

# The Square Root Agreement Rule for Incentivizing Truthful Feedback on Online Platforms

Vijay Kamble  
The University of Illinois at Chicago  
*kamble@uic.edu*

Nihar Shah  
Carnegie Mellon University  
*nihars@cs.cmu.edu*

David Marn  
University of California, Berkeley  
*marn@berkeley.edu*

Abhay Parekh  
University of California, Berkeley  
*yahbaa@gmail.com*

Kannan Ramchandran  
University of California, Berkeley  
*kannanr@eecs.berkeley.edu*

A major challenge in obtaining evaluations of products or services on e-commerce platforms is eliciting informative responses in the absence of verifiability. This paper proposes the *Square Root Agreement Rule* (SRA): a simple reward mechanism that incentivizes truthful responses to objective evaluations on such platforms. In this mechanism, an agent gets a reward for an evaluation only if her answer matches that of her peer, where this reward is inversely proportional to a *popularity index* of the answer. This index is defined to be the square root of the empirical frequency at which any two agents performing the same evaluation agree on the particular answer across evaluations of similar entities operating on the platform. Rarely agreed-upon answers thus earn a higher reward than answers for which agreements are relatively more common.

We show that in the many tasks regime, the truthful equilibrium under SRA is strictly payoff-dominant across large classes of natural equilibria that could arise in these settings, thus increasing the likelihood of its adoption. While there exist other mechanisms achieving such guarantees, they either impose additional assumptions on the response distribution that are not generally satisfied for objective evaluations or they incentivize truthful behavior only if each agent performs a prohibitively large number of evaluations and commits to using the same strategy for each evaluation. SRA is the first known incentive mechanism satisfying such guarantees without imposing any such requirements. Moreover, our empirical findings demonstrate the robustness of the incentive properties of SRA in the presence of mild subjectivity or observational biases in the responses. These properties make SRA uniquely attractive for administering reward-based incentive schemes (e.g., rebates, discounts, reputation scores, etc.) on online platforms.

---

## 1. Introduction

Reputation systems, in which people provide feedback for products or services based on their personal experiences, are a critical component of online platforms and marketplaces (Resnick et al. 2000, Jøsang et al. 2007, Tadelis 2016, Luca 2017). These systems improve the overall quality of transactions, increase trust, and thus play a key 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 even exploit these systems for selfish motives, thus undermining their utility. For instance, significant empirical evidence of bias in user ratings has been found on many online platforms (Hu et al. 2017, Filippas et al. 2018, Nosko and Tadelis 2015). This work describes a simple and intuitive reward mechanism that attempts to address this concern.

We consider a setting where an online platform is interested in obtaining responses for a large number of evaluations pertaining to the products or the services being offered on the platform from a pool of customers, whom we refer to as agents. We focus on *objective but unverifiable evaluations*, i.e., evaluations in which the answers can, in principle, be objectively verified, but such verification is infeasible for the platform. 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 plumber show up within 5 mins of your appointed time? (Yes/No)
3. How long did the moving company take to respond with a quote? (1 day/2 days/3 days/more than three days)
4. Did the dimensions of the received product exceed the dimensions given by the seller? (Yes/No)
5. How long did it take for the product to arrive after the purchase was made from the seller? (1 week/2 weeks/3 or more weeks)

In each of these questions, the evaluating agent is being asked to truthfully report an *observation* about the entity being evaluated. The main property of such evaluations is that each evaluating agent’s observation is an independent sample from an unknown distribution of behaviors specific to the entity being evaluated. In other words, the true responses of agents for a fixed evaluation task are *conditionally independent and identically distributed* (conditional on the unknown distribution of responses). For example, in the first situation, we can assume that each customer experiences an independently sampled waiting time from a common unknown distribution specific to that restaurant. In the second situation, the customer’s experience is sampled from the distribution of whether or not the plumber is punctual. Similarly, in the remaining questions, the customer’s true experience is a sample of the moving company’s or the seller’s business practices. We will refer to this property

as the responses being *homogeneous* for the rest of the paper, informally referring to the fact that the true responses of any set of agents to a fixed evaluation task are statistically exchangeable.

