monotonicity)**.** Mechanism $\mathcal{M} = \langle q, z \rangle$ satisfies *payment monotonicity*, if for all $\boldsymbol{p}, \boldsymbol{b}_{-i}, b_i, b_i'$ we have

$$z_i(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) \geq z_i(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p}) \iff q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) \succeq_i q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p}).$$

It is a natural incentive constraint because if the payment function is not monotone, then bidders will likely manipulate their bids in order to induce better distribution with lower payment. The following is a useful implication of payment monotonicity. Given $\boldsymbol{p}, \boldsymbol{b}_{-i}$, if $i$'s two bids $b_i, b_i'$ lead to the same aggregated distribution, i.e., $q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) = q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p})$, then as a consequence of payment monotonicity the payment must also be the same.

The second incentive constraint is about the consistency of the aggregation function. Intuitively, the relative preference between two bids should not be influenced by the bids of other agents.

**Definition 3.2** (Consistent Aggregation)**.** The distribution aggregation function $q(\boldsymbol{b}, \boldsymbol{p})$ is said to be *consistent* if it admits consistent ordering across all $\boldsymbol{b}_{-i}$. Formally, if $q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) \succ_i q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p})$ for some $\boldsymbol{b}_{-i}$, then for all $\boldsymbol{b}_{-i}'$, $q(b_i, \boldsymbol{b}_{-i}', \boldsymbol{p}) \succeq_i q(b_i', \boldsymbol{b}_{-i}', \boldsymbol{p})$.

Similar to payment monotonicity, this requirement of consistent aggregation is imposed to avoid bidders' concerns that the same bid can lead to better or worse distributions, depending on the opponents' bids. To rule out such behavior, consistency requires that whenever bids $b_i, b_i'$ are such that $b_i$ is strictly preferred over $b_i'$ for some opposing bids $\boldsymbol{b}_{-i}$, then it better be that $b_i$ is weakly preferred over $b_i'$ for all opposing bids $\boldsymbol{b}_{-i}'$.

## 3.2 Monotone Aggregation Functions

Next we show a "revelation principle" type of result, stating that if one is interested in mechanisms satisfying the desirable incentive properties stated above (Definition 3.1 and Definition 3.2), then one can without loss of generality focus on *monotone aggregation functions* as captured in the following definition for any partial order $\succeq_i$.

**Definition 3.3** (Monotone Aggregation Function). The distribution aggregation function $q(\boldsymbol{b}, \boldsymbol{p})$ is called *monotone* if any higher bid from any agent $i$ leads to a more preferred aggregated distribution for $i$. Formally, for all $\boldsymbol{p}$, $\boldsymbol{b}_{-i}$ and $b_i \geq b'_i$:

$$q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) \succeq_i q(b'_i, \boldsymbol{b}_{-i}, \boldsymbol{p}).$$

We are now ready to state our main finding in this subsection with the following definition of *strategic equivalence* from one mechanism $\mathcal{M}$ to another $\tilde{\mathcal{M}}$. In words, the aggregated distribution outcome and all agents' payments will be the same under mechanism $\mathcal{M}$ and $\tilde{\mathcal{M}}$ after each agent $i$ applies some strategy mapping $\pi_i$. Formally,

**Definition 3.4** (Strategic Equivalence). A mechanism $\mathcal{M} = \langle q, z \rangle$ is *strategically equivalent* to another mechanism $\tilde{\mathcal{M}} = \langle \tilde{q}, \tilde{z} \rangle$, if $\forall \boldsymbol{p} \in \Delta(T)^n$, there exists a profile $\pi$ of strategy mappings with $\pi_i : \mathbb{Q}_+ \to \mathbb{Q}_+$ for every agent $i$ (i.e., $\pi(\boldsymbol{b}) = (\pi_1(b_1), \ldots, \pi_n(b_n))$), such that $\forall \boldsymbol{b} \in \mathbb{Q}_+^n$, $q(\boldsymbol{b}, \boldsymbol{p}) = \tilde{q}(\pi(\boldsymbol{b}), \boldsymbol{p})$ and $z(\boldsymbol{b}, \boldsymbol{p}) = \tilde{z}(\pi(\boldsymbol{b}), \boldsymbol{p})$.

**Theorem 3.5** (Revelation Principle). *Any mechanism $\mathcal{M} = \langle q, z \rangle$ with a consistent distribution aggregation function $q$ and a monotone payment function $z$ is strategically equivalent to a mechanism $\tilde{\mathcal{M}} = \langle \tilde{q}, \tilde{z} \rangle$ which has a monotone distribution aggregation function $\tilde{q}$ and a monotone payment function $\tilde{z}$.*

Theorem 3.5 can be viewed as a revelation principle in the sense that it simplifies the design choice of aggregation functions. Monotone aggregation functions are a strict subset of consistent aggregation functions since monotonicity directly implies a total order over possible aggregated distributions $Q(\boldsymbol{b}_{-i}, \boldsymbol{p}) = \{q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) : b_i \in \mathbb{Q}_+\}$, with the order naturally determined by the real-numbers' order on $i$'s bid $b_i$ hence this order will be consistent across different $\boldsymbol{b}_{-i}$ and $\boldsymbol{p}$. In this sense, one might think that consistency — which does not impose any restriction when $q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p})$ is not strictly preferred over $q(b'_i, \boldsymbol{b}_{-i}, \boldsymbol{p})$ — might be a significantly weaker requirement on aggregation functions than monotonicity which requires a total and consistent order. Theorem 3.5 shows that this is not the case — they are essentially the same as long as the natural incentive requirement of payment monotonicity is also imposed.

The proof of Theorem 3.5 hinges on the following two lemmas, Lemma 3.6 and Lemma 3.7. Together these two lemmas imply the existence of a strategy mapping, under which the resulting aggregation function becomes monotone. The proof of the theorem is completed by applying the same mapping to the payment function, and noting that this ensures payment monotonicity. We defer the formal proofs of these results to Appendix A.

**Lemma 3.6.** *For any distribution aggregation function $q$, there exists a payment function $z$ such that mechanism $\mathcal{M} = \langle q, z \rangle$ is payment-monotone if and only if $\succeq_i$ establishes a total order over $Q(\boldsymbol{b}_{-i}, \boldsymbol{p}) = \{q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) : b_i \in \mathbb{Q}_+\}$ for any fixed $\boldsymbol{b}_{-i}$ and $\boldsymbol{p}$.*

**Lemma 3.7.** *Consider any consistent aggregation function $q$. Suppose $\succeq_i$ defines a total order over the aggregation set $Q(\boldsymbol{b}_{-i}, \boldsymbol{p}) = \{q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) : b_i \in \mathbb{Q}_+\}$ induced by agent i's bid for any fixed $\boldsymbol{b}_{-i}$ and $\boldsymbol{p}$, then there exist a profile $\pi$ of strategy mappings such that $q(\boldsymbol{b}, \boldsymbol{p}) = \tilde{q}(\pi(\boldsymbol{b}), \boldsymbol{p})$ for some monotone aggregation function $\tilde{q}(\cdot, \cdot)$.*

We conclude this subsection by giving two natural examples used in today's machine learning practice: an example of a monotone aggregation function and a non-monotone one, both with respect to robust preferences (see Definition 2.1).

**Example 3.8** (Linear Aggregation)**.** Consider $q_{\mathsf{KL}}(\boldsymbol{b}, \boldsymbol{p})$ defined as

$$q_{\mathsf{KL}} = \tfrac{1}{B} \sum_{i \in [n]} b_i \cdot p_i, \text{ where } B = \sum_{i \in [n]} b_i.$$

It is easy to verify that this is a monotone aggregation function with respect to robust preferences.

**Example 3.9** (Log-linear Aggregation)**.** Consider the aggregation function $\bar{q}_{\mathsf{KL}}(\boldsymbol{b}, \boldsymbol{p})$ defined by the following equations:

$$\forall t \in T, \ \ln \bar{q}_{\mathsf{KL}}(t) = \tfrac{1}{B} \sum_{i \in [n]} b_i \cdot \ln p_i(t) - C,$$

where $B = \sum_{i \in [n]} b_i$ and $C = \ln \sum_{t \in T} e^{\frac{1}{B} \sum_{i \in [n]} b_i \cdot \ln p_i(t)}$.

The following two-agent example shows that $\bar{q}_{\mathsf{KL}}$ is not monotone: $p_1 = (.5, .4, .1)$ and $p_2 = (.5, .1, .4)$. When $b_1 = b_2$, $\bar{q}_{\mathsf{KL}} = (\sqrt{.25}, \sqrt{.04}, \sqrt{.04})/.9 = (5/9, 2/9, 2/9)$. Fix $b_2 = 1$, either $b_1 \to 0$ or $b_1 \to \infty$, $\bar{q}_{\mathsf{KL}}(t_1) = .5 < 5/9$. Hence $\bar{q}_{\mathsf{KL}}$ is not monotone with respect to robust preferences.

