

# A Truth Serum for Large-scale Evaluations

Vijay Kamble<sup>1</sup>, David Marn<sup>3</sup>, Nihar Shah<sup>2</sup>, Abhay Parekh<sup>3</sup>, and Kannan Ramachandran<sup>3</sup>

<sup>1</sup>University of Illinois at Chicago, *kamble@uic.edu*

<sup>2</sup>Carnegie Mellon University, *nihars@cs.cmu.edu*

<sup>3</sup>University of California, Berkeley, *{marn, parekh, kannanr}@eecs.berkeley.edu*

May 10, 2018

## Abstract

A major challenge in obtaining large-scale evaluations of products or services in reputation systems is that of eliciting honest responses from agents in the absence of verifiability. We propose a new reward mechanism with strong incentive properties applicable in a wide variety of such settings. This mechanism has a simple and intuitive output agreement structure: an agent gets a reward only if her response for an evaluation matches that of her peer. But instead of the reward being the same across different answers, it is inversely proportional to a popularity index of each answer. This index is a second order population statistic that captures how frequently two agents performing the same evaluation agree on the particular answer. Rare agreements thus earn a higher reward than agreements that are relatively more common.

In the regime where there are a large number of evaluation tasks, we show that truthful behavior is a strict Bayes-Nash equilibrium of the game induced by the mechanism. Further, we show that the truthful equilibrium is approximately optimal in terms of expected payoffs to the agents across all symmetric equilibria, where the approximation error vanishes in the number of evaluation tasks. Moreover, under a mild condition on strategy space, we show that any symmetric equilibrium that gives a higher expected payoff than the truthful equilibrium must be close to being fully informative if the number of evaluations is large. These last two results are driven by a new notion of an agreement measure that is shown to be monotonic in information loss. This notion and its properties are of independent interest.

## 1 Introduction

Feedback and reputation systems, in which people provide ratings and reviews for products or services based on their personal experiences, are a critical component of online platforms and marketplaces [Luc17]. These systems improve the overall quality of transactions, increase trust, and thus play a major role in determining the success of these platforms in the long run. A major practical challenge in these systems is that of eliciting truthful and high-quality responses from the agents. In the absence of appropriate incentives, agents could shirk investing effort, provide uninformative feedback, or could even try to exploit these systems for selfish motives, thus undermining their utility. For instance, several recent works have found significant empirical evidence of bias in user ratings on many online platforms [HG15, NT15]. In this work, we design a new and simple reward mechanism with strong incentive properties that attempts to address this concern in a wide variety of situations.

We consider a general setup, where a principal is interested in obtaining responses for a large number of evaluation tasks from a pool of agents. Our central assumption is that the population of agents is statistically *homogeneous*, which means that in each evaluation task, a) the true response of each evaluating agent is an independent sample from an unknown distribution of answers that is common across agents and b) this response is independent of any personally known individual characteristics or preferences of the agent. This is the case for evaluations comprising of questions like:

1. What was your waiting time to get a table in the restaurant? (Less than 15 mins/Between 15-30 mins/More than 30 mins)
2. Did the handyman show up on time for your appointment? (Yes/No)
3. Did the received product match the description given by the seller? (Yes/No)

In the first situation, we can assume that each customer experiences an independently sampled waiting time from a common unknown distribution that is specific to that restaurant. In the second situation, the customer’s experience is sampled from the distribution of whether or not the handyman is punctual. Similarly, in the third question, the customer’s true experience is a sample of the selling practices of the seller. In all three cases, the true responses are independent of the individual characteristics of the customers.

When the principal either has access to the true answers for some of the evaluation tasks or can verify the answers accurately, she can score the agents based on their performance on these tasks. Proper scoring rules [Bri50, GR07, Sav71, LS09] provide a precise and elegant framework to induce truthfulness in such cases. This approach, coupled with the anonymity of these special tasks, incentivizes truthful behavior in all evaluations as long as the threat of being scored on a verifiable task is high enough relative to the incentive to shirk effort or to lie [GWL16].

On the other hand, such an approach is impractical in cases where the true responses are either intrinsically unknowable to the principal because they are based on personal experiences (as is the case for the questions above), or are difficult to obtain in significant numbers. In such cases, the alternative is to score the agents’ responses based on comparisons with the responses of other agents who have provided answers for the same or similar questions [MRZ05]. The situation is then inherently strategic, in which one hopes to sustain truthful reporting as an equilibrium of a game that the scoring mechanism induces. The present work falls into this category.

In these settings, information elicitation mechanisms typically leverage the property that the true response of any agent is correlated with the response of some other agent for the same question. With a homogeneous population, the structure of this correlation between an agent’s response and the response of a typical agent in the population is identical across agents. This contrasts with the case when the population is *heterogeneous*, i.e., when the agents’ responses strongly depend on their characteristics that vary widely across the population. For instance, consider the question: Did you like this movie? The answer to this question would depend on the personal preferences of the agent – one agent could strongly prefer action movies, while some other agent could strongly prefer dramas. If a majority of the population is expected to be comprised of people who like action movies, then an action lover’s response is positively correlated with the typical response, while a drama lover’s response is negatively correlated with the typical response. Moreover, agents typically know their type and this fact. This creates problems for mechanisms that try to incentivize both types of agents to be truthful without knowing their type. In these cases, designing mechanisms that

induce truthfulness is inherently harder, and without obtaining requisite fine-grained information about agent heterogeneity, even impossible [RF15].

In the homogeneous population setting, the pioneering work [MRZ05] described the *Peer Prediction* mechanism that incentivizes truthful answers to a single evaluation task. The main requirement is that there is a commonly held prior on the unknown distribution from which the agents’ answers are sampled, and this prior is known to the principal. Truthfulness is achieved by scoring an agent’s prediction of her peer’s answer, as implied by her own answer, using a proper scoring rule (hence the name). Another influential design in this domain, the Bayesian Truth Serum (BTS) [Pre04], and its subsequent refinements [WP12b, RF13], preserved the common prior assumption but relaxed the requirement that the principal needs to know this prior.<sup>1</sup> But these mechanisms instead require the agents to make extraneous reports about their beliefs in addition to their answers. Such extraneous reports of beliefs, although undesirable, are indispensable in this setting; it has been shown that it is impossible to design mechanisms that incentivize truthfulness without obtaining some information about the agents’ belief structure [JF11]. Mechanisms that do not assume that the principal knows the details of the agents’ belief structure a priori have been referred to as *detail-free* in literature.

A key feature of online platforms is that they host a large number of similar products or services, although every buyer or customer interacts with only subset of these entities. For instance, there are thousands of similar restaurants (with similar existing ratings) listed on review platforms like Yelp that users rate. Online marketplaces like Amazon or eBay would like to obtain reviews for a large number of existing sellers on these platforms. Online labor platforms like Thumbtack and Handy would like to obtain performance metrics for thousands of workers that operate on these platforms.

The presence of multiple similar evaluation tasks hints at an approach for designing detail-free mechanisms that do not require extraneous belief reports from agents (mechanisms that do not require such extraneous reports have been called *minimal* in literature). This information can instead be replaced by consistent statistical estimates of the distribution of agent responses obtained from the response data across multiple tasks. Recently, several designs have successfully exploited this possibility [JF11, WP13, DG13, SAFF16, RF15, RFJ16, LC17]; some even under heterogeneous population settings. A potential concern with this approach is that it assumes that the responses are truthful. But it turns out that in many of these situations, truthfulness becomes a self-fulfilling prophecy – truthful behavior is an equilibrium in the induced game when the mechanism simply assumes that these reports are truthful. This is the basic principle underlying our design.

Our mechanism builds upon the structure of *output agreement mechanisms* [VAD08, VAD04] that are simple, intuitive, and have been quite popular in practice, except they suffer from a critical drawback of not incentivizing truthful responses in general. In an output agreement mechanism, two agents answer the same question, and they are both rewarded if their answers match. Under such a scheme, the agents tend to gravitate towards answers that are more likely to be popular rather than submitting their true responses. Our mechanism overcomes this drawback by giving proportionately lower rewards for answers that turn out to be more popular. This is achieved by inversely scaling the rewards for agreement by a *popularity index* for each answer. This is not a new idea: such biased output agreement schemes have been explored earlier in prior works; for instance, a well-known scheme of this type is the ‘Peer Truth Serum’ [JF11, RFJ16].

The key innovation in our design relative to these works is in the way these popularity indices are defined. In fact, all the strong incentive properties of our mechanism trace their origin to this

---

<sup>1</sup>[WP12a] later relaxed the common prior assumption for the case of binary evaluations.

novel definition. In our mechanism, these indices are certain second order population statistics that capture how frequently two people performing the same task agree on a particular answer on average across all tasks. Formally, the popularity index of an answer is the square root of the estimate for the probability of agreement on that answer obtained from response data. Thus rare agreements receive higher rewards than agreements that are relatively common. As the number of tasks increases, the accuracy of these indices improves and truthfulness is obtained as a Bayes-Nash equilibrium when the number of tasks is large enough. A common prior is not necessary for this result; it should just be common knowledge that the prior satisfies a certain non-degeneracy property (this property is related to the ‘stochastic relevance’ condition that appears frequently in this domain; for instance, in [MRZ05]).

An important goal for any reward mechanism is to *strictly* incentivize truthfulness. That is, it is not sufficient to ensure that truthful behavior is simply an equilibrium in the resulting game. It is also important to ensure that at equilibrium, each agent gets a strictly higher reward by being truthful than by adopting any other strategy. Mechanisms that achieve this are referred to as *strictly proper* in literature. Without this requirement, trivial mechanisms like the one that gives a fixed payment to each agent already, in principle, weakly incentivize truthfulness. Moreover, it is also important to ensure that at the truthful equilibrium, the difference in the payoffs to an agent resulting from the truthful strategy and any strategy in which an agent’s reported response is independent of her true response, is bounded away from zero. Such strict incentives allow the principal to account for any costs that the agents may incur for their evaluation effort by appropriately scaling the rewards in the mechanism.

Another important concern in these mechanisms is that the induced game may possess multiple equilibria. In such cases, there needs to be an adequate rationale for the truthful equilibrium to be selected. This issue has increasingly been brought into focus recently. It is known that mechanisms for a single task with no extraneous reporting (which includes [MRZ05]) possess uninformative equilibria that give a higher expected payoff to each agent than in the truthful equilibrium [JF05]. The Bayesian Truth Serum demonstrated that this issue can be overcome by requiring extraneous reports of beliefs; the truthful equilibrium under BTS gives the highest expected payoff to an agent across all equilibria. Mechanisms that satisfy this property have been referred to as *focal* in literature.

A few recent works have discovered that such strong truthfulness properties can be obtained in the multi-task setting without requiring extraneous reports from the agents [DG13, SAFF16, KS16]. [KS16] expounds a general information theoretic analysis of incentive mechanisms in this space, and has shown that most mechanisms achieve such properties by (implicitly or explicitly) connecting the loss in the agents’ expected payoff relative to the truthful equilibrium, to some form of mutual information loss or correlation loss in the population due to deviation from truthfulness. Both, extraneous reports of beliefs about the play, or statistical estimates of the play across multiple tasks, in some sense allow these mechanisms to obtain a handle on an appropriate mutual information measure.

Addressing these concerns, we show that under the non-degeneracy assumptions that we mentioned earlier, truthful behavior is a strict Bayes-Nash equilibrium under our mechanism, i.e., our mechanism is strictly proper. Also, at this equilibrium, the difference in the payoffs to an agent under the truthful strategy and under any strategy in which an agent’s reported response is independent of her true response, is bounded away from zero.

Moreover, we obtain a vanishing uniform upper bound (in the number of tasks) on the difference between the expected payoff obtained by the truthful equilibrium and that obtained under any other

symmetric equilibrium (equilibrium in which all players choose the same strategy). Asymptotically, the limiting expected payoff under a truthful strategy profile is higher than that under any other symmetric strategy profile. Under a mild assumption on the strategy spaces, we also show that any symmetric equilibrium that gives the highest reward to an agent across all symmetric equilibria must be close to being *fully informative* (in a precise sense) when the number of tasks is large. A fully informative strategy is simply a permutation map from the set of answers to itself; this includes truthfulness. In line with the framework in [KS16], these results are obtained by showing that the expected payoff of an agent under any symmetric strategy profile is connected to a novel notion of an *agreement measure* between two independent responses that is monotonic in information loss. Both the measure and this monotonicity property are of independent interest.

**Comparison with other mechanisms.** Our mechanism adds to a growing set of elicitation mechanisms that lie on the frontier of the various tradeoffs in verification-free elicitation. Table 1 summarizes the differences between different popular mechanisms at a high-level. The following are some key distinctions.