In such scenarios, we hope to achieve the following informal goals through the design of an effective incentive mechanism: (a) incentivize agents to participate in the provision of feedback in online platforms, i.e., improve response rates, and (b) conditional on participation, incentivize agents to report true observations while overcoming any observation or reporting bias. The second goal is arguably more critical and challenging since an easy way of achieving the first goal is to give everyone a fixed reward for participation. As one can imagine, such a naïve reward scheme may not necessarily lead to a high quality of responses.

If the platform could verify the responses to the evaluations, it can simply reward the agents based on whether or not they reported their true observations. But such verification is infeasible for questions such as the ones mentioned above since these evaluations are based on interactions that take place outside the platform. In these cases, inducing truthful behavior is a challenging problem. A common approach to this problem, first described in the pioneering work of Miller et al. (2005), is to reward the agents' responses based on comparisons with the responses of other agents who have performed the same evaluation task. Such mechanisms have come to be referred to as *peer-prediction mechanisms* in subsequent literature (after the original mechanism called the *peer-prediction method* described in Miller et al. (2005)). Informally, such reward mechanisms leverage the property that the true response of any agent is correlated with the response of some other agent for the same question.<sup>1</sup> The situation is then inherently strategic, in which one hopes to sustain truthful reporting as an equilibrium of the game that the reward mechanism induces. It is additionally desirable that such an equilibrium is preferable to the agents over other non-truthful equilibria that may arise in the game.

**Our contribution.** In this paper, we propose the *Square Root Agreement Rule* (SRA): a new peer-prediction mechanism for online platforms that truthfully elicits objective but unverifiable responses at equilibrium. In the setting of our interest, i.e., elicitation on online platforms, we show that the truthful equilibrium under SRA satisfies a key dominance property; namely, it yields the agents the highest payoff amongst all symmetric equilibrium payoffs in the system limit where there are a large number of evaluation tasks. In addition, the truthful equilibrium payoff is *strictly* higher than that under any symmetric equilibrium strategy profile that incurs any degree of information loss in the

<sup>1</sup> With homogeneous responses, the structure of the correlation between an agent's true response and the true response of a typical agent in the population is identical across agents. This feature contrasts with the case when the responses are *heterogeneous*, i.e., when the agents' true responses strongly depend on their characteristics that vary widely across the population. In these cases, designing mechanisms without obtaining requisite fine-grained information about agent heterogeneity or without making any regularity assumptions on the agent responses, e.g., 'self-predicting responses' (Radanovic et al. 2016) or 'categorical responses' (Dasgupta and Ghosh 2013, Shnayder et al. 2016) (these assumptions are discussed in Section 2), is known to be impossible (Radanovic and Faltings 2015).

reports. In keeping with the existing terminology in the literature, we refer to this property as SRA being asymptotically *strongly truthful across symmetric equilibria*. Such a dominance property is crucial in these settings since it hinders the emergence of low-effort equilibria with poorly informative reports – such as everyone reporting the same answer irrespective of their true evaluation – known to plague many other incentive mechanisms, e.g., the peer-prediction method of Miller et al. (2005).

Moreover, under certain additional assumptions satisfied in applications such as crowdsourcing,<sup>2</sup> we show that the truthful equilibrium under SRA gives the highest payoff amongst *all* equilibrium payoffs (and not just symmetric equilibrium payoffs) in the large system limit, i.e., SRA is asymptotically *strongly truthful* under these assumptions.