## 3.3 Second Price Payment Rules

Next, we investigate whether monotone aggregation functions admit a "second-price" payment rule, capturing the Vickrey auction's concept of the "minimum bid to win". In the Vickrey auction, the payment corresponds to the critical bid where an agent transitions from losing to winning. To port this notion to our setting, we will show in Theorem 3.12 that, under robust preferences, any monotone aggregation rule can be written as a distribution over deterministic allocations from bids to tokens such that there is a critical bid where the allocation transitions from a less preferred to a more preferred token. Such a critical bid becomes then a natural candidate for the payment.

We will show that besides being payment monotone it has other desirable properties, such as a Myerson-style characterization (Myerson, 1981) in terms of the total variation distance between the preferred distribution and the outcome. To define our payment rule, we first introduce the notion of stable sampling.

### 3.3.1 Stable Sampling

To design the payment for bidder $i$, we analyze a monotone distribution aggregation function $q(\boldsymbol{b}, \boldsymbol{p})$ from $i$'s perspective, fixing the distributions $\boldsymbol{p}$ and the bids of other agents $\boldsymbol{b}_{-i}$. To simplify the notation, we refer to $q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p})$ as $q(b_i)$.

We define an implementation of $q(\cdot)$ as a function $\sigma$ that maps $b_i$ and an exogenous random variable $r \sim \mathcal{R}$ (independent of $\boldsymbol{b}$ and $\boldsymbol{p}$) to a token $t \in T$. Here, $r$ is the *random source* used for implementing the random sampling of token $t$. Hence when $r$ is fixed, $\sigma(b_i, r)$ is fully deterministic. We say that $\sigma : \mathbb{Q}_+ \times \mathcal{R} \rightarrow T$ is a valid implementation of $q(b_i)$ if:

$$\Pr_{r \sim R}[\sigma(b_i, r) = t] = q_t(b_i), \forall t \in T.$$

Next we define what it means for an implementation to be a "stable" sampling. It is useful to partition the token set $T$ into $T_+ = \{t \in T : q_t(0) \leq (p_i)_t\}$ and $T_- = \{t \in T : q_t(0) > (p_i)_t\}$ corresponding to undersampled ($T_+$) and oversampled ($T_-$) tokens. The monotonicity of aggregation function $q$ turns out to be equivalent to the monotonicity of $q_t(b_i)$, as formalized in the following lemma (proof in Appendix B).

**Lemma 3.10.** *Under agents' robust preferences, a distribution aggregation function $q$ is monotone, if and only if for every agent $i$ and $b_i \in \mathbb{R}_+$,*

1. *$\forall t \in T_+$, $q_t(b_i) \leq (p_i)_t$ and $q_t$ weakly increases in $b_i$;*

2. *$\forall t \in T_-$, $q_t(b_i) \geq (p_i)_t$ and $q_t$ weakly decreases in $b_i$.*

We can now define the key notion of stable sampling, and state and prove our main result below.

**Definition 3.11** (Stable Sampling)**.** Let $q_i(b_i)$ be an aggregation function obtained by fixing $\boldsymbol{b}_{-i}$ and $\boldsymbol{p}$, and let $T_+$ and $T_-$ be the sets of undersampled and oversampled tokens for agent $i$. Then we say that an implementation $\sigma$ is stable with respect to aggregation function $q$ if for any $r \in \mathcal{R}$ there are two tokens $u_r \in T_+$ and $o_r \in T_-$ and a threshold $\theta_r \in \mathbb{Q}_+ \cup \{\infty\}$ such that:

$$\sigma(b_i, r) = o_r, \text{ if } b_i < \theta_r; \text{ and } \sigma(b_i, r) = u_r, \text{ if } b_i \geq \theta_r.$$

Here is an illustrative example of stable sampling for a simple distribution supported on two tokens: one undersampled token $u$ with probability $q_u(b_i)$ and an overampled token $o$ with probability $q_o(b_i) = 1 - q_u(b_i)$. Monotonicity of $q$ implies $q_u(b_i)$ weakly increases in $b_i$ and $q_o(b_i)$ weakly decreases. Sample $r$ uniformly from $[0, 1]$. For any $b_i$, if $q_o(b_i) > r$ assign $o$; assign $u$ otherwise. Note that this implementation is: (i) deterministic for fixed $r$; (ii) transits from token $o$ to $u$ as $b_i$ increases (hence $q_o(b_i)$ decreases) under any fixed $r$; (iii) matches the probabilities of $q(b_i)$ in expectation for any $b_i$. The following theorem generalizes this to any number of tokens (proof in Appendix B).

**Theorem 3.12.** *Given a monotone distribution aggregation function $q$ then for any agent $i$ with robust preferences and fixed $\boldsymbol{b}_{-i}$ and $\boldsymbol{p}$ there always exists a stable implementation $\sigma$ of $q$.*

### 3.3.2 A Second Price Rule via Stable Sampling

A stable implementation of a monotone aggregation rule naturally leads to a second-price payment rule: for any given randomness $r$, if the oversampled token $o_r$ is sampled, the agent pays zero as it is the same token that would have been sampled under randomness $r$ if their bid was zero. If the agent's bid was high enough to move from the oversampled $o_r$ to the undersampled token $u_r$, then the agent pays the critical bid $\theta_r$.

Perhaps surprisingly, the expected payment induced by the second price rule described above does not depend on the actual implementation of the sampling algorithm so long as it is stable in the sense of Definition 3.11. The expected payment is proportional to the extent to which the sampling shifts the distribution $q(b_i)$ towards the desired distribution $p_i$ as $b_i$ increases and, interestingly, admits an expression simillar to the standard Myersonian payment as a function of allocations (Myerson, 1981) in which the allocation probability is replaced by the total variation distance between the agent's preferred distribution $p_i$ and the implemented distribution $q(b_i)$.

**Proposition 3.13.** *Consider an arbitrary stable sampling implementation of $q(b_i)$ and the induced second price rule using the critical bid $\theta_r$ as described above. Then the expected payment satisfies*

$$z_i(b_i) = \frac{1}{2} \int_0^{b_i} \left[ \|q(b_i) - p_i\|_1 - \|q(b') - p_i\|_1 \right] db', \qquad \forall b_i \in \mathbb{Q}_+.$$

*Proof.* We again omit the terms $\boldsymbol{b}_{-i}, \boldsymbol{p}$ since they are fixed in each context. If a token $t \in T_-$ is

sampled, the payment is naturally $\zeta_i(b_i, t) = 0$. For a token $t \in T_+$ we have:

$$\zeta_i(b_i, t)q_t(b_i) = \mathbb{E}_r[\theta_r \cdot 1\{\sigma(b_i, r) = t\}]$$
$$= \mathbb{E}_r \int_0^{b_i} \left[1\{\sigma(b_i, r) = t\} - 1\{\sigma(b', r) = t\}\right] \mathrm{d}b'$$
$$= \int_0^{b_i} [q_t(b_i) - q_t(b')] \mathrm{d}b'.$$

Hence:

$$z_i(b_i) = \sum_t \zeta_i(b_i, t)q_t(b_i) = \int_0^{b_i} \left[ \sum_{t \in T_+} q_t(b_i) - \sum_{t \in T_+} q_t(b') \right] \mathrm{d}b'$$
$$= \frac{1}{2} \int_0^{b_i} \left[ \|q(b_i) - p_i\|_1 - \|q(b') - p_i\|_1 \right] \mathrm{d}b',$$

as claimed. □

**Counterfactuals.** The practical advantage of a stable sampling implementation coupled with a second price rule is to offer advertisers a description where it is clear that they only pay if they can improve the outcome. Moreover, with a fixed $r$, advertisers can easily evaluate counterfactuals, as each token can only have one of two possible outcomes, simplifying the comparison.

**Universally Stable Sampling.** We chose to define stable sampling as a single-agent algorithm with fixed $\boldsymbol{b}_{-i}$. One may define a universal stable implementation as $\sigma^{\mathsf{univ}} : \mathbb{Q}^n \times \mathcal{R} \to T$ such that

$$\forall \boldsymbol{b} \in \mathbb{Q}_+^n, t \in T, \ \Pr_{r \sim R}[\sigma^{\mathsf{univ}}(\boldsymbol{b}, r) = t] = q_t,$$

and for any $i$ and $\boldsymbol{b}_{-i}$, $\sigma^{\mathsf{univ}}(\cdot, \boldsymbol{b}_{-i}, r)$ is stable. In Appendix C, we provide counter-examples where such $\sigma^{\mathsf{univ}}$ do not always exist.

# 4 Design of Aggregation Functions

In the previous section, we discussed payment schemes and incentive properties assuming we have an aggregation function. Here we investigate the design of aggregation functions. We adopt two guiding principles: (1) We first define an aggregated loss function to model the overall satisfaction of the agents with the final distribution $q$, weighted by their bids $b_i$. The aggregated loss function has the form:

$$\mathrm{Loss}(\boldsymbol{p}, \boldsymbol{b}, q) = \sum_i b_i \rho(p_i, q),$$

where $\rho : \Delta(T) \times \Delta(T) \to \mathbb{R}$ indicates how different the distribution $q$ is from the preferred $p_i$. (2) The second is to define the function $\rho$ based on the typical loss functions in LLM training.