1. The mechanisms proposed in both [SAFP16] (also known as the *Correlated Agreement* mechanism) and [KS16] require each agent to perform a large number of tasks, a large number of which have to be shared in common with some other agent. While suitable for settings like crowdsourcing where the principal controls task allocations and every agent is expected to perform a large number of evaluations, this requirement is impractical on online platforms, where agents typically evaluate few and disparate products/services of their choice. Our mechanism on the other hand only requires a large number of evaluation tasks in *total* – each agent can perform as few as a single task as long as each task is performed by at least two agents (which is typically satisfied in practice). On the flip side, the incentive properties of our mechanism only hold under the homogeneity assumption, while these mechanisms are applicable in general. Moreover these mechanisms are focal (asymptotically in the number of tasks), where as our mechanism is focal only across symmetric equilibria. These are essentially the costs of not requiring the agents perform a large number of tasks; we will have more comments on this point in Section 4.
2. The mechanism described in [DG13], while not requiring each agent to perform a large number of tasks, is truthful and focal only if the agent responses are “categorical”, which means that conditional on an agent’s answer, the posterior probability of all other answers from another agent must reduce relative to the prior. That is, if  $Pr(y)$  is the prior probability of an answer  $y$  and  $Pr(y|y')$  is the conditional probability of some other agent has answer  $y$  given that one agent has answer  $y'$ , then  $Pr(y'|y) \leq Pr(y')$  for all  $y' \neq y$ . Except for the case of binary evaluations, this condition is not satisfied in general in the homogeneous population setting.
3. Similar to our mechanism, the Peer Truth Serum [JF11, RFJ16] has a biased output agreement structure where the popularity score of each answer is the estimate of the prior probability of an agent reporting that answer (in our definition, it is the square root of the estimate of the probability of agreement on that answer). Although this mechanism is not focal in general (for reasons similar to ours), by an appropriate choice of parameters, it can be made focal across symmetric equilibria. But in order to obtain these properties, it requires that the agent responses satisfy a “self-predicting” condition, which says that  $Pr(y|y)/Pr(y) \geq Pr(y'|y)/Pr(y')$  for any  $y' \neq y$ . Again, except for the case of binary evaluations, this condition is also not satisfied in general in the homogeneous population setting.

The remainder of the paper is organized as follows. Section 2 presents a formal description of

Mechanism	Detail-Free	Minimal	Evaluations per person	Conditions for truthfulness	Focal
Peer Prediction Method	No	Yes	1	Homogeneous settings, common prior	No
Bayesian Truth Serum	Yes	No	1	Homogeneous settings, common prior	Yes
Correlated Agreement	Yes	Yes	Large	Any setting	Yes
[KS16]	Yes	Yes	Large	Any setting	Yes
Peer Truth Serum	Yes	Yes	1	Responses are “self-predicting”	Across symmetric equilibria
[DG13]	Yes	Yes	2	Responses are “categorical”	Yes
Our	Yes	Yes	1	Homogeneous settings	Across symmetric equilibria

Table 1: Properties of different mechanisms.

the model considered in the paper. Section 3 presents our mechanism and our main results. We discuss some comparisons and relations with existing mechanisms in greater detail in Section 4. We finally summarize our contributions and conclude in Section 5. The proofs of all of our results can be found in the Appendix.

## 2 Model

We consider a setting with  $N$  evaluation tasks by the set  $\mathcal{N}$  and labeled as  $i = 1, \dots, N$ . Let  $\mathcal{M}$  denote the population of agents, labeled as  $j = 1, \dots, M$ . Let  $\mathcal{M}_i \subseteq \mathcal{M}$  denote the subset of agents that perform task  $i$ , and let  $\mathcal{N}_j$  be the set of tasks that an agent  $j$  performs. We assume that the sets  $\mathcal{M}_i$  and  $\mathcal{N}_j$  are exogenously specified. The set of answers in each evaluation task is assumed to be finite and denoted as  $\mathcal{Y}$ . The unknown distribution of answers in population  $\mathcal{M}_i$  is parameterized by a type  $X_i$  that takes values in the set  $\mathcal{X}$ , also finite, for all populations. The distribution of the answers in the population as a function of type  $x \in \mathcal{X}$  is denoted as  $\mathbf{p}(x) = (p_y(x); y \in \mathcal{Y})$ . The answer  $Y_j^i$  of each agent  $j$  in  $\mathcal{P}_i$  is independently drawn from  $\mathbf{p}(X_i)$ . Further, the types of different tasks are independently sampled from a common distribution  $P_X$ . We further assume that the different answers of any person  $j$  for all the evaluation tasks in  $\mathcal{N}_j$  are mutually independent.<sup>2</sup> The answers of different agents for a single task  $i$  need not be independent unless conditioned on  $X_i$ . Finally, we assume that from the perspective of any agent  $j$ , there are no other observable features of the evaluation task  $i$  except the true answer  $Y_j^i$ .<sup>3</sup>

The probability distribution over types,  $P_X$ , and the function  $\mathbf{p}$  together form a *generating model*, denoted as the pair  $(P_X, \mathbf{p})$ . In particular, this pair fully specifies a joint distribution on the underlying types of the different evaluation tasks and the answers of the different agents across tasks.

Our goal is to design a payment mechanism that elicits true responses from the agents. We are specifically interested in the case where  $N$  is large, although  $\mathcal{N}_j$  is relatively small for each  $j$ . The principal is not assumed to know  $P_X$  or  $\mathbf{p}$ . We do not make any additional assumptions on the agents’ knowledge of  $(P_X, \mathbf{p})$  other than the assumption that every agent knows the structure of the underlying generating model, i.e., the existence of some  $P_X$  and the function  $\mathbf{p}$  that is common across tasks. In particular, this means that all the agents know that all the tasks are statistically similar and the population of agents performing each evaluation is homogeneous.

**Example 2.1.** Consider a situation where a labor platform wants to obtain feedback on punctuality

<sup>2</sup>This assumption precludes the possibility of dependence induced by lack of knowledge of some hidden information about an agent, e.g., if an agent’s mood is bad on a particular day, there may be a bias in all her evaluations.

<sup>3</sup>The presence of extraneous features in the evaluation tasks can introduce equilibria of the form “If the background color of image is green, then report ‘cat’, otherwise report truthfully”. These type of equilibria trouble most payment mechanisms in this domain; it has been shown that it is impossible to elicit every feature under a payoff-dominant truthful equilibrium [GWL16]

of handymen that operate on the platform. In this case, each handyman could be of 2 types  $\mathcal{X} = \{\text{Punctual}, \text{Not Punctual}\}$ , with  $P_X(\text{Punctual}) = P_X(\text{Not Punctual}) = 0.5$ . The question is “Did the handyman show up on time for your appointment?”. The two answers are  $\mathcal{Y} = \{\text{Yes}, \text{No}\}$ . And the distributions of answers as a function of type are  $\mathbf{p}(\text{Punctual}) = (0.95, 0.05)$  and  $\mathbf{p}(\text{Not Punctual}) = (0.5, 0.5)$ .

Note that the above example is simply illustrative. For obtaining our results, it doesn't matter what the specifics of the generating model are, as long as every agent believes that there is a generating model with this structure.

An agent  $j$ 's strategy is a set of mappings  $\{\mathbf{q}^{ij} : i \in \mathcal{N}_j\}$  where  $\mathbf{q}^{ij}(y) = (q_{y'}^{ij}(y); y' \in \mathcal{Y})$  is the probability distribution over answers for evaluation task  $i \in \mathcal{N}_j$  conditional on the true response being  $y$ . We do not consider reporting strategies in which the reported answer of an evaluation depends not just on the true answer for that evaluation, but also on the true answers to all the other evaluations that the agent performs. This restriction is simply for the ease of exposition. It will become clear that all the incentive properties continue to hold for our proposed mechanism even if such strategies are allowed. For any  $j \in \mathcal{M}$ , let  $r_j^i$  denote agent  $j$ 's reported answer for task  $i$ . We define a payment mechanism as follows.

**Definition 2.1.** A payment (or reward/scoring) mechanism is a set of functions  $\{\tau_j : j \in \mathcal{M}\}$ , one for each person in the population, that map the reports  $\{r_j^i : j \in \mathcal{M}, i \in \mathcal{N}_j\}$  to a real valued payment (or score).

We define the notion of a Bayes-Nash equilibrium in the game induced by a payment mechanism. The definition assumes that the generating model is commonly known to the agents, but this assumption will be appropriately relaxed later.

**Definition 2.2.** Given a generating model  $(P_X, \mathbf{p})$  that is common knowledge amongst the agents, we say that a strategy profile  $\{\mathbf{q}^{ij} : j \in \mathcal{M}, i \in \mathcal{N}_j\}$  comprises a Bayes-Nash equilibrium in the game induced by the payment mechanism if for each  $j \in \mathcal{M}$ ,

$$\begin{aligned} & \mathbb{E} \left[ \tau_j(\{\mathbf{q}^{ij'}(Y_{j'}^i) : j' \in \mathcal{M}, i \in \mathcal{N}_{j'}\}) \right] \\ & \geq \mathbb{E} \left[ \tau_j(\{\bar{\mathbf{q}}^{ij}(Y_j^i) : i \in \mathcal{N}_j\} \cup \{\mathbf{q}^{ij'}(Y_{j'}^i) : j' \in \mathcal{M}, j' \neq j, i \in \mathcal{N}_{j'}\}) \right], \end{aligned} \quad (1)$$

for each  $\{\mathbf{q}^{ij} : i \in \mathcal{N}_j\} \neq \{\bar{\mathbf{q}}^{ij} : i \in \mathcal{N}_j\}$ , where the expectation is with respect to the joint distribution on the responses of the population specified by the generating model  $(P_X, \mathbf{p})$ . We say that the strategy profile is a strict Bayes-Nash equilibrium if the above inequality is strict.

In words, this says that assuming all the other agents adhere to the reporting strategy profile  $\{\mathbf{q}^{ij} : j \in \mathcal{M}, i \in \mathcal{N}_j\}$ , each agent maximizes her expected reward by also adhering to the prescriptions of the strategy profile. Next, we define Bayes-Nash incentive compatibility, which is the property that truthful reporting is a Bayes-Nash equilibrium of the game induced by the reward mechanism.

**Definition 2.3.** We say that a payment mechanism is Bayes-Nash incentive compatible with respect to the generating model  $(P_X, \mathbf{p})$  if the truthful reporting strategy profile, i.e., where  $q_{y'}^{ij}(y) = \mathbf{1}_{\{y'=y\}}$  for all  $j \in \mathcal{M}$  and  $i \in \mathcal{N}_j$ , is a Bayes-Nash equilibrium. If this equilibrium is strict, we say that the mechanism is strictly Bayes-Nash incentive compatible.

### 3 Main results

The simplest description of the core idea of our mechanism is obtained in the case of a single evaluation task, in which the generating model  $(P_X, \mathbf{p})$  is *commonly known* to the principal and the agents (this is the setting in [MRZ05]). In this setting, consider the following output agreement scheme. Each agent  $j$  is paired with another randomly chosen agent  $j'$  and their responses are compared. If their responses don't match, then  $j$  gets no reward. If their responses match and this common response is  $y \in \mathcal{Y}$ , then  $j$  gets a reward  $K/\sqrt{P(Y_j = Y_{j'} = y)} = K/\sqrt{\sum_{x \in \mathcal{X}} P_X(x)p_y^2(x)}$ , where  $K$  is some positive constant. That is, she gets a reward that is inversely proportional to the square root of the probability that both the agents form the same response  $y$ . Or to put it simply, *the reward for an agreement is inversely proportional to the square root of the probability of that agreement*. Thus an agreement that is more probable earns a lower reward than an agreement that is relatively less probable.

Let us see why truthful behavior is a Bayes-Nash equilibrium in this mechanism. Consider an agent  $j$ , and suppose that all other agents are truthful. Then if  $j$ 's true response is  $y$ , her expected reward for a truthful report is,

$$K \frac{P(Y_{j'} = y | Y_j = y)}{\sqrt{P(Y_j = Y_{j'} = y)}} = K \frac{\sqrt{P(Y_{j'} = Y_j = y)}}{P(Y_j = y)}. \quad (2)$$

Similarly, her reward for any other report  $y'$  is,

$$K \frac{P(Y_{j'} = y' | Y_j = y)}{\sqrt{P(Y_j = Y_{j'} = y')}} = K \frac{P(Y_{j'} = y', Y_j = y)}{P(Y_j = y)\sqrt{P(Y_j = Y_{j'} = y')}}. \quad (3)$$

Thus being truthful gives a higher reward if,

$$\sqrt{P(Y_{j'} = Y_j = y)}\sqrt{P(Y_j = Y_{j'} = y')} \geq P(Y_{j'} = y', Y_j = y),$$

i.e., if,

$$\sqrt{\sum_{x \in \mathcal{X}} P_X(x)p_y(x)^2} \sqrt{\sum_{x \in \mathcal{X}} P_X(x)p_{y'}(x)^2} \geq \sum_{x \in \mathcal{X}} P_X(x)p_y(x)p_{y'}(x). \quad (4)$$

But this is precisely the Cauchy-Schwarz inequality.