While such strong truthfulness guarantees are satisfied by existing mechanisms, they either impose conditions on the response distributions that are not satisfied in general for objective evaluations (Dasgupta and Ghosh 2013, Radanovic et al. 2016) or they require a prohibitively large number of evaluations from each agent and assume that the agent uses the same reporting strategy for each evaluation (Kong and Schoenebeck 2019, Kong 2020). SRA is the first known mechanism that achieves this guarantee for objective evaluations without imposing any such constraints; in particular, SRA is the first mechanism satisfying strong truthfulness guarantees for objective evaluations that incentivizes truthful behavior even among agents who perform a *single* evaluation.

This result is arguably non-trivial. The dominant existing framework for designing strongly truthful peer-prediction mechanisms is due to Kong and Schoenebeck (2019), which incentivizes truthful behavior only if the agents perform multiple evaluations (at least twice the number of possible responses to any evaluation; see Kong (2020)) and additionally commit to using the same reporting strategy for each evaluation. In Section 2 and Section D.4 in the Appendix, we show that even if one leverages the homogeneous responses property satisfied by objective evaluations, it is not possible to generically adapt the approach of Kong and Schoenebeck (2019) to incentivize truthfulness in a single evaluation. By designing SRA, we nevertheless demonstrate that this is indeed possible. In showing this result, we make novel information-theoretic contributions that are of interest beyond this work.

The fact that SRA strongly incentivizes even a single response is vital since requiring multiple evaluations from any agent to incentivize truthfulness is impractical in platforms where customers interact with the marketplace relatively rarely, e.g., in vacation rental platforms such as Airbnb. In these settings, requiring a single additional evaluation from each user may prohibitively increase the evaluation period’s duration. Such delay and resulting changes in market characteristics over time

<sup>2</sup> Crowdsourcing on labor platforms such as Amazon Mechanical Turk is an important means for sourcing the large-scale execution of information-oriented micro-tasks, such as obtaining labeled data for training machine learning algorithms. Incentivizing truthful, high-quality responses from participants is a key concern in these applications.

undermine the platform’s ability to procure the most up-to-date feedback information reliably. While there are mechanisms that incentivize truthfulness in a single evaluation under the homogeneous responses assumption (see Section 2 for a discussion), none of these mechanisms satisfy the crucial *strong* truthfulness guarantees that SRA satisfies.

**Background.** To appropriately position our contribution, we first present a brief discussion of the main existing approaches in incentive design for the elicitation of unverifiable responses.

In a pioneering work in this domain, Miller et al. (2005) considered the case of a single evaluation task and homogeneous responses and described the so-called *peer-prediction method* that incentivizes truthful answers. The main requirement is that there is a commonly known prior on the unknown distribution from which the agents’ true observations are sampled, and this prior is known to the principal (who, in our case, is the platform). Truthfulness is achieved by rewarding/scoring an agent’s posterior probability distribution of her peer’s answer conditioned on her own answer, using a *Proper Scoring Rule* (PSR). PSRs are a well-known class of payment/scoring rules that incentivize truthful elicitation of probabilities of events that can be observed at a future date (Brier 1950, Gneiting and Raftery 2007, Savage 1971). This approach is infeasible in platforms since a prior on the distribution of evaluations is typically not a priori available, and even if it is, it may not be common knowledge across all agents and the platform. Moreover, this mechanism is known to induce uninformative equilibria that yield the agents a higher payoff than the truthful equilibrium payoff (Jurca and Faltings 2005).

Another influential design in this domain, the Bayesian Truth Serum (BTS) (Prelec 2004), and its subsequent refinements and generalizations (Witkowski and Parkes 2012, Radanovic and Faltings 2013, Schoenebeck and Yu 2020), preserved the common prior assumption but relaxed the requirement that the principal needs to know this prior. These mechanisms instead require the agents to make extraneous reports about their beliefs in addition to their answers. In particular, they are asked to report a prediction of the empirical distribution of answers reported by other agents. Again, requiring customers on platforms to provide such extraneous information about their beliefs is a tall order given the already low response rates seen for simpler forms of feedback, e.g., ratings. Unfortunately, such extraneous reports of beliefs, although undesirable, are indispensable for incentivizing a single evaluation; it has been shown that it is impossible to design mechanisms that incentivize truthfulness without obtaining some information about the prior distribution (Jurca and Faltings 2011) (which is obtained via agents’ reports of their beliefs in BTS).