## 4.1 Review of LLM Training

So far we have assumed that an LLM $f : T^* \to \Delta(T)$ is already trained. In order to discuss the training process, it is useful to recall that an LLM is a neural network parameterized by a vector of weights $W \in \mathbb{R}^N$ in a very high dimensional space. To discuss training, it is useful to think of $f$ as a function of both input and weights:

$$f : T^* \times \mathbb{R}^N \to \Delta(T).$$

We represent the second argument by a superscript $W$. Training is to optimize $W$ such that $f^W(\cdot)$ minimizes a certain loss function over a dataset. A dataset is a sequence of pairs $(x_i, y_i)$ with $x_i \in T^*$ (input sequence) and $y_i \in T$ (label), and a loss function is a function $\ell : T \times \Delta(T) \to \mathbb{R}$. A network typically seeks to minimize:

$$\min_W \sum_i \ell(y_i, f^W(x_i)).$$

A widely used loss function in the first stage is the KL-divergence:

$$\ell_{\mathsf{KL}}(y, x) = -\ln[f^W(y|x)],$$

where we use the notation $f(y|x)$ to represent the probability that a token $y \in T$ is sampled from $f(x) \in \Delta(T)$.

It is useful to think of this problem in the limit where the size of the dataset grows to infinity and it can be effectively thought of as a full-support distribution over $T^* \times T$. In that setting, we can represent the dataset as a distribution $\mu \in \Delta(T^* \times T)$ over input-label pairs $(x, y) \in T^* \times T$. We will write $\mu(x) = \sum_y \mu(x, y)$ to denote the marginal distribution on $x$ and $\mu(\cdot|x)$ to denote the conditional distribution of labels given an input $x$. Then

$$\min_W \mathcal{L}^\mu_{\mathsf{KL}}(f^W), \quad \mathcal{L}^\mu_{\mathsf{KL}}(f) := \sum_{x \in T^*} \mu(x) \cdot D_{\mathsf{KL}}(\mu(\cdot|x) \| f(x)), \tag{KL}$$

where $D_{\mathsf{KL}}(p\|q) = \sum_t p(t) \ln \frac{p(t)}{q(t)}$ is the KL-divergence.

LLMs are typically trained through successive refinement of weights: $W^{\mathsf{pre}} \to W^{\mathsf{SFT}} \to W^{\mathsf{RL}}$. In pre-training we compute $W^{\mathsf{pre}}$ by solving problem (KL) on a generic dataset $\mu^{\mathsf{pre}}$ via stochastic

gradient descent. In the second stage, we initialize the weights as $W = W^{\mathsf{pre}}$ and run stochastic gradient descent to solve the same problem (KL) on a more specialized dataset $\mu^{\mathsf{SFT}}$. In other words, we solve the same problem on two different datasets.

The problem in the RLHF stage is different. The dataset only contains $x$ (still represented by $\mu$) and for any $y$, we have a function $r(x, y)$ giving the reward of mapping $x$ to $y$. Then $W^{\mathsf{RL}}$ is obtained by maximizing reward while minimizing the distance to the function $f^{\mathsf{SFT}}$ from previous stages (the PPO algorithm (Ouyang et al., 2022; Schulman et al., 2017)):

$$\max_{W} \mathcal{L}_{\mathsf{RL}}^{\mu,r}(f^W), \tag{RL}$$

$$\mathcal{L}_{\mathsf{RL}}^{\mu,r}(f) := \sum_{x \in T^*} \mu(x) \Big[ \sum_{y} r(x,y) f(y|x) - \beta D_{\mathsf{KL}}(f(x) \| f^{\mathsf{SFT}}(x)) \Big].$$

**KL-divergence and Entropy.** For the next propositions it is useful to recall that the entropy of a distribution $p \in \Delta(T)$ is $H(p) = -\sum_{t \in T} p(t) \ln p(t)$. Given two distributions $p, q \in \Delta(T)$, the cross entropy of $q$ relative to $p$ is $H(p, q) = -\sum_{t \in T} p(t) \ln q(t)$. Hence we can write $D_{\mathsf{KL}}(p \| q) = H(p, q) - H(p)$. We will also use Gibbs' inequality $H(p) \leq H(p, q), \forall p, q \in \Delta(T)$.

## 4.2 KL-inspired Aggregation

The first aggregation method is based on the (KL) program. When trying to aggregate LLMs $f_i$ according to bids $b_i$, we will design a function $q$ that mimics the outcome of the following thought experiment. We will imagine that each LLM was obtained by solving the (KL) on a dataset represented by $\mu_i$, where the marginal over inputs are the same $\mu_i(x) = \mu(x), \forall i$ but potentially differ on the marginals on the labels $\mu_i(y|x)$. In this thought experiment, we will combine their LLMs by re-training an LLM on the combined labels weighted by the bids. In other words, we will solve the problem:

$$\min_W \ \mathcal{L}_{\mathsf{KL}}^{\bar{\mu}}(f^W), \qquad \text{where } \bar{\mu} = \frac{\sum_i b_i \mu_i}{\sum_i b_i}. \tag{3}$$

The next proposition morally says that we can obtain a solution to the (KL) problem on the aggregated dataset (3) by combining its solutions on individual datasets. The proof appears in Appendix D.

**Proposition 4.1.** *Consider datasets $\mu_i$ such that $\mu_i(x) = \mu(x), \forall i, x$ and let $\bar{\mu}$ be their weighted average. Let $f_i$ be the minimizer of $\mathcal{L}_{\mathsf{KL}}^{\mu_i}$ and $f^*$ be the minimizer of $\mathcal{L}_{\mathsf{KL}}^{\bar{\mu}}$, the loss on the aggregated dataset. Then $f^*$ is the solution to:*

$$\min_f \sum_x \mu(x) \sum_i b_i D_{\mathsf{KL}}(f_i(x) \| f(x)). \tag{4}$$

19

Proposition 4.1 motivates the following aggregated loss function:

$$\text{LOSS}_{\text{KL}} = \sum_{i \in [n]} b_i \cdot \rho_{\text{KL}}(p_i, q) = \sum_{i \in [n]} \sum_{t \in T} b_i \cdot p_i(t) \cdot \ln \frac{p_i(t)}{q(t)}.$$

Now, we characterize the aggregation function that minimizes $\text{LOSS}_{\text{KL}}$. The proof in Appendix D uses Gibb's inequality.

**Lemma 4.2.** *The aggregation function that minimizes* $\text{LOSS}_{\text{KL}}$ *is the linear combination of* $p_i$*:*

$$\forall t \in T, \ q_{\text{KL}}(t) = \frac{\sum_i b_i \cdot p_i(t)}{\sum_i b_i}.$$

Besides being monotone and aligned with how LLMs are trained, this aggregation function has the advantage that we only need to call a single LLM to sample each token. We can choose an index $i$ proportionally to the bids and then sample a token from the $i$-th LLM. To compute second price payments, however, we still need to query the LLMs for all agents.

## 4.3 RL-inspired Aggregation

Now, consider a different thought experiment where all agents use the same pre-trained and fine-tuned model of weights $W^{\text{SFT}}$, but employs a different reward function $r_i(x, y)$ for RLHF. We combine their LLMs by solving the (RL) problem on the combined reward functions weighted by the bids. We will solve the problem:

$$\max_W \mathcal{L}_{\text{RL}}^{\mu, \bar{r}}(f^W), \quad \bar{r} = \frac{\sum_i b_i r_i}{\sum_i b_i}.$$

Analogously to Proposition 4.1, we show that we can obtain the solution to the aggregated problem by combining the solutions on each dataset. We defer the proof of Proposition 4.3 to Appendix D.

**Proposition 4.3.** *Consider datasets* $\mu, r_i$ *and let* $f^{\text{SFT}}$ *be the solution of program* (KL) *with data* $\mu$*. If* $f_i$ *is the maximizer of* $\mathcal{L}_{RL}^{\mu, r_i}$*, let* $f^*$ *be the maximizer of* $\mathcal{L}_{RL}^{\mu, \bar{r}}$ *where* $\bar{r}$ *is the weighted average of the reward functions, then* $f^*$ *is the function minimizing:*

$$\min_f \sum_x \mu(x) \sum_i b_i D_{\text{KL}}(f(x) \| f_i(x)). \tag{5}$$

Similar to what we did in the previous subsection, we can also define an aggregated loss function inspired by Proposition 4.3. Namely:

$$\overline{\text{Loss}}_{\text{KL}} = \sum_{i \in [n]} b_i \cdot \rho_{\text{KL}}(q, p_i) = \sum_{i \in [n]} \sum_{t \in T} b_i \cdot q(t) \cdot \ln \frac{q(t)}{p_i(t)}.$$

**Lemma 4.4.** *The aggregation function that minimizes* $\overline{\text{Loss}}_{\text{KL}}$ *is the log-linear combination of* $p_i$:

$$\forall t \in T, \ \ln \bar{q}_{\text{KL}}(t) = C + \frac{\sum_i b_i \ln \cdot p_i(t)}{\sum_i b_i},$$

*where C is a normalization constant such that* $\sum_t \bar{q}_{\text{KL}}(t) = 1$.

The proof follows directly from the proof of Proposition 4.3. While not monotone (Example 3.9), the log-linear aggregation function is a reasonable choice assuming that the agents' preferences are aligned with the KL-divergence loss for the RL-stage training.