Now there are three concerns. First, we need truthfulness to be a strict Bayes-Nash equilibrium. This means that we need the above inequality to be strict. It turns out that this is true under natural assumptions on the generating model. Second, in our original setting, the principal does not know the generating model. This can be addressed by replacing the required agreement probabilities by consistent statistical estimates computed from reports obtained across multiple tasks. As the number of tasks grows, the accuracy of these estimates improves and for a large enough  $N$ , truthfulness is recovered as a strict equilibrium. Finally, we need to tackle the issue of multiple equilibria. We elaborate on all these aspects one by one.

#### 3.1 Obtaining strictness

For truthfulness to be a strict equilibrium, we need the Cauchy-Schwarz inequality in (4) to be strict for every pair  $y, y' \in \mathcal{Y}$ . It will be useful to define the following ‘‘inequality gap’’.



**Definition 3.1.** For a generating model  $(P_X, \mathbf{p})$  defined on  $\mathcal{X}$  and  $\mathcal{Y}$ , define

$$\delta(P_X, \mathbf{p}) = \min_{y, y' \in \mathcal{Y}, y \neq y'} \sqrt{\left( \sum_{x \in \mathcal{X}} P_X(x) p_y(x)^2 \right) \left( \sum_{x \in \mathcal{X}} P_X(x) p_{y'}(x)^2 \right) - \sum_{x \in \mathcal{X}} P_X(x) p_y(x) p_{y'}(x)}.$$

By the Cauchy-Schwarz inequality,  $\delta(P_X, \mathbf{p}) \geq 0$ . If  $\delta(P_X, \mathbf{p}) > 0$  for some generating model  $(P_X, \mathbf{p})$ , then truthfulness is a strict Nash equilibrium in the game induced by the mechanism we described earlier for the case where the principal knows this generating model. To understand whether this is a reasonable assumption, a little demystification of this condition is in order.

For any answer  $y \in \mathcal{Y}$ , define the vector

$$\mathbf{v}(y) \triangleq (\sqrt{P_X(x)} p_y(x); x \in \mathcal{X}). \quad (5)$$

Then the Cauchy-Schwarz inequality says that for any two answers  $y$  and  $y'$ , the magnitude (in the Euclidean norm) of the projection of the vector  $\mathbf{v}(y)$  on the unit vector in the direction  $\mathbf{v}(y')$  is less than the magnitude of the vector  $\mathbf{v}(y)$  itself (one can reverse the roles of  $y$  and  $y'$ ), i.e.,

$$\frac{|\mathbf{v}(y) \cdot \mathbf{v}(y')|}{\|\mathbf{v}(y)\|} \leq \|\mathbf{v}(y')\|,$$

or

$$|\mathbf{v}(y) \cdot \mathbf{v}(y')| \leq \|\mathbf{v}(y)\| \|\mathbf{v}(y')\|. \quad (6)$$

Let  $\theta(\mathbf{u}, \mathbf{v})$  denote the angle in radians between two non-zero vectors  $\mathbf{u}$  and  $\mathbf{v}$ , defined as

$$\theta(\mathbf{u}, \mathbf{v}) \triangleq \arccos \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}. \quad (7)$$

Then it is clear that the inequality in (6) is strict if and only if the angle between the vectors  $\mathbf{v}(y)$  and  $\mathbf{v}(y')$  is positive and their magnitude is non-zero. In fact, under the condition that  $\|\mathbf{v}(y)\| \leq 1$  for all  $y \in \mathcal{Y}$ , which holds in our case, we can show that the gap in (6) is bounded away from 0 if and only if the angle between the vectors  $\mathbf{v}(y)$  and  $\mathbf{v}(y')$ , and their magnitudes are all bounded away from 0. The following proposition gives a precise statement.

**Proposition 1.** For a generating model  $(P_X, \mathbf{p})$  defined on  $\mathcal{X}$  and  $\mathcal{Y}$ , the following two conditions are equivalent.

1. There is some  $\alpha > 0$  such that  $\delta(P_X, \mathbf{p}) > \alpha$ .

2. There is some  $\tau > 0$  such that

- (a)  $\sum_{x \in \mathcal{X}} P_X(x) p_y(x) > \tau$  for each  $y \in \mathcal{Y}$ , and,
- (b)  $\theta(\mathbf{v}(y), \mathbf{v}(y')) > \tau$  for all  $y, y' \in \mathcal{Y}$  such that  $y \neq y'$ .

Condition 2(a) says that the probability of an agent forming any response  $y \in \mathcal{Y}$  is bounded away from zero. This condition is naturally satisfied in practice. If this doesn't hold for some answer, then one can simply eliminate that answer from the admissible set.

Condition 2(b) says that the angle between  $\mathbf{v}(y)$  and  $\mathbf{v}(y')$  is bounded away from zero for any  $y \neq y'$ . If this is not true for some  $y$  and  $y'$  then there is a  $C \in \mathbb{R}$  such that  $p_y(x) = C p_{y'}(x)$  for each  $x \in \mathcal{X}$  such that  $P_X(x) > 0$ . But in this case, the responses  $y$  and  $y'$  need not be distinguished

at all, since they contain the same information about  $X$ , and hence about the rest of the random quantities. In particular,  $P(X_i = x|Y_j^i = y) = P(X_i = x|Y_j^i = y')$  for each  $x \in \mathcal{X}$ . Hence, the principal can simply ask the agents to map both these responses to a single response.

In the context of our model, the condition  $\theta(\mathbf{v}(y), \mathbf{v}(y')) > 0$  for any  $y \neq y'$  is weaker than the stochastic relevance condition that is imposed to obtain strictness in several works in this domain, starting from [MRZ05]. An agent's answer to a question is a stochastically relevant random variable if no two answers induce the same conditional distribution on the answers of some other agent who has answered the same question. Clearly, if  $\theta(\mathbf{v}(y), \mathbf{v}(y')) = 0$ , then stochastic relevance is violated, and thus stochastic relevance implies that  $\theta(\mathbf{v}(y), \mathbf{v}(y')) > 0$ .

### 3.2 Relaxing the knowledge assumption: the large $N$ regime

Next, we tackle the problem of the principal not knowing the generating model. In this case, the principal can simply replace the parameters of the mechanism by consistent statistical estimates obtained from the reports across multiple tasks, assuming these reports are truthful. We thus present our main mechanism.

**Definition 3.2 (Main mechanism. Assumes  $\mathcal{M}_i \geq 2$  for all  $j \in \mathcal{M}$ ).** *The responses of agents for the different evaluation tasks are solicited. Let these be denoted by  $\{r_j^i : j \in \mathcal{M}, i \in \mathcal{N}_j\}$ . An agent  $j$ 's payment is computed as follows:*

- For each population  $\mathcal{M}_i$  such that  $i \notin \mathcal{N}_j$ , choose any two agents  $j_1(i), j_2(i) \in \mathcal{M}_i$ , and for each possible evaluation  $y \in \mathcal{Y}$ , compute the quantity

$$\bar{f}_j(y) = \frac{1}{N - |\mathcal{N}_j|} \sum_{i \in \mathcal{N} \setminus \mathcal{N}_j} \mathbf{1}_{\{r_{j_1(i)}^i = y\}} \mathbf{1}_{\{r_{j_2(i)}^i = y\}}.$$

- For each answer  $y$ , fix a payment  $e_j(y)$  defined as

$$e_j(y) = \begin{cases} \frac{K}{\sqrt{\bar{f}_j(y)}} & \text{if } \bar{f}_j(y) \neq 0, \\ 0 & \text{if } \bar{f}_j(y) = 0, \end{cases}$$

where  $K > 0$  is any positive constant.  $\sqrt{\bar{f}_j(y)}$  is the popularity index of answer  $y$ .

- For computing agent  $j$ 's payment for evaluation task  $i \in \mathcal{W}_j$ , choose another agent  $j' \in \mathcal{M}_i$ , who will be called  $j$ 's peer for task  $i$ . If their responses match, i.e., if  $r_j^i = r_{j'}^i = y$ , then  $j$  gets a reward of  $e_j(y)$ . If the responses do not match, then  $j$  gets 0 payment for that task.

Observe that if everyone except agent  $j$  is truthful, then  $E[\bar{f}_j(y)] = P[Y_{j_1(i)}^i = Y_{j_2(i)}^i = y]$ , i.e.,  $\bar{f}_j(y)$  is a consistent estimate of  $P[Y_{j_1(i)}^i = Y_{j_2(i)}^i = y]$ . In fact, as  $N$  grows large, assuming  $|\mathcal{N}_j|$  remains bounded,  $\bar{f}_j(y)$  almost surely converges to  $P[Y_{j_1(i)}^i = Y_{j_2(i)}^i = y]$  by the strong law of large numbers. For an  $N$  large enough, we can show that the quality of the estimate is sufficiently high to ensure that truthfulness is sustained as a strict equilibrium under appropriate conditions on the generating model. Following is our main result.

**Theorem 2.** *Consider a generating model  $(P_X, \mathbf{p})$  such that  $\delta(P_X, \mathbf{p}) > \alpha$  for some  $\alpha > 0$ . Further, suppose that  $\mathcal{N}_j \leq n$  for all  $j \in \mathcal{M}$ . Then for any  $\omega \in (0, \alpha K(|\mathcal{Y}| - 1))$ , there exists a positive integer*

$N_0$  that depends only on  $\omega$ ,  $\alpha$ ,  $|\mathcal{Y}|$ ,  $n$ , and  $K$  such that if the number of evaluation tasks  $N > N_0$ , then

- Our mechanism is strictly Bayes-Nash incentive compatible with respect to  $(P_X, \mathbf{p})$ , and,
- At the truthful Bayes-Nash equilibrium, the expected payoff to an agent under the truthful strategy is at least  $\omega$  higher than the expected payoff under any reporting strategy where the agent's reported response is independent of her true response.

This result implies that if it is common knowledge amongst agents that the generating model  $(P_X, \mathbf{p})$  is such that  $\delta(P_X, \mathbf{p}) > \alpha$  for some  $\alpha > 0$ , then irrespective of their beliefs about the specifics of the generating model, truthful reporting is a strict Bayes-Nash equilibrium in the game induced by the mechanism for a large enough  $N$ . Moreover, the truthful equilibrium can be sustained even when the agents incur a bounded cost for their evaluation effort by scaling the rewards appropriately.

### 3.3 Equilibrium selection

Next, we address the issue of multiplicity of equilibria. First, observe that if truthful behavior is an equilibrium, then so is any *symmetric fully informative strategy profile* where all agents apply a common permutation map to the responses they receive. And all such equilibria are payoff-equivalent. But the significantly higher degree of coordination needed for the agents to play a fully informative equilibrium other than truthful behavior, coupled with the fact that there is no real benefit in doing so, makes it unlikely that such equilibria will emerge in practice. Thus full informativeness shall be our benchmark as we focus on other equilibria that may emerge.

The equilibria that give high expected payoffs are arguably the most attractive for the agents and thus can be assumed to have an increased likelihood of being chosen. In what follows, we show that for a large  $N$ , the truthful equilibrium is approximately payoff-optimal across all symmetric equilibria, with an approximation error that vanishes in  $N$ . In the limit, any symmetric fully informative strategy profile gives a strictly higher expected payoff to any agent than *any* other symmetric strategy profile. We also show that under certain assumption on the strategy spaces, any symmetric equilibrium that results in the highest expected payoff to an agent across all symmetric equilibria cannot be too “uninformative” when  $N$  is large, where uninformativeness is a precise notion that we will define in due course.