Mechanisms that do not assume that the principal knows the prior on the distribution of answers have been referred to as *detail-free* in the literature and those that do not require any extraneous reporting from the agents apart from the evaluations are called *minimal*. Ideally, for reputation systems, we need incentive mechanisms that are *both* detail-free and minimal, and that do not rely on

the assumption of the existence of a commonly known prior across agents. In light of the impossibility result mentioned above, this seems like a challenging task, if not entirely impossible. The earliest known minimal and detail-free incentive mechanism for a single evaluation task was designed by Jurca and Faltings (2008), in which respondents arrive sequentially and the distributional knowledge of responses is obtained and leveraged by the mechanism in an online fashion; however, this mechanism critically relies on the assumption of binary evaluations.

This is where a key feature of online platforms can be exploited: they typically host a *large number of similar products or services*. For instance, there are thousands of similar restaurants listed on review platforms like Yelp that users rate. Online marketplaces like Amazon or eBay would like to obtain reviews for many existing sellers on these platforms. Online labor platforms like Thumbtack and Handy would like to get performance metrics for thousands of workers and service providers that operate on these platforms.

The presence of multiple similar evaluation tasks hints at an approach for designing detail-free mechanisms that are also minimal: the missing information about the prior can be obtained from consistent statistical estimates of the distribution of agent responses derived from the response data across multiple tasks. Witkowski and Parkes (2013) first explored such a possibility in the context of crowdsourcing, for eliciting binary (e.g., yes/no) responses. A potential concern with this approach is that it assumes that the response data is truthfully generated. But it turns out that in these situations, with careful design, truthfulness can become a *self-fulfilling prophecy* – truthful behavior is an equilibrium in the induced game when the mechanism assumes that these reports are truthful. This is the basic principle underlying the design of Witkowski and Parkes (2013) that forms the foundation of our design, resulting in the fact that SRA is both minimal and detail-free. Mechanisms exploiting this principle are commonly referred to as *multi-task* peer prediction mechanisms in the literature.

**Structure of SRA.** SRA is a *multi-task* peer-prediction mechanism that builds upon the structure of *output agreement mechanisms* (Von Ahn and Dabbish 2008, 2004), which are simple and intuitive mechanisms that have been quite popular in crowdsourcing practice, except they suffer from a critical drawback of being susceptible to strategic manipulations. In an output agreement mechanism, two agents answer the same question, and they are both rewarded if their answers match. There are two critical drawbacks of this scheme: (a) truthful behavior may not necessarily be an equilibrium (see Section 4.1 for an example) and (b) there is always an undesirable equilibrium in the game it induces, in which every person reports the same answer irrespective of their true evaluation. This equilibrium guarantees each person the highest possible payoff rewarded by the mechanism. Our mechanism overcomes these drawbacks by giving proportionately lower rewards for answers that turn out to be more popular on other similar evaluation tasks. This is achieved by

inversely scaling the rewards for agreement by a *popularity index* for each answer, thus discouraging blind convergence on a single answer. This is not a new idea: such a biased output agreement scheme, called the Peer Truth Serum (PTS), was first introduced by Jurca and Faltings (2011), and was further refined by Radanovic et al. (2016) and Faltings et al. (2017).

The key innovation in our design is in the way these popularity indices are defined. All the strong incentive properties of our mechanism trace their origin to this 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 of 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 amongst agents that the prior satisfies a certain non-degeneracy property.