# 5 Demonstration

We implement the aggregation functions proposed in Section 4 and discuss the examples they produce. Off-the-shelf LLMs generate full text passages. In our case, we need to peek at the internal states of LLMs (the probability distributions over tokens) at each token generation stage. Therefore, we use a custom version of the Google Bard model with a modified inference method that allows access to the token distributions.

Starting from the same base model, we customize it for different agents by prompt-tuning. We start with a base model $f : T^* \to \Delta(T)$ and for each agent $i$ we come up with a "prompt" $s_0^i \in T^*$ and now we define for each agent $i$ the LLM $f_i : T^* \to \Delta(T)$ as:

$$f_i(s) = f(s_0^i \oplus s)$$

Therefore if $\tau_1, \ldots, \tau_{k-1}$ are the first $k-1$ tokens generated, then the preferred distribution of agent $i$ over the $k$-th token is given by:

$$p_i = f(s_0^i \oplus s \oplus \tau_1 \oplus \cdots \oplus \tau_{k-1}),$$

A key advantage of simulating LLM agents with different prompts is the ability to use a single LLM, making multiple queries with different prompts instead of serving multiple LLMs concurrently. Because of their large sizes, serving multiple LLMs can be very costly and practically challenging. As one of the key strengths of LLMs, the flexibility to accomplish various tasks with properly designed prompts sheds light to the possibility of training one universal LLM that can, for

example, generate different ads according to agent-specific prompts. That is, the universal advertising LLM, plus an advertiser-specific prompt, behaves like an advertiser-specific LLM through the online in-context few-shot learning.

## 5.1 Setups

We illustrate our method with a co-marketing example as well as a competing brands example. In the setup of co-marketing, the two agents would like to advertise for their brands, "Alpha Airlines" and "Beta Resort" respectively, regarding a shared topic "Hawaii." We intentionally choose fictitious brands in order to avoid the model directly retrieving any existing ads. We use the brand names "Alpha" and "Beta" that do not have strong meanings on their own to minimize any potential hallucination, as we are using a common purposed LLM that is not optimized for our task. Each agent is given the following prompt:

> "You are an expert of writing texts that naturally combines two ads together. Your choice of words and sentences is full of artistic flair.
>
> Write a one-sentence ad for ―――――."

Agent $A$ uses *"a flight to Hawaii using [**Alpha Airlines**]"* to fill the blank, while agent $B$ uses *"a vacation in Hawaii at the [**Beta Resort**]"*. The first two sentences in the prompt aim to improve the quality of the ad generation through *assigning roles* (see, for example, (Wu et al., 2023)). A natural question is whether the proposed method can adjust the combining strategy according to the context. Since in both the linear aggregation rule $q_{\mathsf{KL}}$ and the log-linear aggregation rule $\bar{q}_{\mathsf{KL}}$, there is only one degree of freedom, we parameterize the response by $\lambda = {}^{b_1}/_{(b_1+b_2)}$.

The demonstration for competing brands uses the same setup, except that the two brands are "Beta Resort" and "Gamma Hotel". We seek to understand whether competition will lead LLMs to avoid advertising for both advertisers in the same creative.

## 5.2 Results

The results for the co-marketing example are listed in Table 2, where from top to bottom, the value of $\lambda$ decreases from 1 to 0. As we can see for both aggregation functions, the generated texts roughly follow the pattern of "only Alpha Airlines" → "both Alpha Airlines and Beta Resort" → "only Beta Resort" when $\lambda$ goes from 1 to 0. This is expected, as $\lambda$ going from 1 to 0 corresponds to $b_2$ increasing from 0 to $\infty$ with $b_1$ fixed (or $b_1$ decreasing from $\infty$ to 0 with $b_2$ fixed). The thresholds of pattern changes are 0.75 and 0.4 for the linear aggregation, and 0.5 and 0.45 for the

| $\lambda$ | Linear aggregation function $q_{\text{KL}}$ | Log-linear aggregation function $\bar{q}_{\text{KL}}$ |
|---|---|---|
| 1 | **Alpha Airlines**: Your ticket to paradise. | |
| 0.75 | **Alpha Airlines**: Fly to Hawaii and experience the beauty of the islands with aloha. **Beta Resorts**: Stay at our resorts and enjoy the best of Hawaii. | Feel the magic of Hawaii with a flight on **Alpha Airlines**, now offering 20% off all flights when you book with us today! |
| 0.6 | **Alpha Airlines** flies you to Hawaii, where you can enjoy a week-long stay at the **Beta resort** for just $1000. | Experience the magic of Hawaii with a flight on **Alpha Airlines**, now offering 20% off all flights when you book with us today! |
| 0.55 | **Alpha Airlines** flies you to Hawaii, where you can enjoy a week-long stay at the **Beta resort**. | Escape to the tropical paradise of Hawaii with [**Alpha Airlines**]! |
| 0.5 | **Alpha Airlines** flies you to Hawaii, where you can enjoy a beautiful sunset on the beach. Stay 3 nights and get the 4th free at the **Beta Resort**. | Experience the magic of Hawaii with a stay at the luxurious [**Beta Resort**] and a refreshing flight on [**Alpha Airlines**]. |
| 0.45 | Fly **Alpha Airlines** to sunny Hawaii and enjoy the secluded beaches and private lagoons of the Royal Hawaiian **Beta Resort**. | Experience the magic of Hawaii with a stay at the luxurious [**Beta Resort**] and a special flight offer from [**Alpha Airlines**]. |
| 0.4 | Fly **Alpha Airlines** to sunny Hawaii and enjoy the first-class treatment that awaits you at **Beta Resort**, all for one low price. | Experience the magic of Hawaii at the [**Beta Resort**], where you'll feel like you're in a tropical paradise. |
| 0.25 | Experience the magic of Hawaii at the **Beta Resort**, where the sun shines brighter and the waves crash louder — book your stay today with our exclusive 20% off discount! | Experience the magic of Hawaii at the **Beta Resort**, where you'll be pampered like royalty and surrounded by breathtaking beauty. |
| 0 | Hawaii's **Beta Resort**: a paradise where the sun shines brighter, the waves sing sweeter, and the sand feels softer. | |

Table 2: Text generation from two aggregation functions with different $\lambda = b_1/(b_1 + b_2)$.

log-linear aggregation. We emphasize that the example is generated with a general purposed LLM, and it is reasonable to believe that the performance can be improved with proper fine-tuning for the specific task at hand.

The results for the example with competing brands is listed in Table 3. As we can see for both aggregation functions, the generated texts roughly follow the pattern of "only Gamma Hotel" → "only Beta Resort" as $\lambda$ goes from 1 to 0. The thresholds of pattern changes are $0.5 \sim 0.45$ for both the linear aggregation and the log-linear aggregation. We note that for $\lambda = 0.45$ and $\lambda = 0.5$, the generated results from the linear aggregation seem to get confused on the brands of "Gamma Hotel" and "Beta Resort" to generate non-existing "Gamma Resort" and "Beta Hotel". We expect such phenomenons can be addressed by additional task specific training to enforce the

| $\lambda$ | Linear aggregation function $q_{\mathsf{KL}}$ | Log-linear aggregation function $\bar{q}_{\mathsf{KL}}$ |
|---|---|---|
| 1 | For those who want a relaxing vacation in Hawaii, the **Gamma Hotel** is the perfect place to stay. | |
| 0.75 | Experience the magic of Hawaii at the **Gamma Hotel**, where the sun shines brighter, the waves crash louder, and the nights are filled with endless possibilities. | Experience the magic of Hawaii at the [**Gamma Hotel**], where you'll be pampered by our luxurious amenities and stunning views of the Pacific Ocean. |
| 0.6 | For those who want a truly unforgettable vacation, come to Hawaii and stay at the [**Gamma Hotel**]. | Plan your dream vacation in Hawaii with our special rates at the luxurious [**Gamma Hotel**]. |
| 0.55 | For those who want a truly unforgettable vacation, come to Hawaii and stay at the [**Gamma Hotel**]. | Plan your dream vacation in Hawaii with our special rates at the luxurious [**Gamma Hotel**]! |
| 0.5 | For those who want a relaxing vacation in Hawaii, the *Gamma Resort* is the perfect place to stay. | **Gamma Hotel**: a touch of heaven in the heart of Hawaii. |
| 0.45 | Hawaii is a beautiful place to vacation, and the *Beta Hotel* is the perfect place to stay. | Plan your dream vacation in Hawaii with our special rates at the luxurious [**Beta Resort**]. |
| 0.4 | Escape to the beautiful beaches of Hawaii and stay at the luxurious **Gamma Hotel**, where you'll be treated like royalty. | Escape to the tropical paradise of Hawaii and stay at the luxurious [**Beta Resort**], where you'll enjoy stunning views of the Pacific Ocean and easy access to all the island has to offer. |
| 0.25 | Escape to the tropical paradise of Hawaii and stay at the luxurious **Beta Resort**, where you'll be surrounded by lush greenery and stunning ocean views. | Escape to the tropical paradise of Hawaii and stay at the luxurious **Beta Resort**, where you'll be pampered from the moment you arrive. |
| 0 | Escape to the tropical paradise of Hawaii and stay at the luxurious **Beta Resort**, where you'll be pampered from the moment you arrive. | |

Table 3: Text generation from two aggregation functions with different $\lambda = b_1/(b_1 + b_2)$ in the competing setup.

extra requirement on using the exact terms on brands.