#### 3.3.1 Truthfulness vs. symmetric equilibria in the large $N$ regime

Before we discuss the result for a large but finite  $N$ , let us first discuss the result in the limiting case as  $N \rightarrow \infty$ , which is easier to obtain, and sheds light on the core idea. Consider a symmetric strategy profile in which every agent adopts a reporting strategy  $\mathbf{q}$ , where  $\mathbf{q}(y) = (q_{y'}(y); y' \in \mathcal{Y})$  is the distribution over the reported response conditional on the true response. Let us denote the reported responses under this strategy by the random variables  $\{Z_j^i; i = 1, \dots, N, j \in \mathcal{M}_i\}$ . Under the truthful strategy profile (or equivalently, any symmetric fully informative strategy profile), in the limit as  $N \rightarrow \infty$ , the expected reward of each agent performing task  $i$  converges to (see Equation 2),

$$\sum_{y \in \mathcal{Y}} \mathbb{P}(Y_j^i = y) K \frac{\sqrt{\mathbb{P}(Y_{j'}^i = Y_j^i = y)}}{\mathbb{P}(Y_j^i = y)} = K \sum_{y \in \mathcal{Y}} \sqrt{\mathbb{P}(Y_{j'}^i = Y_j^i = y)} = K \sum_{y \in \mathcal{Y}} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_y(x)^2}. \quad (8)$$

Whereas, under any other symmetric strategy profile, the expected reward of each agent converges to,

$$K \sum_{y \in \mathcal{Y}} \sqrt{P(Z_{j'}^i = Z_j^i = y)} = K \sum_{y \in \mathcal{Y}} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) \left( \sum_{y' \in \mathcal{Y}} p_{y'}(x) q_y(y') \right)^2}. \quad (9)$$

It turns out that the quantity in (9) is in general lower than the quantity in (8). How much lower depends on the ‘uninformativeness’ of the strategy  $\mathbf{q}$ : more uninformative the strategy  $\mathbf{q}$ , the higher is the difference. Speaking informally, a reporting strategy is more uninformative if it frequently maps multiple true responses to a single reported response, the extreme case being when a report is chosen independently of the true response. The following definition formalizes this notion.

**Definition 3.3 (An uninformativeness measure).** *The uninformativeness of a reporting strategy  $\mathbf{q}$  is defined as*

$$\Omega(\mathbf{q}) = \frac{1}{|\mathcal{Y}|(|\mathcal{Y}| - 1)} \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}; y' \neq y''} \sqrt{q_y(y') q_y(y'')}. \quad (10)$$

We say that a strategy  $\mathbf{q}$  is  $\omega$ -uninformative if  $\Omega(\mathbf{q}) \geq \omega$ .

Clearly,  $\Omega(\mathbf{q}) = 0$  if and only if  $(\mathbf{q}(y); y \in \mathcal{Y})$  have disjoint supports across all  $y \in \mathcal{Y}$ , i.e., if and only if  $\mathbf{q}$  is fully informative. On the other hand  $\Omega(\mathbf{q})$  attains its highest value of 1, if and only if  $\mathbf{q}(y) = \mathbf{q}(y')$  for any  $y \neq y'$ , i.e., if the report is chosen independent of the true answer. To see this, observe that,

$$\begin{aligned} & \frac{1}{|\mathcal{Y}|(|\mathcal{Y}| - 1)} \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} \sqrt{q_y(y') q_y(y'') \mathbf{1}_{y' \neq y''}} \\ & \stackrel{(a)}{\leq} \frac{1}{|\mathcal{Y}|(|\mathcal{Y}| - 1)} \sum_{y \in \mathcal{Y}} \sqrt{\left( \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} q_y(y') \mathbf{1}_{y' \neq y''} \right) \left( \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} q_y(y'') \mathbf{1}_{y' \neq y''} \right)} \\ & = \frac{1}{|\mathcal{Y}|(|\mathcal{Y}| - 1)} \sum_{y \in \mathcal{Y}} \sqrt{(|\mathcal{Y}| - 1)^2 \left( \sum_{y' \in \mathcal{Y}} q_y(y') \right)^2} \\ & = \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}} q_y(y') \\ & = 1. \end{aligned}$$

Here, (a) follows from the Cauchy-Schwarz inequality. We will see that higher the uninformativeness of a strategy, the higher is the gap between the payoffs in equations (8) and (9). It is useful to describe this phenomenon more generally, since it could be of independent interest beyond this work. We define the following notion of an agreement measure.

**Definition 3.4 (An agreement measure).** *Consider a generating model  $(P_X, \mathbf{p})$  defined over  $\mathcal{X}$  and  $\mathcal{Y}$ , and consider two random responses  $Y_1$  and  $Y_2$  drawn from this model. Then the agreement measure between  $Y_1$  and  $Y_2$  is defined as*

$$\Gamma(Y_1, Y_2) = \sum_{y \in \mathcal{Y}} \sqrt{P(Y_1 = Y_2 = y)} = \sum_{y \in \mathcal{Y}} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_y(x)^2}$$

Under a symmetric strategy profile in which every agent adopts a reporting strategy  $\mathbf{q}$ , the

expected payoff to each agent in the limit as  $N \rightarrow \infty$  is  $K$  times the agreement measure between the reported responses (see Equation 9). The agreement measure has the following properties:

1.  $\Gamma(Y_1, Y_2) \geq 1$ . To see this, note that Jensen's inequality implies that

$$\sum_{y \in \mathcal{Y}} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_y(x)^2} \geq \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} P_X(x) p_y(x) = 1.$$

In fact  $\Gamma(Y_1, Y_2) = 1$  only when  $Y_1$  and  $Y_2$  are independent.

2.  $\Gamma(Y_1, Y_2) \leq \sqrt{|\mathcal{Y}|}$ . To see this, note that Jensen's inequality implies that

$$\begin{aligned} \sum_{y \in \mathcal{Y}} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_y(x)^2} &\leq |\mathcal{Y}| \sqrt{\frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} P_X(x) p_y(x)^2} \\ &\leq |\mathcal{Y}| \sqrt{\frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} P_X(x) p_y(x)} = \sqrt{|\mathcal{Y}|}. \end{aligned}$$

In fact  $\Gamma(Y_1, Y_2) = \sqrt{|\mathcal{Y}|}$  only when  $Y_1$  and  $Y_2$  are identical and they are distributed uniformly, i.e.,  $Y_1 = Y_2$  and  $P(Y_1 = y) = 1/|\mathcal{Y}|$  for all  $y \in \mathcal{Y}$ .

The following information monotonicity property is key to our results.

**Proposition 3 (A monotonicity property).** *Consider a generating model  $(P_X, \mathbf{p})$  defined over  $\mathcal{X}$  and  $\mathcal{Y}$ , and consider two random responses  $Y_1$  and  $Y_2$  drawn from this model. Also, consider two random responses  $Z_1$  and  $Z_2$  obtained by applying a reporting strategy  $\mathbf{q}$  independently to  $Y_1$  and  $Y_2$  respectively. Then,*

$$\Gamma(Z_1, Z_2) \leq \Gamma(Y_1, Y_2) - \frac{\delta(P_X, \mathbf{p}) \Omega(\mathbf{q})^2 (|\mathcal{Y}| - 1)}{2\sqrt{|\mathcal{Y}|}}. \quad (11)$$

Clearly,  $\Gamma(Y_1, Y_2) = \Gamma(Z_1, Z_2)$  if  $\mathbf{q}$  is a permutation map., i.e., it is fully informative. The proposition implies that if  $\delta(P_X, \mathbf{p}) > 0$ , then  $\Gamma(Y_1, Y_2) = \Gamma(Z_1, Z_2)$  *only if*  $\Omega(\mathbf{q}) = 0$ , i.e., only if  $\mathbf{q}$  is fully informative. Thus we immediately conclude that if  $\delta(P_X, \mathbf{p}) > 0$ , then in the limit as  $N \rightarrow \infty$ , any fully informative strategy profile gives a strictly higher payoff than any other symmetric strategy profile that is not fully informative.

Now we turn to the finite  $N$  setting. In this case, the expected payoffs under the fully informative strategy and under any other symmetric strategy will not have converged to  $\Gamma(Y_1, Y_2)$  and  $\Gamma(Z_1, Z_2)$  respectively. But for any fixed  $N$ , one can obtain concentration bounds on how far the expected payoffs will be from these target values. This, in turn, gives us vanishing bounds on how much lower the payoff under the truthful equilibrium could be as compared to any other symmetric equilibrium.

**Theorem 4.** *Consider a generating model  $(P_X, \mathbf{p})$  such that  $\delta(P_X, \mathbf{p}) > 0$ . Suppose that this generating model is common knowledge amongst agents. Further, suppose that  $\mathcal{N}_j \leq n$  for all  $j \in \mathcal{M}$ . Then for any  $\omega > 0$ , there is a positive integer  $N_0$  that depends on  $\delta(P_X, \mathbf{p})$ ,  $\omega$ ,  $n$ ,  $K$ , and  $|\mathcal{Y}|$ , such that for any  $N > N_0$ ,*

1. *Any symmetric fully informative strategy profile is a strict Bayes-Nash equilibrium, and,*
2. *It gives an expected payoff that is at most  $\omega$  less than any other symmetric Bayes-Nash equilibrium strategy profile.*

Note that the  $N_0$  in the statement above suffices to ensure that the truthful equilibrium gives a payoff no less than  $\omega$  compared to *any* symmetric equilibrium strategy profile. This is significant, since for a large but fixed  $N$ , it is not possible to obtain uniformly vanishing concentration bounds on  $E(e_j(y))$  (which involves an inverse) across all symmetric reporting strategies. This is because there could be a symmetric strategy profile for which the probability of agreement for an answer  $y \in \mathcal{Y}$  gets arbitrarily close to 0. To overcome this issue, we utilize the fact that under a mixed equilibrium, since the problem of computing the best response is a linear optimization problem, a fixed agent is indifferent between multiple deterministic reporting strategies. This allows us to choose a best-response strategy for a single agent in a way that ensures that the probability of agreement on any answer  $y$  is bounded away from zero, while ensuring that the expected payoff is same as that under the given symmetric equilibrium. Note that the case where the probability of agreement on any answer  $y$  is 0 is not problematic since in this case  $e_j(y)$  is defined to be 0.

Can we say anything about the informativeness of the symmetric equilibrium profile that gives the highest expected payoff across all symmetric equilibria? Intuitively, bounds on  $E(e_j(y))$  for a large  $N$ , coupled with the “inequality gap” characterized in Proposition 3 should result in an upper bound on the uninformative nature of any symmetric reporting strategy that gives a higher expected payoff than a fully informative strategy. It turns out that in doing so, the same difficulty that we described earlier arises in obtaining the requisite concentration bounds when the symmetric strategies could lead to probabilities of agreement arbitrarily close to 0. In this case, we cannot use the trick we used earlier and instead, we show the following result.

**Theorem 5.** *Consider a generating model  $(P_X, \mathbf{p})$  such that  $\delta(P_X, \mathbf{p}) > 0$ . Suppose that this generating model is common knowledge amongst agents. Further, suppose that  $N_j \leq n$  for all  $j \in \mathcal{M}$ . Also, suppose that for each evaluation task, each agent restricts herself to using a reporting strategy in which the probability of reporting any answer  $y \in \mathcal{Y}$  is either 0 or at least  $\eta$  for some  $\eta \in (0, \delta(P_X, \mathbf{p})^2]$ . Then for any  $\omega > 0$ , there is a positive integer  $N_0$  that depends on  $\omega$ ,  $\delta(P_X, \mathbf{p})$ ,  $\eta$ ,  $n$ ,  $K$ , and  $|\mathcal{Y}|$ , such that for any  $N > N_0$ , any symmetric strategy profile that gives a higher expected payoff to an agent than the truthful strategy profile is at most  $\omega$ -uninformative.*

This restriction on the reporting strategies is not implausible since agents are expected to use simple strategies in practice. The lower bound  $\eta$ , for instance, could arise from the desire to produce a marginal distribution of answers consistent with the prior; recall that the prior probability of any answer is bounded away from 0 as long as  $\delta(P_X, \mathbf{p}) > 0$ . The condition  $\eta \leq \delta(P_X, \mathbf{p})^2$  ensures that the fully informative reporting strategy is contained in the restricted strategy space.

## 4 Discussion

### 4.1 Small vs. large number of evaluations per agent

The Correlated Agreement mechanism of [SAFP16] and the mechanism of [KS16] are both minimal, focal, and detail-free for a general distribution of population responses. Their design operates on every pair of agents and requires an estimate of certain functionals of the joint distributions of their evaluations, which in turn requires these agents to perform a large number of tasks (all of the above properties are obtained asymptotically). [AMPS17] combines correlated agreement with agent clustering to reduce the sample complexity of obtaining these estimates in a heterogeneous population setting. These approaches are suitable for crowdsourcing settings where the principal controls task allocations and it is natural for agents to perform a large number of tasks. But they

are impractical in reputation systems of online platforms where agents are expected to perform a small number of evaluations of their choice.