**Strong truthfulness.** A crucial concern in any reward mechanism 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. Indeed, the theory of equilibrium selection, i.e., justifying certain equilibria as more likely to arise than others, occupies an important position in game theory; see Harsanyi et al. (1988) and Van Damme (2002), and references therein. It is known that elicitation mechanisms for a single evaluation task with no extraneous reporting (which includes Miller et al. (2005)) possess uninformative equilibria that give a higher expected payoff to each agent than in the truthful equilibrium (Jurca and Faltings 2005). Moreover, these equilibria involve simple strategies such as every agent reporting the same answer, due to which these mechanisms are particularly vulnerable to uninformative feedback. The Bayesian Truth Serum demonstrated that this issue could be overcome by requiring extraneous reports of beliefs; the truthful equilibrium under BTS gives the highest expected payoff to an agent across all equilibria, and in particular, this payoff is strictly higher than that under any equilibrium strategy profile that is not fully informative. Mechanisms that satisfy this property are called *strongly truthful* mechanisms in the literature. Dasgupta and Ghosh (2013) first showed that such strong truthfulness properties could be obtained in the multi-task setting without requiring extraneous reports from the agents. Kong and Schoenebeck (2019) describe a general information-theoretic analysis of incentive mechanisms in this space, and show 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.

We show that SRA achieves 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 reporting strategy). As the number of tasks grows, asymptotically, the expected payoff in the limit under a truthful strategy profile is higher than that under any other symmetric strategy profile. Moreover, this payoff is strictly higher than that under any symmetric equilibrium strategy profile that is not *fully informative*. A fully informative strategy profile is one where each agent applies a common permutation map to her observation to generate her report (essentially amounting to relabeling the set of responses). In other words, SRA is asymptotically *strongly truthful across all symmetric equilibria*. As a dual to this property, under a mild assumption on the strategy spaces, we also show that any symmetric equilibrium that gives the highest expected reward to an agent across all symmetric equilibria must be close, in a well-defined sense, to being fully informative when the number of tasks is large. Such properties hinder the rise of potential “obviously attractive” symmetric equilibria such as all agents reporting the same answer for every evaluation.

The restriction to symmetric equilibrium payoffs in the equilibrium dominance property of SRA may seem undesirable. However, this restriction stems from a crucial difference in our setting compared to the settings considered by other mechanisms that are mainly motivated by crowdsourcing applications. In our setting, identifying information for the different entities to be evaluated is available to the agents (captured by the task number in our formal setup). Moreover, allocations of evaluation tasks to the agents are exogenously specified. Thus, in our setting, the agents are free to choose reporting strategies that depend on the identities of the tasks they perform to coordinate their reports with other agents, in addition to the (desirable) coordination that can be achieved by reporting their observations truthfully. Due to the possibility of such extraneously achieved coordination, it is well known that it is impossible to elicit truthful evaluations under a payoff-dominant truthful equilibrium in general (Gao et al. 2019).<sup>3</sup> Other mechanisms get around this difficulty by making certain assumptions that eliminate the possibility of the agents choosing reporting strategies that depend on task identity. Although such assumptions may be justifiable in applications such as crowdsourcing, we do not rely on such assumptions since they are inappropriate in the context of eliciting feedback on platforms; see Section 5.3 for a discussion. In Section 5.3, we additionally show that if such assumptions are made, then SRA is indeed asymptotically strongly truthful.

In a similar spirit as the framework of Kong and Schoenebeck (2019), these strong truthfulness guarantees are obtained by showing that the expected payoff of an agent under a particular strategy

<sup>3</sup> The argument for this impossibility is the following: suppose that there is a mechanism that ensures that the payoffs that agents obtain by truthfully signaling an extraneous feature of an entity being evaluated via their reports are always lower than those obtained by truthfully reporting the feature that they are supposed to report. Then one obtains a contradiction by exchanging the role of the extraneous feature and the feature to be reported.

profile under SRA is proportional to a novel notion of a *square root agreement measure* (SRAM) between two independent responses, which we show to be monotonically decreasing in unilateral information loss in the responses. Both the measure and this monotonicity property are new and of independent interest. Moreover, in Section 2, we show that a generic adaptation of the framework of Kong and Schoenebeck (2019) to the homogeneous response setting along the lines of SRA does not yield truthfulness under any arbitrary mutual information measure. SRAM is thus, arguably, the “right” notion of an agreement measure for objective evaluations. We discuss this aspect in more detail in Section 2.