# 6 Conclusion

In this work, we put forward a mechanism-design approach to the problem of aggregating the output of several LLMs into one. We have proposed an auction format, the *token auction model*, which operates on a token-by-token basis and allows the LLM agents to influence the output through single-dimensional bids. We explored mechanism design with such LLMs through a robust lens,

which makes minimal assumptions about the agents' preferences. Working under this paradigm, we have shown that natural requirements on the incentive properties of the auction mechanism, imply that it should feature a monotone aggregation function. We then showed that under robust preferences, any monotone aggregation function enables second-price style payments (akin to those in the Vickrey auction or the Generalized-Second Price auction). We also explored the design of aggregation functions inspired by common loss functions used in LLM training, such as KL-divergence or PPO (in RLHF). As a "proof of concept" for our designed mechanism, we demonstrated promising outcomes of our aggregation methods by implementing these aggregation functions in a real-world state-of-the-art LLM using prompt tuning.

# References

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Chris Olah, Benjamin Mann, and Jared Kaplan. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *CoRR* abs/2204.05862 (2022). https://doi.org/10.48550/arXiv.2204.05862 6

Michiel A. Bakker, Martin J. Chadwick, Hannah Sheahan, Michael Henry Tessler, Lucy Campbell-Gillingham, Jan Balaguer, Nat McAleese, Amelia Glaese, John Aslanides, Matt M. Botvinick, and Christopher Summerfield. 2022. Fine-tuning language models to find agreement among humans with diverse preferences. In *NeurIPS 2022*. 38176–38189. 6

Dirk Bergemann and Stephen Morris. 2005. Robust mechanism design. *Econometrica* (2005), 1771–1813. 7

Dirk Bergemann and Stephen Morris. 2012. *Robust mechanism design: The role of private information and higher order beliefs*. Vol. 2. World Scientific. 7

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz

Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. *NeurIPS 2020*, 1877–1901. 2, 6, 8, 9

Jeremy Bulow and John Roberts. 1989. The simple economics of optimal auctions. *Journal of Political Economy* 97, 5 (1989), 1060–1090. 34

Gabriel Carroll. 2015. Robustness and linear contracts. *American Economic Review* 105, 2 (2015), 536–563. 7

Kumar Avinava Dubey, Zhe Feng, Rahul Kidambi, Aranyak Mehta, and Di Wang. 2024. Auctions with LLM Summaries. *CoRR* abs/2404.08126 (2024). arXiv:2404.08126 https://doi.org/10.48550/arXiv.2404.08126 7

Shaddin Dughmi. 2011. A truthful randomized mechanism for combinatorial public projects via convex optimization. In *EC 2011*. 263–272. 7

Paul Dütting, Tim Roughgarden, and Inbal Talgam-Cohen. 2019. Simple versus Optimal Contracts. In *Proceedings of the 2019 ACM Conference on Economics and Computation*. 369–387. 7

Benjamin Edelman, Michael Ostrovsky, and Michael Schwarz. 2007. Internet Advertising and the Generalized Second-Price Auction: Selling Billions of Dollars Worth of Keywords. *American Economic Review* 97(1) (2007), 242–259. 1

Rupert Freeman, David M Pennock, Dominik Peters, and Jennifer Wortman Vaughan. 2019. Truthful aggregation of budget proposals. In *Proceedings of the 2019 ACM Conference on Economics and Computation*. 751–752. 7

Ashish Goel, Anilesh K Krishnaswamy, Sukolsak Sakshuwong, and Tanja Aitamurto. 2019. Knapsack Voting for Participatory Budgeting. *ACM Transactions on Economics and Computation* 7, 2 (2019). 7

Google, Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave,

Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea Hu, Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li, Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu, Frederick Liu, Marcello Maggioni, Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Moussalem, Zachary Nado, John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alex Castro Ros, Aurko Roy, Brennan Saeta, Rajkumar Samuel, Renee Shelby, Ambrose Slone, Daniel Smilkov, David R. So, Daniel Sohn, Simon Tokumine, Dasha Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, Pidong Wang, Zirui Wang, Tao Wang, John Wieting, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav Petrov, and Yonghui Wu. 2023. PaLM 2 Technical Report. arXiv:2305.10403 [cs.CL] 2, 8, 9

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The Curious Case of Neural Text Degeneration. In *International Conference on Learning Representations*. 3

Tommy Mordo, Moshe Tennenholtz, and Oren Kurland. 2024. Sponsored Question Answering. In *Proceedings of the ACM SIGIR International Conference on the Theory of Information Retrieval*. https://openreview.net/forum?id=JVh7rytFjh 8

Roger B Myerson. 1981. Optimal auction design. *Mathematics of Operations Research* 6, 1 (1981), 58–73. 15, 16

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *NeurIPS 2022*. 27730–27744. 6, 9, 19

Christos Papadimitriou, Michael Schapira, and Yaron Singer. 2008. On the hardness of being truthful. In *FOCS 2008*. 250–259. 7

Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. In *ICML 2021*. 8821–8831. 3

Tim Roughgarden and Inbal Talgam-Cohen. 2016. Optimal and robust mechanism design with interdependent values. *ACM Transactions on Economics and Computation* 4 (3) (2016), 1–34. 7

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal Policy Optimization Algorithms. *CoRR* abs/1707.06347 (2017). http://arxiv.org/abs/1707.06347 19

Amartya Sen. 2018. *Collective choice and social welfare*. Harvard University Press. 30

Ermis Soumalias, Michael J. Curry, and Sven Seuken. 2024. Truthful Aggregation of LLMs with an Application to Online Advertising. *CoRR* abs/2405.05905 (2024). https://doi.org/10.48550/ARXIV.2405.05905 arXiv:2405.05905 7

Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. 2019. Videobert: A joint model for video and language representation learning. In *ICCV 2019*. 7464–7473. 3

Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch, Marc Pickett, Kathleen S. Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Agüera-Arcas, Claire Cui, Marian Croak, Ed H. Chi, and Quoc Le. 2022. LaMDA: Language Models for Dialog Applications. *CoRR* abs/2201.08239 (2022). https://arxiv.org/abs/2201.08239 2, 8

Hal R. Varian. 2007. Position auctions. *International Journal of Industrial Organization* 25 (6) (2007), 1163–1178. 1

William Vickrey. 1961. Counterspeculation, auctions, and competitive sealed tenders. *The Journal of Finance* 16, 1 (1961), 8–37. 5

Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned Language Models are Zero-Shot Learners. In *International Conference on Learning Representations*. 6, 9

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In *NeurIPS 2022*. 24824–24837. 6

Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, and Tengyu Ma. 2023. Larger language models do in-context learning differently. *CoRR* abs/2303.03846 (2023). https://doi.org/10.48550/arXiv.2303.03846 6

Ning Wu, Ming Gong, Linjun Shou, Shining Liang, and Daxin Jiang. 2023. Large language models are diverse role-players for summarization evaluation. *arXiv preprint arXiv:2303.15078* (2023). 22

Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. 2022. Scaling Autoregressive Models for Content-Rich Text-to-Image Generation. *Transactions on Machine Learning Research* (2022). 3

# A  Omitted Proofs from Section 3.2

Our proofs make use of the following technical lemma, which shows that any complete preference relation on a countable set can be expressed via a function that maps elements of the countable set into the non-negative rationals.

**Lemma A.1.** *Let $X$ be a countable set, and let $\succeq$ be a complete preference relation on $X$. Then there is a function $f : X \to \mathbb{Q}_+$ such that $x \preceq y$ if and only if $f(x) \leq f(y)$.*

*Proof.* Since $X$ is countable, we can arrange it in a sequence $\{x_1, x_2, \ldots\}$. We describe a procedure to define $f(x_i)$ for each $x_i$ in turn. Let $f(x_1) = 1$. Suppose we have defined $f(x_1)$, $f(x_2)$, ..., $f(x_n)$ in such a way that all order relations are preserved. That is, for all $i, j \leq n$, we have $x_i \preceq x_j$ if and only if $f(x_i) \leq f(x_j)$. We want to define $f(x_{n+1})$. If $x_{n+1} \succeq x_i$ and $x_{n+1} \preceq x_i$ (i.e., $x_{n+1} \sim x_i$) for some $i \leq n$, then we can let $f(x_{n+1}) = f(x_i)$ and we have extended the function to one more element in a way that preserves all order relations until (including) $x_{n+1}$. Otherwise, we can partition the set $\{x_1, \ldots, x_n\}$ into two sets

$$L = \{x_i : i \leq n, \ x_i \prec x_{n+1}\} \quad \text{and} \quad H = \{x_i : i \leq n, \ x_i \succ x_{n+1}\}.$$

In the set $\mathbb{Q}_+$ of all positive rational numbers, every element in $f(L) = \{f(x_i)\}_{x_i \in L}$ is strictly smaller than every element in $f(H)$. Choose $q \in \mathbb{Q}_+$ strictly larger than all the elements of $f(L)$ and strictly smaller than all the elements of $f(H)$. For each $x_i$ with $i \leq n$ the relationship between $x_i$ and $x_{n+1}$ is the same as the relationship between $f(x_i)$ and $q$. Therefore, letting $f(x_{n+1}) = q$ extends the function to one more element, preserving all order relations, as required. The resulting function defined on all of $X$ is thus an isomorphism from $X$ to the image of $f$. $\qquad\square$

We remark that the assumption of $X$ being countable in Lemma A.1 is crucial as the lemma does not hold for uncountable set $X$. A canonical counter example is $X = \mathbb{R}^2$ with the lexicographic ordering: $x > x'$ if and only if either $x_1 > x_1'$ or $x_1 = x_1' \wedge x_2 \geq x_2'$. It is well-known that there is no utility function (continuous or non-continuous) that can induce this complete order (Sen, 2018). This is also why we restrict bids to be non-negative rational numbers, which is certainly practical but also has underlying technical reasons.