In the view of the impossibility of designing mechanisms in the heterogeneous settings that do not obtain any information about agent heterogeneity [RF15], it is unlikely that it is possible to design minimal, detail-free mechanisms that only obtain a few evaluations per agent. It seems necessary to make certain ‘regularity’ assumptions on the distribution of agent responses so that the correlation structure between pairs of agents cannot vary arbitrarily. Recall that the mechanism in [DG13] requires the responses to be *categorical* while the Peer Truth Serum [JF11, RFJ16] requires the responses to be *self-predicting*. Both obtain truthfulness without requiring the agents to perform a large number of tasks. The former also obtains that the truthful equilibrium is focal.

We, on the other hand, impose the homogeneity assumption, which does not imply either of the two conditions in general. While we obtain truthfulness, we obtain that the truthful equilibrium is focal only across symmetric equilibria. There could be other non-symmetric equilibria where some agents may get a higher payoff than the truthful or a symmetric fully informative equilibrium. For instance, suppose  $y^*$  is the answer with the lowest probability of agreement, i.e.,  $y^* = \arg \min_{y \in \mathcal{Y}} P(Y_1 = Y_2 = y)$ . Then a strategy profile where all agents report  $y^*$  for some task  $i$ , and are truthful for all other tasks is an equilibrium for a large enough  $N$ . This equilibrium gives a higher reward to the agents that evaluate task  $i$  than the fully truthful equilibrium. We do not expect such asymmetric equilibria to arise given the anonymity of agents and tasks in a large platform. It nevertheless remains an open question whether we can design a truthful and focal mechanism with a small number of evaluations per agent under population homogeneity.

We finally mention that in settings like crowdsourcing where the principal controls task allocations and agents are naturally expected to perform a large number of evaluations, we can utilize our agreement measure to design a (asymptotically) truthful *and focal* mechanism under the [KS16] framework. This is facilitated by the following inequality satisfied by the agreement measure, which generalizes Proposition 3 without the characterizing the inequality gap.

**Proposition 6 (A general monotonicity property).** *Consider a generating model  $(P_X, \mathbf{p})$  defined over  $\mathcal{X}$  and  $\mathcal{Y}$ , and consider two random responses  $Y_1$  and  $Y_2$  drawn from this model. Also, consider two random responses  $Z_1$  and  $Z_2$  obtained by applying a reporting strategies  $\mathbf{q}$  and  $\mathbf{q}'$  independently to  $Y_1$  and  $Y_2$  respectively. Then,*

$$\sum_{y \in \mathcal{Y}} \sqrt{P(Z_1 = Z_2 = y)} \leq \Gamma(Y_1, Y_2). \quad (12)$$

In the task allocation scheme, agents are paired together, and each pair performs the same set of tasks. The payment mechanism is defined for each pair of agents as follows.

**Definition 4.1 (An alternative mechanism).** *Let the agents be denoted as 1 and 2, and let  $\mathcal{N}$  be the set of tasks that they perform, labeled as  $i = 1, \dots, N$ . The responses of the agents for the different evaluation tasks are solicited. Let these be denoted by  $\{r_1^i : i \in \mathcal{N}\}$  and  $\{r_2^i : i \in \mathcal{N}\}$ . Then agents’ payment is computed as follows:*

- For each task  $i$  and for each possible evaluation  $y \in \mathcal{Y}$ , compute the quantity

$$\bar{f}^i(y) = \frac{1}{N-1} \sum_{i' \in \mathcal{N} \setminus i} \mathbf{1}_{\{r_1^{i'}=y\}} \mathbf{1}_{\{r_2^{i'}=y\}}.$$

- For each answer  $y$ , fix a payment  $e^i(y)$  defined as

$$e^i(y) = \begin{cases} \frac{K}{\sqrt{\bar{f}^i(y)}} & \text{if } \bar{f}^i(y) \neq 0, \\ 0 & \text{if } \bar{f}^i(y) = 0, \end{cases}$$

where  $K > 0$  is any positive constant.

- For computing the agents' payments for task  $i$ : if their responses match, i.e., if  $r_1^i = r_2^i = y$ , then both gets a reward of  $e^i(y)$ . If the responses do not match, then they get 0 payment for that task.

We will assume that each agent restricts herself to choosing a single common reporting strategy across all evaluation tasks. This assumption is common in literature in this setting, informally stated as “workers treat the tasks equally” [DG13]. For a fixed number of evaluation tasks  $N$ , let  $G^N(\mathbf{q}, \mathbf{q}')$  denote the expected payoff to either of the agents for evaluating any single task (the expected payoff will be the same to both agents) when one agent plays strategy  $\mathbf{q}$  and other plays  $\mathbf{q}'$ . Also, suppose that  $\mathbf{t}$  is the truthful reporting strategy. Then

$$G^N(\mathbf{q}, \mathbf{q}') = \sum_{y \in \mathcal{Y}} \mathbb{E}(e^i(y)) \sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x) p_{y_2}(x) q_y(y_1) q'_y(y_2). \quad (13)$$

But then  $\lim_{N \rightarrow \infty} \mathbb{E}(e^i(y)) = \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x) p_{y_2}(x) q_y(y_1) q'_y(y_2)}$ , and thus,

$$\lim_{N \rightarrow \infty} G^N(\mathbf{q}, \mathbf{q}') = \sum_{y \in \mathcal{Y}} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x) p_{y_2}(x) q_y(y_1) q'_y(y_2)}. \quad (14)$$

Thus, we can use Proposition 6 to conclude that,

$$\lim_{N \rightarrow \infty} G^N(\mathbf{q}, \mathbf{q}') \leq \lim_{N \rightarrow \infty} G^N(\mathbf{t}, \mathbf{t}). \quad (15)$$

for any  $\mathbf{q}$  and  $\mathbf{q}'$ . Taking  $\mathbf{q}$  (or  $\mathbf{q}'$ ) to be equal to  $\mathbf{t}$ , this shows that the proposed mechanism is Bayes-Nash incentive compatible “in the limit”. Further, it is also focal “in the limit”, i.e., asymptotically, the payoff under the truthful strategy is no less than the payoff under any other strategy profile. Obtaining finite  $N$  results is trickier in this case and beyond the scope of our work.

## 4.2 Peer Prediction Method with estimation

The original Peer Prediction method of [MRZ05] is defined for the homogeneous population setting, but it requires the principal to know the joint distribution of answers of any two agents for implementation. Using an approach similar to ours, it is possible to relax this knowledge assumption when there are a large number of evaluation tasks, by obtaining the required estimates from the reports for a large number of tasks. While truthfulness and strict properness (assuming stochastic relevance) could be achieved asymptotically in such a mechanism, in general, it is not true that truthful behavior gives a higher reward than any other symmetric strategy profile. We demonstrate this for the spherical scoring rule [GR07]. Suppose that the conditional distribution of the answers of a peer if the answer of an agent is  $y$  is denoted as  $\mathbf{m}(y) = (m_{y'}(y); y' \in \mathcal{Y})$ . Under the Peer



Prediction method utilizing the spherical scoring rule, if the agent reports  $y$  and the peer reports  $y'$ , then the reward of the agent is:

$$R(y, y') = \frac{m_{y'}(y)}{\sqrt{\sum_{y'' \in \mathcal{Y}} m_{y''}(y)^2}}. \quad (16)$$

Let  $Y_1$  and  $Y_2$  be the reports of  $j$  and  $j'$  under the truthful strategy, and let  $Z_1$  and  $Z_2$  be the reports under some other symmetric strategy profile  $\mathbf{q}$ . Then from (16), assuming that the principal has obtained accurate estimates of the required distributions, it is easy to see that the expected reward under the truthful strategy is,

$$\sum_{y \in \mathcal{Y}} P(Y_1 = y) \sqrt{\sum_{y' \in \mathcal{Y}} P(Y_2 = y' | Y_1 = y)^2}, \quad (17)$$

and that under the other symmetric strategy profile is,

$$\sum_{y \in \mathcal{Y}} P(Z_1 = y) \sqrt{\sum_{y' \in \mathcal{Y}} P(Z_2 = y' | Z_1 = y)^2}. \quad (18)$$

Now, it is not true in general that (18) is strictly smaller than (17). For instance, suppose that  $Y_1$  and  $Y_2$  are i.i.d. with a uniform distribution over  $\mathcal{Y}$  (i.e.,  $\mathcal{X}$  is a singleton). Then the reward under the truthful strategy profile is  $1/\sqrt{|\mathcal{Y}|}$ . But the non-truthful strategy of always reporting a fixed answer gives an expected reward of 1, which is larger. Under our framework, both these strategies give an agreement measure of 1.

## 5 Conclusion

We presented a new payment mechanism with strong incentive properties for obtaining truthful reports from agents in settings where a large number of evaluations are to be made and each agent performs a subset of these evaluations. Our main assumption is the homogeneity of the population performing a single task. The mechanism has the simple structure of output agreement mechanisms, which are often adopted in crowdsourcing applications. Although the naive output agreement scheme does not necessarily incentivize truthfulness, the basic insight from our work is that inversely scaling the agreement rewards by the square roots of the corresponding probabilities of agreement appropriately aligns incentives. In the process, we made new contributions to the calculus of peer-prediction through the information-theoretic notion of the agreement measure and the notion of uninformativeness of reporting strategies.

Effective feedback and reputation systems are fundamental to the efficient functioning of online platforms. The impact of user feedback and peer-reviews on customer decisions is evident in the success of independent reputation systems like Yelp and TripAdvisor, which are used by millions of people across the world. But as has been recently shown, these systems are currently fraught with several operational as well as behavioral and strategic concerns [HG15, NT15]. We believe that appropriate incentive mechanisms that are simple and intuitive can go a long way in addressing some of these concerns, and hence we think that our mechanism has strong practical significance. We emphasize here that rather than thinking of our mechanism as a fully specified solution in any setting, it is more useful to think of it as a framework that provides conceptual guidelines for

platform designers as they undertake their design decisions.

As we argued in the introduction, population homogeneity is a reasonable assumption for objective queries about personal experiences (e.g., Questions (1) to (3)). Considering that subjective reviews and ratings are fraught with biases, we believe that transitioning to such objective questions, or at least complementing subjective reviews with such questions may be a sound design decision. Moreover, from the perspective of the platform interested in curating the composition of the pool of agents that transact on the platform, answers to such questions are no less important than obtaining subjective ratings. For instance, it is important for a labor platform to know if a particular handyman is habitually late for appointments, or a cleaner does not stick to established standards. As Nosko and Tadelis have elegantly argued in [NT15], there are reputational externalities caused by such experiences – a bad encounter can cause a customer to attribute that experience to the platform rather than a single bad agent and thus may leave the platform entirely. Many platforms have already started eliciting fine-grained feedback along multiple dimensions, thereby reducing its subjectivity. For instance, the short term lodging and apartment rental marketplace, Airbnb, elicits renter feedback on objective dimensions like accuracy of description, cleanliness, ease of check-in, etc.

Our work presents many avenues for future exploration. For instance, in our model, we assume that the task allocations are exogenously specified. But for a platform that is interested in learning the underlying distributions of responses for each task, some of these distributions may be more difficult to learn than others, and thus may need more evaluations. Moreover, the agents may be willing to strategically respond to differences in potential rewards across tasks by choosing which tasks to evaluate. In such situations, it is important to understand the fundamental tradeoffs faced by dynamic mechanisms that balance incentives with different statistical accuracy objectives. We are optimistic that our framework and insights from this paper can be used as building blocks in this pursuit.

## References

- [AMPS17] Arpit Agarwal, Debmalya Mandal, David C. Parkes, and Nisarg Shah. Peer prediction with heterogeneous users. In *Proceedings of the 2017 ACM Conference on Economics and Computation*, EC '17, pages 81–98, New York, NY, USA, 2017. ACM.
- [Bri50] Glenn W Brier. Verification of forecasts expressed in terms of probability. *Monthly weather review*, 78(1):1–3, 1950.
- [DG13] Anirban Dasgupta and Arpita Ghosh. Crowdsourced judgement elicitation with endogenous proficiency. In *Proceedings of the 22nd international conference on World Wide Web*, pages 319–330. International World Wide Web Conferences Steering Committee, 2013.
- [GR07] Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007.
- [GWL16] Alice Gao, James R Wright, and Kevin Leyton-Brown. Incentivizing evaluation via limited access to ground truth: Peer-prediction makes things worse. *arXiv preprint arXiv:1606.07042*, 2016.