**Robustness.** We finally perform numerical experiments on synthetic as well as real data to test SRA’s robustness in incentivizing truthful behavior in finite-data settings featuring deviations from the homogeneous responses assumption. This is practically important since, despite making a faithful effort to obtain and report true observations, agents may have residual biases in their observations and reports. Such biases could also capture mild subjectivity in responding to objective evaluations. Hence, it is desirable that SRA incentivizes each agent to be truthful even when the agent accounts for such biases in other agents. We find that SRA exhibits a high degree of robustness to these practical concerns and generates strong incentives for truthful behavior.

**Reward mechanisms in practice.** Non-monetary rewards for incentivizing informative feedback, e.g., coupons, badges, or some form of a reputation score, are commonly seen in crowdsourced review forums like Yelp, Tripadvisor, etc. A prominent example of monetary incentives is the “Rebate for Feedback (RFF)” program that was launched by Taobao.com (one of the world’s largest e-commerce websites), on March 1, 2012.<sup>4</sup> In this program, sellers can set a rebate amount in the form of cashback or a store coupon for any items they sell, as a reward for a buyer’s feedback after purchasing that item. If a seller opts for RFF, then Taobao ensures that the rebate is transferred from the seller’s account to a buyer who leaves high-quality feedback. The feedback quality is determined by a machine learning algorithm depending on attributes like the length of the feedback, whether or not certain key features of the item (e.g., longevity, whether or not it is true to size, etc.) are mentioned, etc. The main contention of the present work is that strategic considerations are paramount in incentivizing *informative* feedback. For example, it is easy to give untruthful feedback that appears to be of high quality to a machine learning algorithm; this is especially a concern for objective evaluations with a fixed, finite set of answers. SRA can thus be an effective approach to administer such rebate schemes in a manner that is robust to strategic behavior.

**Organization of the paper.** The remainder of the paper is organized as follows. In Section 2, we discuss related mechanisms and their comparisons with SRA. Section 3 presents a formal description

<sup>4</sup> See <https://bit.ly/2GVntzC>, and also Li et al. (2020).

Mechanism	Incentivizes truthful homogeneous responses	Incentivizes single evaluations	Dominance property for truthful equilibrium
<b>SRA</b>	✓	✓	Strongly truthful across symmetric equilibria
Vanilla output agreement	✗	✓	None
Witkowski and Parkes (2013)	✗ (needs binary responses)	✓	None
Dasgupta and Ghosh (2013)	✗ (needs categorical responses)	✓	Strongly truthful across symmetric equilibria
Peer Truth Serum for Crowdsourcing (PTSC)	✗ (needs self-predicting responses)	✓	Strongly truthful across symmetric equilibria
Kong and Schoenebeck (2019)	✓	✗	Strongly truthful* (see Section 3)
Correlated Agreement (CA)	✓	✗ ( $\geq 2$ )	Informed truthful across symmetric equilibria
CA-HR (Appendix Section D.1)	✓	✓	Informed truthful across symmetric equilibria
Multi-task Peer Prediction Method	✓	✓	None
Radanovic and Faltings (2015)	✓	✓	None

**Table 1** Properties of different multi-task peer-prediction mechanisms in the many tasks regime in our setting.

of the model considered in the paper. Section 4 presents the SRA mechanism and its main incentive property. In Section 5 we address the issue of equilibrium selection. In Section 6, we perform numerical experiments in a practically motivated experimental setup to test the robustness of SRA in incentivizing truthful behavior to deviations from our main assumptions. We finally summarize our contributions and conclude in Section 7. The proofs of all of our results can be found in the Appendix.