## Proof of Lemma 3.6

*Proof.* We first prove the "only if" ("$\implies$") direction. That is, suppose $\mathcal{M} = \langle q, z \rangle$ is payment-monotone, then it must imply a total order over $Q(\boldsymbol{b}_{-i}, \boldsymbol{p}) = \{q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) : b_i \in \mathbb{Q}_+\}$ for any fixed $\boldsymbol{b}_{-i}$ and $\boldsymbol{p}$.

Fix any $\boldsymbol{b}_{-i}$ and $\boldsymbol{p}$. For any $q, q' \in Q(\boldsymbol{b}_{-i}, \boldsymbol{p}) = \{q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) : b_i \in \mathbb{Q}_+\}$ such that $q = q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p})$, $q' = q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p})$. Let $z_i = z_i(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p})$ and $z_i' = z_i(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p})$ be the corresponding payment given by $\mathcal{M}$. Without loss of generality, suppose that $z_i \geq z_i'$. Then by payment-monotonicity of the mechanism $\mathcal{M}$, we have $q \succeq_i q'$. In other words, $\succeq_i$ establishes a total order over $Q$.

Next we show the "if" ("$\impliedby$") direction. That is, given a total order over $Q(\boldsymbol{b}_{-i}, \boldsymbol{p}) = \{q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) : b_i \in \mathbb{Q}_+\}$ for any fixed $\boldsymbol{b}_{-i}$ and $\boldsymbol{p}$, we can construct a payment rule $z$ such that $\mathcal{M} = \langle q, z \rangle$ is payment-monotone.

Fix any $\boldsymbol{b}_{-i}$ and $\boldsymbol{p}$. Since $\succeq_i$ establishes a total order over $Q(\boldsymbol{b}_{-i}, \boldsymbol{p})$ and $Q(\boldsymbol{b}_{-i}, \boldsymbol{p})$ is countable, by Lemma A.1 there exists a function $f_{i, \boldsymbol{b}_{-i}, \boldsymbol{p}} : Q(\boldsymbol{b}_{-i}, \boldsymbol{p}) \to \mathbb{Q}_+$ such that $\forall q, q' \in Q(\boldsymbol{b}_{-i}, \boldsymbol{p})$,

$$q \succeq_i q' \iff f_{i, \boldsymbol{b}_{-i}, \boldsymbol{p}}(q) \geq f_{i, \boldsymbol{b}_{-i}, \boldsymbol{p}}(q').$$

Letting $z_i(\boldsymbol{b}, \boldsymbol{p}) = f_{i, \boldsymbol{b}_{-i}, \boldsymbol{p}}(q(\boldsymbol{b}, \boldsymbol{p}))$, we obtain a payment-monotone mechanism $\langle q, z \rangle$. $\qquad\square$

## Proof of Lemma 3.7

*Proof.* Given the consistency of the distribution aggregation function $q$, we can define a preference relation $\succeq_{i,p}$ on $\mathbb{Q}_+$ according to the total orders over $Q(\boldsymbol{b}_{-i}, \boldsymbol{p})$ for all possible $\boldsymbol{b}_{-i}$ (assumed by the lemma) such that, $\forall b_i, b_i' \in \mathbb{Q}_+$:

- If there exists $\boldsymbol{b}_{-i}$ such that $q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) \succ_i q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p})$, then $b_i$ is preferred to $b_i'$, i.e., $b_i \succ_{i,p} b_i'$;

- If $\forall \boldsymbol{b}_{-i}$, $q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) = q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p})$, then $b_i$ and $b_i'$ are indifferent, i.e., $b_i \sim_{i,p} b_i'$.

Note that the validity of this construction is guaranteed by the definition of consistency (Definition 3.2):

- For any $b_i, b_i'$ assigned $b_i \succ_{i,p} b_i'$, there exists some $\boldsymbol{b}_{-i}$ such that $q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) \succ_i q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p})$. By the definition of consistency, $\forall \boldsymbol{b}_{-i}'$, $q(b_i, \boldsymbol{b}_{-i}', \boldsymbol{p}) \succeq_i q(b_i', \boldsymbol{b}_{-i}', \boldsymbol{p})$. So we never assign the contradicting order ($b_i' \succ_{i,p} b_i$).

- Meanwhile, for any $b_i, b_i'$ assigned $b_i \sim_{i,p} b_i'$, we have that $\forall \boldsymbol{b}_{-i}$, $q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) = q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p})$, which implies that for all $\boldsymbol{b}_{-i}$, $q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) \not\succ_i q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p})$. Therefore, we never assign $b_i \succ_{i,p} b_i'$.

By the assumption of the lemma, $\succeq_i$ establishes a total order over $Q(\boldsymbol{b}_{-i}, \boldsymbol{p})$. We argue that this implies that $\succeq_{i,p}$ is a complete preference relation on $\mathbb{Q}_+$. Concretely, for every pair $b_i, b_i' \in \mathbb{Q}_+$ we have $q(b_i), q(b_i') \in Q(\boldsymbol{b}_{-i}, \boldsymbol{p})$. So either for some $\boldsymbol{b}_{-i}$, $q(b_i) \succ_i q(b_i')$ (or $q(b_i') \succ_i q(b_i)$), or for all $\boldsymbol{b}_{-i}$, $q(b_i) = q(b_i')$, due to the lemma's assumption of total order over $Q(\boldsymbol{b}_{-i}, \boldsymbol{p})$. Hence every pair $b_i, b_i' \in \mathbb{Q}_+$ has an order under $\succeq_{i,p}$ (i.e., $\succ_{i,p}$ or $\sim_{i,p}$), so the preference relation $\succeq_{i,p}$ is complete.

By Lemma A.1, since $\mathbb{Q}_+$ is countable, there exists a function $f_{i,p} : \mathbb{Q}_+ \rightarrow \mathbb{Q}_+$ that represents the preference relation $\succeq_{i,p}$, i.e.,

$$\forall b_i, b_i' \in \mathbb{Q}_+, \ b_i \succeq_{i,p} b_i' \iff f_{i,p}(b_i) \geq f_{i,p}(b_i').$$

Hence we can construct $\tilde{q}$ as:

$$\tilde{q}(f_{i,p}(b_i), \boldsymbol{b}_{-i}, \boldsymbol{p}) = q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}).$$

Note that for any $b_i, b_i'$ such that $f_{i,p}(b_i) = f_{i,p}(b_i')$, we have $b_i \sim_{i,p} b_i' \implies q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) = q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p})$ for every $\boldsymbol{b}_{-i}, \boldsymbol{p}$. Hence, the construction of $\tilde{q}$ never assigns different values to the same input.

To complete the construction of $\tilde{q}(\cdot, \boldsymbol{b}_{-i}, \boldsymbol{p})$, we expand its definition from the image set of $f_{i,\boldsymbol{p}}$ to the entire bid space while preserving monotonicity. The existence of the expansion is guaranteed by the monotonicity of $\tilde{q}(\cdot, \boldsymbol{b}_{-i}, \boldsymbol{p})$ on the image set of $f_{i,\boldsymbol{p}}$ and the fact that the space of aggregations is compact.

Letting $\pi_i(b) = f_{i,\boldsymbol{p}}(b)$, we have $q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) = \tilde{q}(\pi_i(b_i), \boldsymbol{b}_{-i}, \boldsymbol{p})$ and $\tilde{q}$ is a monotone distribution aggregation function for agent $i$. Applying the same argument and the relabeling procedure for every $i$ from 1 to $n$, completes the proof. $\qquad\square$

## Proof of Theorem 3.5

*Proof.* Applying Lemma 3.6, we know that the payment monotonicity of $\mathcal{M}$ implies a total order over $Q(\boldsymbol{b}_{-i}, \boldsymbol{p})$ for any fixed $\boldsymbol{b}_{-i}$ and $\boldsymbol{p}$. Then applying Lemma 3.7, we know that for any $\boldsymbol{p}$ there exist strategy mappings $\pi_i : \mathbb{Q}_+ \to \mathbb{Q}_+$ for each $i$ such that $q(\boldsymbol{b}, \boldsymbol{p}) = \tilde{q}(\pi(\boldsymbol{b}), \boldsymbol{p})$ and $\tilde{q}(\cdot, \boldsymbol{p})$ is a monotone aggregation function.