- [HG15] John Horton and Joseph Golden. Reputation inflation: Evidence from an online labor market. 2015.
- [JF05] Radu Jurca and Boi Faltings. Enforcing truthful strategies in incentive compatible reputation mechanisms. In *WINE*, volume 3828, pages 268–277. Springer, 2005.
- [JF11] Radu Jurca and Boi Faltings. Incentives for answering hypothetical questions. In *Workshop on Social Computing and User Generated Content, EC-11*, number EPFL-CONF-197783, 2011.
- [KS16] Yuqing Kong and Grant Schoenebeck. A framework for designing information elicitation mechanisms that reward truth-telling. *arXiv preprint arXiv:1605.01021*, 2016.
- [LC17] Yang Liu and Yiling Chen. Machine-learning aided peer prediction. In *Proceedings of the 2017 ACM Conference on Economics and Computation*, pages 63–80. ACM, 2017.
- [LS09] Nicolas Lambert and Yoav Shoham. Eliciting truthful answers to multiple-choice questions. In *ACM conference on Electronic commerce*, pages 109–118, 2009.
- [Luc17] Michael Luca. Designing online marketplaces: Trust and reputation mechanisms. *Innovation Policy and the Economy*, 17(1):77–93, 2017.
- [MRZ05] Nolan Miller, Paul Resnick, and Richard Zeckhauser. Eliciting informative feedback: The peer-prediction method. *Management Science*, 51(9):1359–1373, 2005.
- [NT15] Chris Nosko and Steven Tadelis. The limits of reputation in platform markets: An empirical analysis and field experiment. Technical report, National Bureau of Economic Research, 2015.
- [Pre04] Dražen Prelec. A Bayesian truth serum for subjective data. *Science*, 306(5695):462–466, 2004.
- [RF13] Goran Radanovic and Boi Faltings. A robust bayesian truth serum for non-binary signals. In *Proceedings of the 27th AAAI Conference on Artificial Intelligence, AAAI 2013*, number EPFL-CONF-197486, pages 833–839, 2013.
- [RF15] Goran Radanovic and Boi Faltings. Incentives for subjective evaluations with private beliefs. In *Proceedings of the 29th AAAI Conference on Artificial Intelligence (AAAI’15)*, 2015.
- [RFJ16] Goran Radanovic, Boi Faltings, and Radu Jurca. Incentives for effort in crowdsourcing using the peer truth serum. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 7(4):48, 2016.
- [SAFP16] Victor Shnayder, Arpit Agarwal, Rafael Frongillo, and David C Parkes. Informed truthfulness in multi-task peer prediction. In *Proceedings of the 2016 ACM Conference on Economics and Computation*, pages 179–196. ACM, 2016.
- [Sav71] Leonard J Savage. Elicitation of personal probabilities and expectations. *Journal of the American Statistical Association*, 66(336):783–801, 1971.
- [VAD04] Luis Von Ahn and Laura Dabbish. Labeling images with a computer game. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 319–326. ACM, 2004.

- [VAD08] Luis Von Ahn and Laura Dabbish. Designing games with a purpose. *Communications of the ACM*, 51(8):58–67, 2008.
- [WP12a] Jens Witkowski and David C Parkes. Peer prediction without a common prior. In *Proceedings of the 13th ACM Conference on Electronic Commerce*, pages 964–981. ACM, 2012.
- [WP12b] Jens Witkowski and David C Parkes. A robust Bayesian truth serum for small populations. In *AAAI*, 2012.
- [WP13] Jens Witkowski and David C Parkes. Learning the prior in minimal peer prediction. In *3rd Workshop on Social Computing and User Generated Content at the ACM Conference on Electronic Commerce*, 2013.

## A Proofs

*Proof of Proposition 1.* To see that 2 implies 1, note that  $\theta(\mathbf{v}(y), \mathbf{v}(y')) \geq \tau$  implies that:

$$\frac{\mathbf{v}(y) \cdot \mathbf{v}(y')}{\|\mathbf{v}(y)\| \|\mathbf{v}(y')\|} \leq \cos \tau,$$

Multiplying throughout by  $\|\mathbf{v}(y)\| \|\mathbf{v}(y')\|$ , we have:

$$\|\mathbf{v}(y)\| \|\mathbf{v}(y')\| - \mathbf{v}(y) \cdot \mathbf{v}(y') \geq (1 - \cos \tau) \|\mathbf{v}(y)\| \|\mathbf{v}(y')\| \geq (1 - \cos \tau) \tau^2 > 0.$$

Here in the last inequality, we use that fact that

$$\|v(s)\| = \sqrt{\sum_{h \in \mathcal{H}} P_X(h) p(s|h)^2} \geq \sum_{h \in \mathcal{H}} P_X(h) p(s|h) > \tau,$$

which follows from the Jensen's inequality. To show that 1 implies 2 is less straightforward and this is where we need to use the fact that  $\|\mathbf{v}(y)\| \leq 1$  for all  $y \in \mathcal{Y}$ . First of all

$$|\mathbf{v}(y) \cdot \mathbf{v}(y')| \leq \|\mathbf{v}(y)\| \|\mathbf{v}(y')\| - \alpha,$$

implies that either  $\|\mathbf{v}(y)\|$  or  $\|\mathbf{v}(y')\|$  is non-zero. Say  $\|\mathbf{v}(y')\| > 0$ . Then dividing on both sides, we get:

$$\begin{aligned} \frac{|\mathbf{v}(y) \cdot \mathbf{v}(y')|}{\|\mathbf{v}(y')\|} &\leq \|\mathbf{v}(y)\| - \frac{\alpha}{\|\mathbf{v}(y')\|} \\ &\leq \|\mathbf{v}(y)\| - \alpha \end{aligned}$$

where the last inequality holds since  $\|\mathbf{v}(y)\| \leq 1$ . In other words:

$$\|\mathbf{v}(y)\| \cos \theta(\mathbf{v}(y), \mathbf{v}(y')) \leq \|\mathbf{v}(y)\| - \alpha,$$

or

$$\|\mathbf{v}(y)\| (1 - \cos \theta(\mathbf{v}(y), \mathbf{v}(y'))) \geq \alpha.$$

Since  $\|\mathbf{v}(y)\| \leq 1$  and  $(1 - \cos\theta(\mathbf{v}(y), \mathbf{v}(y'))) \in [0, 1]$ , this implies both  $\|\mathbf{v}(y)\| \geq \alpha$  and  $(1 - \cos\theta(\mathbf{v}(y), \mathbf{v}(y'))) \geq \alpha$ , i.e.,  $\theta(\mathbf{v}(y), \mathbf{v}(y'))) \geq \arccos(1 - \alpha)$ . Finally we have  $\sum_{x \in \mathcal{X}} P_X(x)p_y(x) \geq \|\mathbf{v}(y)\|^2 \geq \alpha^2$ . Note that  $\alpha \leq 1$  so that  $\arccos(1 - \alpha) \leq \pi/2$ .

□

*Proof of Theorem 2.* First, note that the payments  $e_j(y)$  for the different  $y \in \mathcal{Y}$  are independent of the reports of agent  $j$  for any reporting strategy. This is because  $\{e_j(y) : y \in \mathcal{Y}\}$  are computed only based on evaluation tasks that  $j$  does not perform. Next, suppose that everyone but agent  $j$  is truthful. Recalling the definition of  $\mathbf{v}(y) \triangleq (\sqrt{P_X(x)p_y(x)}; x \in \mathcal{X})$ , we have,

$$\mathbb{E}(\bar{f}_j(y)) = \mathbb{E}\left[\frac{1}{N - |\mathcal{N}_j|} \sum_{i \in \mathcal{N} \setminus \mathcal{N}_j} \mathbf{1}_{\{r_{j_1}^i = y\}} \mathbf{1}_{\{r_{j_2}^i = y\}}\right] = \sum_{x \in \mathcal{X}} P_X(x)p_y(x)^2 = \|\mathbf{v}(y)\|^2 \triangleq g(y).$$

In the proof of Proposition 1, we have seen that  $\delta(P_X, \mathbf{p}) > \alpha$  implies that  $\|\mathbf{v}(y)\| > \alpha$ , and thus we have  $g(y) > \alpha^2 > 0$  for all  $y \in \mathcal{Y}$ . Next, recall that

$$e_j(y) = \mathbf{1}_{\{\bar{f}_j(y) \neq 0\}} \frac{K}{\sqrt{\bar{f}_j(y)}}.$$

Let  $N' = N - |\mathcal{N}_j|$ . Then we have for any  $\epsilon \in (0, 1)$ :

$$\begin{aligned} \mathbb{E}(e_j(y)) &\geq P(\bar{f}_j(y) \in [g(y)(1 - \epsilon), g(y)(1 + \epsilon)]) \frac{K}{\sqrt{g(y)(1 + \epsilon)}} \\ &\stackrel{(a)}{\geq} (1 - 2 \exp(-\epsilon^2 g(y)^2 N')) \frac{K}{\sqrt{g(y)(1 + \epsilon)}} \\ &\geq (1 - 2 \exp(-\epsilon^2 \alpha^4 N')) \frac{K}{\sqrt{g(y)(1 + \epsilon)}} \\ &\geq \frac{K}{\sqrt{g(y)(1 + \epsilon)}} - 2 \exp(-\epsilon^2 \alpha^4 N') \frac{K}{\alpha \sqrt{1 + \epsilon}} \\ &\stackrel{(b)}{\geq} \frac{K}{\sqrt{g(y)}} (1 - \epsilon) - 2 \exp(-\epsilon^2 \alpha^4 N') \frac{K}{\alpha} \\ &\geq \frac{K}{\sqrt{g(y)}} (1 - \epsilon) - 2 \exp(-\epsilon^2 \alpha^4 (N - n)) \frac{K}{\alpha}. \end{aligned} \tag{19}$$

Here (a) follows from Hoeffding's inequality, and (b) is because  $\frac{1}{\sqrt{1 + \epsilon}} \geq 1 - \epsilon$  for every  $\epsilon \in (0, 1)$ . The other inequalities result from the fact that  $g(s) \geq \alpha^2$  and  $|\mathcal{N}_j| \leq n$ . Taking  $\epsilon = (N - n)^{-1/4}$ , we obtain:

$$\mathbb{E}(e_j(y)) \geq \frac{K}{\sqrt{g(s)}} - o(N).$$

Next, we also have,

$$\mathbb{E}(e_j(y)) \leq P(\bar{f}_j(y) \in [g(y)(1 - \epsilon), g(y)(1 + \epsilon)]) \frac{K}{\sqrt{g(y)(1 - \epsilon)}}$$

$$\begin{aligned}
& + E \left( \mathbf{1}_{\{\bar{f}_j(y) \notin \{0\} \cup [g(y)(1-\epsilon), g(y)(1+\epsilon)]\}} \frac{K}{\sqrt{\bar{f}_j(y)}} \right) \\
& \stackrel{(a)}{\leq} \frac{K}{\sqrt{g(y)(1-\epsilon)}} + P \left( \mathbf{1}_{\bar{f}_j(s) \notin \{0\} \cup [g(y)(1-\epsilon), g(y)(1+\epsilon)]} \right) K \sqrt{N'} \\
& \stackrel{(b)}{\leq} \frac{K}{\sqrt{g(y)(1-\epsilon)}} + 2K \sqrt{N'} \exp(-\epsilon^2 g(y)^2 N') \\
& \leq \frac{K}{\sqrt{g(y)(1-\epsilon)}} + 2K \sqrt{N'} \exp(-\epsilon^2 \alpha^4 N') \\
& \stackrel{(c)}{\leq} \frac{K}{\sqrt{g(y)}} \left( 1 + \frac{\epsilon}{2} + w(\epsilon) \right) + 2K \sqrt{N'} \exp(-\epsilon^2 \alpha^4 N') \\
& \leq \frac{K}{\sqrt{g(y)}} + \frac{\epsilon K}{2\alpha} + \frac{|w(\epsilon)|K}{\alpha} + 2K \sqrt{N'} \exp(-\epsilon^2 \alpha^4 N') \\
& \leq \frac{K}{\sqrt{g(y)}} + \frac{\epsilon K}{2\alpha} + \frac{|w(\epsilon)|K}{\alpha} + 2K \sqrt{N} \exp(-\epsilon^2 \alpha^4 (N-n)). \tag{20}
\end{aligned}$$

Here, (a) results from the fact that on the event  $\{\bar{f}_j(y) \neq 0\}$ ,  $\bar{f}_j(y) \geq 1/N'$ . This is because  $\bar{f}_j(y)$  only takes values in the set  $[0, \frac{1}{N'}, \frac{2}{N'}, \dots, 1]$ . (b) follows from Hoeffding's inequality, and (c) follows from the Taylor approximation of the function  $1/\sqrt{1-\epsilon}$ , where  $w(\epsilon) = o(\epsilon)$ . Now choosing  $\epsilon = (N-n)^{-1/4}$ , we get:

$$E(e_j(y)) \leq \frac{K}{\sqrt{g(y)}} + o(N).$$

Thus, we finally have  $|E(e_j(y)) - \frac{K}{\sqrt{g(y)}}| \leq \sigma(N) = o(N)$ , where  $\sigma(N) \geq 0$  is some function of  $N$  that depends only on  $\alpha$ ,  $n$  and  $K$  and not on  $y$ .