Now, following the same argument in the proof of Lemma 3.7, let us further define $\tilde{z}(\pi(\boldsymbol{b}), \boldsymbol{p}) = z(\boldsymbol{b}, \boldsymbol{p})$ for each $\boldsymbol{b} \in \mathbb{Q}_+^n$ and $\boldsymbol{p} \in \Delta(T)^n$. So $\mathcal{M} = \langle q, z \rangle$ is strategically equivalent to $\tilde{\mathcal{M}} = \langle \tilde{q}, \tilde{z} \rangle$ by definition.

It remains to show that the mechanism $\tilde{\mathcal{M}} = \langle \tilde{q}, \tilde{z} \rangle$ is payment-monotone. We know that the original mechanism $\mathcal{M} = \langle q, z \rangle$ satisfies payment monotonicity, meaning that for each $\boldsymbol{p}$, $\boldsymbol{b}_{-i}$, $b_i$, $b_i'$,

$$z_i(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) \geq z_i(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p}) \iff q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) \succeq_i q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p}).$$

But then, with $\boldsymbol{b} = (b_i, \boldsymbol{b}_{-i})$ and $\boldsymbol{b}' = (b_i', \boldsymbol{b}_{-i})$, we also have

$$\tilde{z}_i(\pi(\boldsymbol{b}), \boldsymbol{p}) = z_i(\boldsymbol{b}, \boldsymbol{p}) \geq z_i(\boldsymbol{b}', \boldsymbol{p}) = \tilde{z}_i(\pi(\boldsymbol{b}'), \boldsymbol{p})$$
$$\iff \quad \tilde{q}(\pi(\boldsymbol{b}), \boldsymbol{p}) = q(\boldsymbol{b}, \boldsymbol{p}) \succeq_i q(\boldsymbol{b}', \boldsymbol{p}) = \tilde{q}(\pi(\boldsymbol{b}'), \boldsymbol{p}),$$

so the pair $\tilde{q}$, $\tilde{z}$ satisfies payment monotonicity as needed. $\qquad\square$

# B Omitted Proofs from Section 3.3

## Proof of Lemma 3.10

*Proof.* We first prove the "only if" ("$\Longrightarrow$") direction. Suppose $q$ is a monotone distribution aggregation function. By Definition 3.3, for any agent $i$ and $b_i' \geq b_i \geq 0$, we have

$$q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p}) \succeq_i q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) \succeq_i q(0, \boldsymbol{b}_{-i}, \boldsymbol{p}).$$

For any undersampled token $t \in T_+$, because $q_t(0, \boldsymbol{b}_{-i}, \boldsymbol{p}) \leq (p_i)_t$, then by Definition 2.1, we have

$$q_t(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}), q_t(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p}) \in [q_t(0, \boldsymbol{b}_{-i}, \boldsymbol{p}), (p_i)_t].$$

Hence by $q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p}) \succeq_i q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p})$, we have

$$q_t(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p}) \geq q_t(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}),$$

namely, $q_t(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p})$ weakly increases with $b_i$ and never goes above $(p_i)_t$.

Similarly, we can prove that for any oversampled token $t \in T_-$, $q_t(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p})$ weakly decreases with $b_i$ and never goes below $(p_i)_t$.

Then we prove the "if" ("$\Longleftarrow$") direction. Consider any $b_i' \geq b_i$. For any undersampled token $t \in T_+$, as $q_t(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) \leq (p_i)_t$ weakly increases with $b_i$, we have

$$q_t(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) \leq q_t(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p}) \leq (p_i)_t.$$

Similarly, we have for any oversampled token $t \in T_-$,

$$q_t(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}) \geq q_t(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p}) \geq (p_i)_t.$$

Then by Definition 2.1,
$$q(b_i', \boldsymbol{b}_{-i}, \boldsymbol{p}) \succeq_i q(b_i, \boldsymbol{b}_{-i}, \boldsymbol{p}),$$

which then implies the monotonicity of $q$. $\qquad\qquad\square$

## Proof of Theorem 3.12

*Proof.* Given an aggregation $q(b_i)$ we construct a stable sampling procedure $\sigma$. Let $Q_+(b_i)$ and $Q_-(b_i)$ be the probability of sampling tokens from $T_+$ and $T_-$:

$$Q_+(b_i) = \sum_{t \in T_+} q_t(b_i), \qquad Q_-(b_i) = \sum_{t \in T_-} q_t(b_i).$$

Both are monotone as $q$ is monotone (by Lemma 3.10). Reparameterize functions $q_t(b_i)$ in the range $I = [Q_+(0), Q_+(\infty)]$ by defining:[6]

$$\hat{q}_t(x) = q_t(Q_+^{-1}(x)), \forall t \in T, x \in I.$$

Since $\hat{q}_t(x)$ is monotone, by Lebesgue's Differentiation Theorem, it is differentiable almost everywhere on $I$. We observe that:

$$\sum_{t \in T_+} \hat{q}_t(x) = Q_+(Q_+^{-1}(x)) = x, \quad \sum_{t \in T_-} \hat{q}_t(x) = Q_-(Q_+^{-1}(x)) = 1 - x.$$

Then $\hat{q}'_t(x)$ for $t \in T_+$ forms a probability distribution over $T_+$. Similarly $-\hat{q}'_t(x)$ for $t \in T_-$ forms a probability distribution over $T_-$. Define $\kappa^+(x), \kappa^-(x) \in \Delta(T)$ as $\kappa_t^+(x) = \hat{q}'_t(x)$ for $t \in T_+$, and zero otherwise and $\kappa_t^-(x) = -\hat{q}'_t(x)$ for $t \in T_-$ and zero otherwise.

We also define the vector $q^+, q^- \in \Delta(T)$ such that $q_t^+ = q_t(0)/Q_+(0)$ for $t \in T_+$ and zero for $t \in T_-$. Similarly: $q_t^- = q_t(\infty)/Q_+(\infty)$ for $t \in T^-$ and zero otherwise.

Finally, we define a deterministic function

$$\mathsf{Sampler} : \Delta(T) \times [0, 1] \to T$$

that takes a probability vector $p \in \Delta(T)$ and $r \in [0, 1]$ and outputs an index $t \in T$ such that $\sum_{j<t} p_j < r \le \sum_{j \le t} p_j$.

Now, we are ready to define the stable sampling procedure. Let $\mathcal{R}$ be the uniform distribution on $[0, 1]^2$. Given $r = (r_A, r_B) \sim \mathcal{R}$ we define the output $t = \sigma(b_i, r)$ as follows:

1. if $r_A \le Q_+(0)$, $t = \mathsf{Sampler}(q^+, r_B) \in T_+$;

2. if $Q_+(0) < r_A \le Q_+(b_i)$, $t = \mathsf{Sampler}(\kappa^+(r_A), r_B) \in T_+$;

3. if $Q_+(b_i) < r_A \le Q_+(\infty)$, $t = \mathsf{Sampler}(\kappa^-(r_A), r_B) \in T_-$;

---

[6]Here $Q_+^{-1}(x)$ refers to a generalized inverse (or the quantile function (Bulow and Roberts, 1989)) so that it is properly defined even when $Q_+(x)$ is discontinuous.

4. if $Q_+(\infty) < r_A$, $t = \mathsf{Sampler}(q^-, r_B) \in T_-$.

Since $\sigma(b_i, r)$ is deterministic, for any fixed $r$, either the output can not be influenced by the bid $(r_A < Q_+(0)$ or $r_A > Q_+(\infty))$ or it can only cause the output to shift from an oversampled token $\mathsf{Sampler}(\kappa^-(r_A), r_B)$ to an undersampled token $\mathsf{Sampler}(\kappa^+(r_A), r_B)$.

We now argue that the tokens are sampled with the correct probabilities. For $t \in T_+$, the total probability of getting sampled is:

$$\int_0^{Q_+(0)} q_t^+ \mathrm{d}r_A + \int_{Q_+(0)}^{Q_+(b)} \hat{q}_t'(r_A)\mathrm{d}r_A = q_t(0) + \hat{q}_t(Q_+(b)) - \hat{q}_t(Q_+(0))$$

$$= q_t(0) + q_t(b) - q_t(0) = q_t(b).$$

Similarly for tokens in $T_-$, the probability of being sampled is:

$$\int_{Q_+(b)}^{Q_+(\infty)} -\hat{q}_t'(r_A)\mathrm{d}r_A + \int_{Q_+(\infty)}^1 q_t^- \mathrm{d}r_A$$

$$= \hat{q}_t(Q_+(b)) - \hat{q}_t(Q_+(\infty)) + q_t(\infty) = q_t(b) - q_t(\infty) + q_t(\infty) = q_t(b).$$

This completes the proof. □

# C   Universally Stable Sampling

**Example C.1** (Counterexample 4-token). Consider two agents $\{1, 2\}$ and 4 tokens $\{t_1, t_2, t_3, t_4\}$. Assume that both agents have the same preferred distribution $p_1 = p_2 = (0, 0, .5, .5)$ and the allocation function is such that if both agents bid zero the allocation is $(.5, .5, 0, 0)$. Hence both both agents have the same set of favored tokens $T_+ = \{t_3, t_4\}$ and less favored tokens $T_- = \{t_1, t_2\}$. The aggregation function $q(b_1, b_2, p_1, p_2)$ is given by the following table:

|           | $b_1 = 0$                  | $b_1 = 1$                  |
|-----------|----------------------------|----------------------------|
| $b_2 = 0$ | $q_{00} = (.5, .5, 0, 0)$  | $q_{10} = (0, .5, .5, 0)$  |
| $b_2 = 1$ | $q_{01} = (.5, 0, .5, 0)$  | $q_{11} = (0, 0, .5, .5)$  |

One can verify that the aggregation function is monotone: When either of the agents increase the bid from 0 to 1, exactly $1/2$ of the probability mass moves from $T_- = \{t_1, t_2\}$ to $T_+ = \{t_3, t_4\}$.