Assuming everyone else is truthful, the expected reward of person  $j$  for evaluating object  $i$  if she chooses a reporting strategy  $\mathbf{q}^{ij}$  is,

$$R(\mathbf{q}^{ij}) \triangleq \sum_{y \in \mathcal{Y}} P(Y_{j'}^i = y, Y_j^i = y) E(r_j(y)) = \sum_{y \in \mathcal{Y}} E(r_j(y)) \sum_{x \in \mathcal{X}} P_X(x) p_y(x) \sum_{y' \in \mathcal{Y}} p_{y'}(x) q_y^{ij}(y').$$

Thus the agent solves  $\max_{\mathbf{q}^{ij}} R(\mathbf{q}^{ij})$ . The objective is linear in  $\mathbf{q}^{ij}$ , and further,  $\mathbf{q}^{ij}(y)$  lies on a unit simplex for each  $y \in \mathcal{Y}$ . Thus the optimal reporting strategy chooses  $\mathbf{q}^{ij}(y)$  to be one of the extreme points of the simplex for each  $y \in \mathcal{Y}$ , i.e., the optimal reporting strategy is deterministic. Now let  $\mathbf{t}$  be the truthful strategy, i.e.,  $t_{y'}(y) = \mathbf{1}_{\{y=y'\}}$ . Then for any deterministic reporting strategy  $\mathbf{q}^{ij}$ , we have,

$$\begin{aligned}
R(\mathbf{q}^{ij}) &= \sum_{y \in \mathcal{Y}} E(e_j(y)) \sum_{y' \in \mathcal{Y}} q_y^{ij}(y') \sum_{x \in \mathcal{X}} P_X(x) p_y(x) p_{y'}(x) \\
&\stackrel{(a)}{\leq} \sum_{y \in \mathcal{Y}} E(e_j(y)) \sum_{y' \in \mathcal{Y}} q_y^{ij}(y') \left( \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_y(x)^2} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y'}(x)^2} - \alpha \mathbf{1}_{\{y \neq y'\}} \right) \\
&\leq \sum_{y \in \mathcal{Y}} \left( \frac{K}{\sqrt{g(y)}} + \sigma(N) \right) \sum_{y' \in \mathcal{Y}} q_y^{ij}(y') \left( \sqrt{g(y)g(y')} - \alpha \mathbf{1}_{\{y \neq y'\}} \right)
\end{aligned}$$

$$\leq K \sum_{y' \in \mathcal{Y}} \sqrt{g(y')} - \alpha K \sum_{y' \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} \mathbf{1}_{\{y \neq y'\}} q_y^{ij}(y') + \sigma(N) \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}} q_y^{ij}(y') \sqrt{g(y)g(y')} \quad (21)$$

$$\stackrel{(b)}{\leq} K \sum_{y' \in \mathcal{Y}} \sqrt{g(y')} - \alpha K \mathbf{1}_{\{\mathbf{q}^{ij} \neq \mathbf{t}\}} + |\mathcal{Y}| \sigma(N). \quad (22)$$

Here, (a) follows from the Cauchy-Schwarz inequality, from the definition of  $\delta(P_x, \mathbf{p})$ , and the fact that  $\delta(P_x, \mathbf{p}) > \alpha$ . (b) follows from the fact that  $\mathbf{q}^{ij}$  is deterministic and so is  $\mathbf{t}$ . While we have,

$$\begin{aligned} R(\mathbf{t}) &= \sum_{y \in \mathcal{Y}} \mathbb{E}(e_j(y)) \sum_{x \in \mathcal{X}} P_X(x) p_y(x)^2 \\ &\geq \sum_{y \in \mathcal{Y}} \left( \frac{K}{\sqrt{g(y)}} - \sigma(N) \right) g(y) \\ &\geq \sum_{y \in \mathcal{Y}} \sqrt{g(y)} - |\mathcal{Y}| \sigma(N). \end{aligned}$$

Thus we have,

$$R(\mathbf{q}^{ij}) \leq R(\mathbf{t}) - \alpha K \mathbf{1}_{\{\mathbf{q}^{ij} \neq \mathbf{t}\}} + 2|\mathcal{Y}| \sigma(N)$$

Since  $\sigma(N)$  depends only on  $\delta$  and  $K$  and  $\sigma(N) = o(1)$ , there is an  $N_1$  that depends only on  $\alpha$ ,  $K$ ,  $n$  and  $|\mathcal{Y}|$  such that for all  $N > N_1$ ,  $2|\mathcal{Y}| \sigma(N) < K\alpha$ , which means that truthful behavior is a strict Bayes-Nash equilibrium. To prove the second statement, suppose that  $\mathbf{q}^{ij}$  is a strategy in which reports are chosen independently of the true answers. Denote  $q_y^{ij} \triangleq q_y^{ij}(y')$  since  $q_y^{ij}(y') = q_y^{ij}(y'')$  for all  $y, y', y'' \in \mathcal{Y}$ . Then in (21),

$$\begin{aligned} \alpha K \sum_{y' \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} \mathbf{1}_{\{y \neq y'\}} q_y^{ij}(y') &= \alpha K \sum_{y' \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} \mathbf{1}_{\{y \neq y'\}} q_y^{ij} \\ &= \alpha K (|\mathcal{Y}| - 1). \end{aligned}$$

And thus,

$$R(\mathbf{q}^{ij}) \leq R(\mathbf{t}) - \alpha K (|\mathcal{Y}| - 1) + 2|\mathcal{Y}| \sigma(N).$$

Thus for any  $\omega \in (0, \alpha K (|\mathcal{Y}| - 1))$ , there is a positive integer  $N_2$  depending on  $\omega$ ,  $\alpha$ ,  $K$ ,  $n$  and  $|\mathcal{Y}|$  such that for any  $N > N_2$ ,  $R(\mathbf{q}^{ij}) \leq R(\mathbf{t}) - \omega$ . Choosing  $N_0 = \max(N_1, N_2)$  proves the result.  $\square$

*Proof of Proposition 3.* We have,

$$\begin{aligned} \Gamma(Z_1, Z_2) &= \sum_{y \in \mathcal{Y}} \sqrt{\sum_{x \in \mathcal{X}, y_1 \in \mathcal{Y}, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) q_y(y_1) q_y(y_2)} \\ &= \sum_{y \in \mathcal{Y}} \sqrt{\sum_{y_1 \in \mathcal{Y}, y_2 \in \mathcal{Y}} q_y(y_1) q_y(y_2) \sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x) p_{y_2}(x)} \\ &\stackrel{(a)}{\leq} \sum_{y \in \mathcal{Y}} \sqrt{\sum_{y_1 \in \mathcal{Y}, y_2 \in \mathcal{Y}} q_y(y_1) q_y(y_2) \left( \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x)^2} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_2}(x)^2} - \delta(P_X, \mathbf{p}) \mathbf{1}_{y_1 \neq y_2} \right)} \\ &= \sum_{y \in \mathcal{Y}} \sqrt{\left( \sum_{y_1 \in \mathcal{Y}} q_y(y_1) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x)^2} \right)^2 - \delta(P_X, \mathbf{p}) \sum_{y_1 \in \mathcal{Y}, y_2 \in \mathcal{Y}} q_y(y_1) q_y(y_2) \mathbf{1}_{y_1 \neq y_2}} \end{aligned}$$

$$\begin{aligned}
&\stackrel{(b)}{\leq} \sum_{y \in \mathcal{Y}, y_1 \in \mathcal{Y}} q_y(y_1) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x)^2} - \frac{\delta(P_X, \mathbf{p})}{2} \sum_{y \in \mathcal{Y}} \frac{\left( \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} q_y(y') q_y(y'') \mathbf{1}_{y' \neq y''} \right)}{\sum_{y_1 \in \mathcal{Y}} q_y(y_1) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x)^2}} \\
&\stackrel{(c)}{\leq} \Gamma(Y_1, Y_2) - \frac{\delta(P_X, \mathbf{p})}{2} \sum_{y \in \mathcal{Y}} \frac{\left( \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} q_y(y') q_y(y'') \mathbf{1}_{y' \neq y''} \right)}{\Gamma(Y_1, Y_2)} \\
&\stackrel{(d)}{\leq} \Gamma(Y_1, Y_2) - \frac{\delta(P_X, \mathbf{p})}{2\sqrt{|\mathcal{Y}|}} \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} q_y(y') q_y(y'') \mathbf{1}_{y' \neq y''} \\
&\stackrel{(e)}{\leq} \Gamma(Y_1, Y_2) - \frac{\delta(P_X, \mathbf{p}) \Omega(\mathbf{q})^2 (|\mathcal{Y}| - 1)}{2\sqrt{|\mathcal{Y}|}}.
\end{aligned}$$

Here, (a) follows from the Cauchy-Schwarz inequality and the definition of  $\delta(P_X, \mathbf{p})$ . (b) follows from the fact that for  $a, b > 0$  and  $a > b$ ,  $\sqrt{a-b} \leq \sqrt{a} - b/(2\sqrt{a})$ . (c) follows from the fact that  $q_y(y_1) \leq 1$  and from the definition of  $\Gamma(Y_1, Y_2)$ . (d) follows from the fact that  $\Gamma(Y_1, Y_2) \leq |\mathcal{Y}|$ . (e) holds since, by Jensen's inequality,

$$\begin{aligned}
\Omega(\mathbf{q})^2 &= \left( \frac{|\mathcal{Y}|}{|\mathcal{Y}|^2 (|\mathcal{Y}| - 1)} \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} \sqrt{q_y(y') q_y(y'') \mathbf{1}_{y' \neq y''}} \right)^2 \\
&\leq \frac{1}{|\mathcal{Y}| - 1} \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} q_y(y') q_y(y'') \mathbf{1}_{y' \neq y''}
\end{aligned}$$

□

*Proof of Theorem 4.* The first statement follows from Theorem 2: there is an  $N_1$  such that for all  $N \geq N_1$ , the truthful strategy profile is a Bayes-Nash equilibrium. We focus on the second claim. With some abuse of notation, we denote  $e_j^t(y)$  to be the agreement scores for an agent  $j$  under the truthful equilibrium, and  $e_j^s(y)$  to be the scores under a fixed symmetric equilibrium strategy profile where each agent follows the reporting strategy  $\mathbf{q}$ .

We have shown in the proof of Theorem 2 that if everyone is truthful, then  $|\mathbb{E}(e_j^t(y)) - \frac{K}{\sqrt{g(y)}}| \leq \sigma(N) = o(1)$ , where  $\sigma(N) \geq 0$  is some function of  $N$  that depends only on  $\delta(P_X, \mathbf{p})$ ,  $n$  and  $K$  and not on  $y$ .

Let us denote  $\sum_{x \in \mathcal{X}} P_X(x) (\sum_{y' \in \mathcal{Y}} p_{y'}(x) q_y(y'))^2 \triangleq s(y)$  and denote  $\sum_{x \in \mathcal{X}} P_X(x) \sum_{y' \in \mathcal{Y}} p_{y'}(x) q_y(y') \triangleq b(y)$ . By Jensen's inequality, we have  $s(y) \geq b(y)^2$ . Then using arguments similar to the ones leading up to (20) in the proof of Theorem 2, we can show that for all  $y \in \mathcal{Y}$  such that  $b(y) \geq \delta(P_X, \mathbf{p})^2 / |\mathcal{Y}|$  (and hence,  $s(y) \geq \delta(P_X, \mathbf{p})^4 / |\mathcal{Y}|^2$ ), and for any  $\epsilon \in (0, 1)$ ,

$$\left| \mathbb{E}(e_j^s(y)) - \frac{K}{\sqrt{s(y)}} \right| \leq \sigma'(N),$$

where  $|\sigma'(N)| = o(1)$  depends on  $\delta(P_X, \mathbf{p})$ ,  $K$ ,  $n$  and  $|\mathcal{Y}|$ . Consider the strategy  $\mathbf{q}$  and consider a  $y \in \mathcal{Y}$ , such that  $b(y) > 0$  but  $b(y) < \delta(P_X, \mathbf{p})^2 / |\mathcal{Y}|$ . Then one can construct another strategy  $\mathbf{q}'$  such that a) a fixed agent  $j$  is indifferent between choosing  $\mathbf{q}$  and  $\mathbf{q}'$  assuming everyone else is playing  $\mathbf{q}$ , and, 2) for all  $y$  such that  $b(y) < \delta(P_X, \mathbf{p}) / |\mathcal{Y}|$ ,  $q'_y(y') = 0$  for all  $y' \in \mathcal{Y}$ . To show this, observe that for each  $y'$ ,  $\mathbf{q}(y')$  cannot have support only on those  $y$  for which  $b(y) < \delta(P_X, \mathbf{p})^2 / |\mathcal{Y}|$ .



This is because if that is the case then  $P(Y_j^i = y') = P(Y_j^i = y') \sum_{y \in \mathcal{Y}; b(y) < \delta(P_X, \mathbf{p})^2 / |\mathcal{Y}|} q_y(y') \leq \sum_{y \in \mathcal{Y}; b(y) < \delta(P_X, \mathbf{p})^2 / |\mathcal{Y}|} b(y) < \delta(P_X, \mathbf{p})^2$ , which contradicts the fact that  $P(Y_j^i = y') \geq \delta(P_X, \mathbf{p})^2$  as we have seen in the proof of Proposition 1. So then define  $\mathbf{q}'(y')$  to have support only on the  $y \in \mathcal{Y}$  for which  $b(y) \geq \delta(P_X, \mathbf{p})^2 / |\mathcal{Y}|$  by transferring the probability masses. If we define  $G(\mathbf{q})$  to be the expected payment to a fixed agent  $j$  for a fixed task  $i$  under the symmetric equilibrium under strategy  $\mathbf{q}$ , and define  $G(\mathbf{q}', \mathbf{q}^{-j})$  to be the expected payment to  $j$  if she plays  $\mathbf{q}'$  while others play  $\mathbf{q}$ , then we have  $G(\mathbf{q}) = G(\mathbf{q}', \mathbf{q}^{-j})$ . Let us define  $\sum_{x \in \mathcal{X}} P_X(x) (\sum_{y' \in \mathcal{Y}} p_{y'}(x) q'_{y'}(y'))^2 \triangleq s'(y)$ . Then we have,

$$\begin{aligned}
G(\mathbf{q}) &= G(\mathbf{q}', \mathbf{q}^{-j}) \\
&\leq \sum_{y \in \mathcal{Y}; b(y) \geq \delta(P_X, \mathbf{p})^2 / |\mathcal{Y}|} \mathbb{E}(e_j^s(y)) \sum_{x \in \mathcal{X}} P_X(x) \left[ \sum_{y_1 \in \mathcal{Y}} p_{y_1}(x) q'_{y_1}(y_1) \right] \left[ \sum_{y_2 \in \mathcal{Y}} p_{y_2}(x) q_{y_2}(y_2) \right] \\
&\stackrel{(a)}{\leq} \sum_{y \in \mathcal{Y}; b(y) \geq \delta(P_X, \mathbf{p})^2 / |\mathcal{Y}|} \mathbb{E}(e_j^s(y)) \sqrt{s(y) s'(y)} \\
&\leq \sum_{y \in \mathcal{Y}; b(y) \geq \delta(P_X, \mathbf{p})^2 / |\mathcal{Y}|} \left( \frac{K}{\sqrt{s(y)}} + \sigma'(N) \right) \sqrt{s(y) s'(y)} \\
&\leq \sum_{y \in \mathcal{Y}; b(y) \geq \delta(P_X, \mathbf{p})^2 / |\mathcal{Y}|} K \sqrt{s'(y)} + |\mathcal{Y}| \sigma'(N) \\
&\stackrel{(b)}{=} K \sum_{y \in \mathcal{Y}} \sqrt{s'(y)} + |\mathcal{Y}| \sigma'(N). \tag{23}
\end{aligned}$$

Here (a) follows from the Cauchy-Schwarz inequality and (b) follows from the fact that  $s'(y) = 0$  for all  $y$  such that  $b(y) < \delta(P_X, \mathbf{p})^2 / |\mathcal{Y}|$  by construction of the strategy  $\mathbf{q}'$ . Let  $G(\mathbf{t})$  be the expected payment to agent  $j$  for task  $i$  under the truthful equilibrium. Let  $j'$  be  $j$ 's peer for task  $i$ . Then we have,

$$\begin{aligned}
G(\mathbf{t}) &= \sum_{y \in \mathcal{Y}} \mathbb{E}(e_j^t(y)) g(y) \\
&\geq \sum_{y \in \mathcal{Y}} K \sqrt{g(y)} - \sum_{y \in \mathcal{Y}} \sigma(N) g(y) \\
&\geq K \Gamma(Y_j^i, Y_{j'}^i) - |\mathcal{Y}| \sigma(N) \\
&\stackrel{(a)}{\geq} K \sum_{y \in \mathcal{Y}} \sqrt{s'(y)} - |\mathcal{Y}| \sigma(N). \tag{24}
\end{aligned}$$

Here, (a) follows from Proposition 3. Finally, (24) and (23) together imply that, for a large enough  $N$ ,

$$G(\mathbf{t}) \geq G(\mathbf{q}) - |\mathcal{Y}| (\sigma(N) + \sigma'(N)).$$

Thus for any  $\omega > 0$ , there exists some  $N_2$  such that for any  $N \geq N_2$ , the payoff under the truthful equilibrium is less than that under any other symmetric strategy profile by at most  $\omega$ . Taking  $N_0 = \max(N_1, N_2)$  proves our claim.  $\square$

*Proof of Theorem 5.* As before, we denote  $e_j^t(y)$  to be the agreement scores for an agent  $j$  under a

fully informative equilibrium, and  $e_j^s(y)$  to be the scores under a fixed symmetric strategy profile where each agent follows the reporting strategy  $\mathbf{q}$ . We denote  $\sum_{x \in \mathcal{X}} P_X(x) (\sum_{y' \in \mathcal{Y}} p_{y'}(x) q_y(y'))^2 \triangleq s(y)$  and denote  $\sum_{x \in \mathcal{X}} P_X(x) \sum_{y' \in \mathcal{Y}} p_{y'}(x) q_y(y') \triangleq b(y)$ . By our assumption  $b(y) \geq \eta$ , and since  $s(y) \geq b(y)^2$ , we have  $s(y) \geq \eta^2$ . Then using arguments similar to the ones leading up to (20) in the proof of Theorem 2, we can show that for all  $y \in \mathcal{Y}$ ,  $|\mathbb{E}(e_j^t(y)) - \frac{K}{\sqrt{g(y)}}| \leq \sigma(N) = o(1)$ , and  $|\mathbb{E}(e_j^s(y)) - \frac{K}{\sqrt{s(y)}}| \leq \sigma'(N) = o(1)$ , where  $\sigma(N) \geq 0$  is some function of  $N$  that depends only on  $\delta(P_X, \mathbf{p})$ ,  $n$  and  $K$ , and  $\sigma'(N) \geq 0$  is some function of  $N$  that depends only on  $\delta(P_X, \mathbf{p})$ ,  $\eta$ ,  $n$  and  $K$ . Neither of these functions depend on  $y$ . Let  $G(\mathbf{t})$  and  $G(\mathbf{q})$  be the expected payments to agent  $j$  for task  $i$  under the truthful strategy profile and the symmetric profile  $\mathbf{q}$ . Let  $j'$  be  $j$ 's peer for task  $i$ . Let  $Z_j^i$  and  $Z_{j'}^i$  be the reported answers of  $j$  and  $j'$  for task  $i$  under  $\mathbf{q}$ . Then we have,

$$\begin{aligned} G(\mathbf{q}) &= \sum_{y \in \mathcal{Y}} \mathbb{E}(e_j^s(y)) s(y) \\ &\leq \sum_{y \in \mathcal{Y}} K \sqrt{s(y)} + \sum_{y \in \mathcal{Y}} s(y) \sigma'(N) \\ &\leq K\Gamma(Z_j^i, Z_{j'}^i) + |\mathcal{Y}| \sigma'(N). \end{aligned} \quad (25)$$

Similarly, we can show that

$$\begin{aligned} G(\mathbf{t}) &= \sum_{y \in \mathcal{Y}} \mathbb{E}(e_t^s(y)) g(y) \\ &\geq \sum_{y \in \mathcal{Y}} K \sqrt{g(y)} - \sum_{y \in \mathcal{Y}} g(y) \sigma(N) \\ &\geq \Gamma(Y_j^i, Y_{j'}^i) - |\mathcal{Y}| \sigma(N) \\ &\geq \Gamma(Z_j^i, Z_{j'}^i) + \frac{\delta(P_X, \mathbf{p}) \Omega(\mathbf{q})^2 (|\mathcal{Y}| - 1)}{2\sqrt{|\mathcal{Y}|}} - |\mathcal{Y}| \sigma(N). \end{aligned} \quad (26)$$

Thus if  $G(\mathbf{q}) \geq G(\mathbf{t})$  for any strategy  $\mathbf{q}$ , then this implies that,

$$K\Gamma(Z_j^i, Z_{j'}^i) + \frac{\delta(P_X, \mathbf{p}) \Omega(\mathbf{q})^2 (|\mathcal{Y}| - 1)}{2\sqrt{|\mathcal{Y}|}} - |\mathcal{Y}| \sigma(N) \leq K\Gamma(Z_j^i, Z_{j'}^i) + |\mathcal{Y}| \sigma'(N),$$

which implies that

$$\frac{\delta(P_X, \mathbf{p}) \Omega(\mathbf{q})^2 (|\mathcal{Y}| - 1)}{2\sqrt{|\mathcal{Y}|}} \leq |\mathcal{Y}| (\sigma(N) + \sigma'(N)),$$

or that,

$$\Omega(\mathbf{q}) \leq \sqrt{\frac{2|\mathcal{Y}|^{3/2} (\sigma(N) + \sigma'(N))}{\delta(P_X, \mathbf{p}) (|\mathcal{Y}| - 1)}}. \quad (27)$$

Now the quantity on the right is  $o(1)$ . Thus for any  $\omega > 0$ , there exists some  $N_0$  such that for any  $N \geq N_0$ , any symmetric strategy profile that gives a higher expected payoff to each agent than the truthful strategy profile is at most  $\omega$ -uninformative. Since truthful reporting is a Bayes-Nash equilibrium for a large enough  $N$ , this implies the result.  $\square$

*Proof of Proposition 6.* We have,

$$\begin{aligned}
\sum_{y \in \mathcal{Y}} \sqrt{P(Z_1 = Z_2 = y)} &= \sum_{y \in \mathcal{Y}} \sqrt{\sum_{x \in \mathcal{X}, y_1 \in \mathcal{Y}, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) q_y(y_1) q'_y(y_2)} \\
&= \sum_{y \in \mathcal{Y}} \sqrt{\sum_{y_1 \in \mathcal{Y}, y_2 \in \mathcal{Y}} q_y(y_1) q'_y(y_2) \sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x) p_{y_2}(x)} \\
&\stackrel{(a)}{\leq} \sum_{y \in \mathcal{Y}} \sqrt{\sum_{y_1 \in \mathcal{Y}, y_2 \in \mathcal{Y}} q_y(y_1) q'_y(y_2) \left( \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x)^2} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_2}(x)^2} \right)} \\
&= \sum_{y \in \mathcal{Y}} \sqrt{\left( \sum_{y_1 \in \mathcal{Y}} q_y(y_1) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x)^2} \right) \left( \sum_{y_2 \in \mathcal{Y}} q'_y(y_2) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_2}(x)^2} \right)} \\
&\stackrel{(b)}{\leq} \frac{1}{2} \sum_{y \in \mathcal{Y}, y_1 \in \mathcal{Y}} q_y(y_1) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x)^2} + \frac{1}{2} \sum_{y \in \mathcal{Y}, y_2 \in \mathcal{Y}} q'_y(y_2) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_2}(x)^2} \\
&= \frac{1}{2} \sum_{y_1 \in \mathcal{Y}} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x)^2} + \frac{1}{2} \sum_{y_2 \in \mathcal{Y}} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_2}(x)^2} \\
&= \frac{\Gamma(Y_1, Y_2)}{2} + \frac{\Gamma(Y_1, Y_2)}{2} \\
&= \Gamma(Y_1, Y_2)
\end{aligned}$$

Here, (a) follows from the Cauchy-Schwarz inequality, and (b) results from the fact that the arithmetic mean of two numbers is no less than the geometric mean.  $\square$