Now we show that there does not exist a universally stable sampling algorithm that implements this aggregation function. Suppose there exist one, $\sigma$, let $r_A = \{r | \sigma(q_{00}, r) = t_1\}$ and $r_B = \{r | \sigma(q_{00}, r) = t_2\}$. Because $\sigma$ is stable for bidder 1, when bidder 1 increases bid, the probability mass only transfers from $T_-$ to $T_+$.

In this case, when the bid profile $(b_1, b_2)$ moves from $(0, 0)$ to $(1, 0)$, we must have

$$\sigma(q_{10}, r) = \begin{cases} t_3, & r \in r_A \\ t_2, & r \in r_B \end{cases}.$$

Further apply the same argument when $(b_1, b_2)$ moves from $(1, 0)$ to $(1, 1)$, we have

$$\sigma(q_{11}, r) = \begin{cases} t_3, & r \in r_A \\ t_4, & r \in r_B \end{cases}.$$

However, if we consider $(b_1, b_2)$ moves along the path $(0, 0) \to (0, 1) \to (1, 1)$, we should have

$$\sigma(q_{01}, r) = \begin{cases} t_1, & r \in r_A \\ t_3, & r \in r_B \end{cases}, \qquad \sigma(q_{11}, r) = \begin{cases} t_4, & r \in r_A \\ t_3, & r \in r_B \end{cases}.$$

We end up with a contradiction on the value of $\sigma(q_{11}, r)$ while moving from $(0, 0)$ to $(1, 1)$ along two different paths.

**Example C.2** (Counterexample 3-token)**.** Consider two agents $\{1, 2\}$ and 3 tokens $\{t_1, t_2, t_3\}$, where the agents have different sets of favored (less favored) tokens. In particular, $T_1^+ = \{t_1, t_3\}$, $T_1^- = \{t_2\}$ and $T_2^+ = \{t_3\}$, $T_2^- = \{t_1, t_2\}$. The aggregation function is given by the following table:

|  | $b_1 = 0$ | $b_1 = 1$ |
|---|---|---|
| $b_2 = 0$ | $q_{00} = (.5, .5, 0)$ | $q_{10} = (.5, 0, .5)$ |
| $b_2 = 1$ | $q_{01} = (0, .5, .5)$ | $q_{11} = (.5, 0, .5)$ |

One can verify that the aggregation function is monotone: When $b_1$ increases from 0 to 1, exactly $1/2$ of the probability mass moves from $t_2$ to $t_3$ (when $b_2 = 0$) or $t_1$ (when $b_2 = 1$). When $b_2$ increases from 0 to 1, either $1/2$ of the probability mass moves from $t_1$ to $t_3$ (when $b_1 = 0$) or no move (when $b_1 = 1$).

Similarly, suppose that there exists a universally stable sampling algorithm $\sigma$ that implements $q$. Let $r_A = \{r | \sigma(q_{00}, r) = t_1\}$ and $r_B = \{r | \sigma(q_{00}, r) = t_2\}$.

Following the same argument in Example C.1, consider the bid profile $(b_1, b_2)$ moves along the path $(0, 0) \to (1, 0) \to (1, 1)$, we must have

$$\sigma(q_{10}, r) = \sigma(q_{11}, r) = \begin{cases} t_1, & r \in r_A \\ t_3, & r \in r_B \end{cases}.$$

However, consider the bid profile $(b_1, b_2)$ moves along the path $(0, 0) \to (1, 0) \to (1, 1)$, we must

have

$$\sigma(q_{01}, r) = \begin{cases} t_3, & r \in r_A \\ t_2, & r \in r_B \end{cases}, \qquad \sigma(q_{11}, r) = \begin{cases} t_3, & r \in r_A \\ t_1, & r \in r_B \end{cases}.$$

We end up with a contradiction on the value of $\sigma(q_{11}, r)$.

# D  Omitted Proofs from Section 4

## Proof of Proposition 4.1

*Proof of Proposition 4.1.* We prove the theorem by showing that the loss $\mathcal{L}_{\mathsf{KL}}^{\bar{\mu}}$ and the loss in equation (4) differ by a constant and hence have the same minimizer. For $B = \sum_i b_i$ and a fixed $x$ we will show that:

$$B \cdot D_{\mathsf{KL}}(\sum_i \tfrac{b_i}{B} \mu_i(\cdot|x) \| f^W(x)) - \sum_i b_i D_{\mathsf{KL}}(f_i(x) \| f^W(x)) = \mathsf{const}.$$

For notation simplicity, we omit the parameters $x$ and $W$ when it is clear from the context and write $\sum_i b_i f_i(x)/B = \bar{f}(x)$. Below we treat any term that doesn't depend on $f$ as a constant:

$$
\begin{aligned}
\sum_i b_i D_{\mathsf{KL}}(f_i \| f) &= \sum_i b_i H(f_i) - \sum_y b_i f_i(y|x) \ln f(y|x) \\
&= -B \sum_y \frac{\sum_i b_i f_i(y|x)}{B} \cdot \ln f(y|x) + \sum_i b_i H(f_i) \\
&= -B \cdot H(\bar{f}, f) + \sum_i b_i H(f_i) \\
&= B \cdot D_{\mathsf{KL}}(\bar{f} \| f) - B \cdot H(\bar{f}) + \sum_i b_i H(f_i) \\
&= B \cdot D_{\mathsf{KL}}(\bar{f} \| f) - \mathsf{const}.
\end{aligned}
$$

To complete the proof, observe that if $f_i$ is the unconstrained minimizer of $\mathcal{L}_{\mathsf{KL}}^{\mu_i}(f)$ we must have $f_i(y|x) = \mu_i(y|x)$. $\qquad\square$

## Proof of Lemma 4.2

*Proof of Lemma 4.2.* Let $B = \sum_i b_i$ and consider $\mathsf{LOSS}_{\mathsf{KL}}$:

37

$$\text{Loss}_{\text{KL}} = \sum_{i \in [n]} b_i \cdot (H(p_i, q) - H(p_i))$$

$$= -\sum_{i \in [n]} b_i \cdot \sum_{t \in [T]} p_i(t) \ln q(t) - \sum_{i \in [n]} b_i \cdot H(p_i)$$

$$= -\sum_{t \in [T]} \ln q(t) \sum_{i \in [n]} b_i \cdot p_i(t) - \sum_{i \in [n]} b_i \cdot H(p_i)$$

$$= -B \cdot \sum_{t \in [T]} \frac{\sum_{i \in [n]} b_i \cdot p_i(t)}{B} \ln q(t) - \sum_{i \in [n]} b_i \cdot H(p_i)$$

$$= \sum_{i \in [n]} B \cdot H(q_{\text{KL}}, q) - b_i \cdot H(p_i).$$

By Gibbs' inequality, the cross entropy $H(q_{\text{KL}}, q)$ is minimized if and only if $q = q_{\text{KL}}$. Hence this is also the minimizer of $\text{Loss}_{\text{KL}}$. $\qquad\square$

## Proof of Proposition 4.3

*Proof.* For a fixed $x$, $f^*(y|x)$ can be obtained by solving:

$$\max_{f(\cdot|x)} \sum_y f(y|x)\bar{r}(x, y) - \beta D_{\text{KL}}(f(x) \| f^{\text{SFT}}(x))$$

$$\text{s.t.} \sum_y f(y|x) = 1 \text{ and } f(y|x) \geq 0.$$

By the standard KKT conditions, the solution has the form:

$$f^*(y|x) = f^{\text{SFT}}(y|x) e^{\bar{r}(x,y)/\beta} C_x.$$

where $C_x$ is a normalization constant to ensure $\sum_y f^*(y|x) = 1$. For the same reason, we have that:

$$f_i(y|x) = f^{\text{SFT}}(y|x) e^{r_i(x,y)/\beta} C_{i,x}.$$

If $f^\circ$ is the function minimizing problem (5), then we can apply KKT conditions to obtain:

$$f^\circ(y|x) = \exp\left(\frac{1}{B} \sum_i b_i \ln f_i(y|x)\right) \cdot C_x'$$

for normalization constants $C'_x$ and $B = \sum_i b_i$. Replacing the formula for $f_i(x)$ from the previous line, we obtain that:

$$f^\circ(y|x) = \exp\left(\frac{1}{B}\sum_i b_i\left(r_i(x,y)/\beta + \ln f^{\mathsf{SFT}}(y|x)\right)\right) \cdot C''_x$$

$$= \exp\left(\bar{r}(x,y)/\beta + \ln f^{\mathsf{SFT}}(y|x)\right) \cdot C''_x = f^*(y|x).$$

This completes the proof. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$