## 2. Related literature

As a minimal, detail-free, multi-task mechanism, the key feature of SRA is that it *strongly* incentivizes truthful responses in homogeneous response settings, even among agents who have performed a single evaluation. We now discuss this property in relation to the properties satisfied by other existing multi-task mechanisms. Table 1 summarizes the differences between SRA and these mechanisms at a high-level.

### 2.1. Existing strongly truthful multi-task mechanisms

We first discuss mechanisms that achieve strong truthfulness guarantees, focusing on distinctions from SRA.

1. **Dasgupta and Ghosh (2013)**. Dasgupta and Ghosh (2013) proposed the first known detail-free and minimal multi-task peer-prediction mechanism that is also strongly truthful, assuming that agents do not choose reporting strategies contingent on task identities; in the absence of this assumption, it is strongly truthful across symmetric equilibria. The original paper restricted the setting to binary evaluations; however, it was later shown by Shnayder et al. (2016) that the mechanism is strongly truthful for non-binary responses under the condition that the responses are “categorical.” This condition says that conditional on an agent’s answer, the posterior probability of all other agents’ answers 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 that 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 under homogeneous responses (see Remark 3 in Section 6).

2. **Peer Truth Serum for crowdsourcing (PTSC) (Radanovic et al. 2016)**. PTSC is the multi-task extension of the PTS mechanism, originally defined by Jurca and Faltings (2011) for the case where the prior distribution of responses is known to the principal. Both these mechanisms operate in the homogeneous responses setting. PTS has a biased output agreement structure requiring only one evaluation per agent, where the popularity index of each answer is the prior probability of an agent reporting that answer. In PTSC, this prior probability is replaced by its estimate computed from the response data obtained from a large number of tasks. In order to obtain truthfulness, PTS/PTSC requires that the agent responses satisfy a “self-prediction” condition, which says that  $Pr(y|y)/Pr(y) \geq Pr(y'|y)/Pr(y')$  for any  $y' \neq y$ . This is equivalent to saying that  $Pr(y|y) \geq Pr(y|y')$  for any  $y' \neq y$ . We can show that this condition is weaker than the categorical responses condition required by the mechanism of Dasgupta and Ghosh (2013) for incentivizing truthfulness; see Proposition E.1 in the Appendix. However, except for the case of binary evaluations, this condition is also not satisfied in general under homogeneous responses (see Remark 3 in our numerical evaluations). If this condition is satisfied, PTSC has been shown to be asymptotically strongly truthful while restricting to symmetric strategy profiles, or in other words, it is strongly truthful across symmetric equilibria.

3. **Kong and Schoenebeck (2019)**. The underlying principle leading to the strong truthfulness property of SRA is closely related to the mechanism design paradigm of Kong and Schoenebeck (2019) (KS), who provide an information-theoretic framework for designing strongly truthful mechanisms for general settings with non-homogeneous responses. Their mechanism operates on a pair of agents, and the payment to each agent is defined to be a scaling of an unbiased estimate of some mutual information measure between the two agents’ response distributions, constructed using their reported responses to a common set of evaluation tasks. A variety of mechanisms can be obtained by varying the information measure. Strong truthfulness follows from the fact that the mutual information measure is monotonically decreasing with respect to loss of informativeness in the agents’ responses. We note that strong truthfulness here rests on the assumption that agents use the same reporting strategy for all tasks. However, in our setting, since agents are allowed to choose reporting strategies contingent on task identities, this mechanism is strongly truthful only across equilibria in which agents choose a common reporting strategy for all tasks that they perform (which includes symmetric equilibria).

Until recently, all known mutual information measures required an unboundedly large number of responses per agent (a large fraction of which must be commonly performed by the two agents)

to construct unbiased estimates.<sup>5</sup> Kong (2020) recently proposed a mutual information measure, of which an unbiased estimator can be constructed using a finite number of per-agent responses. Nevertheless, this construction still requires a number of per-agent responses at least twice the number of possible answers (e.g., 4 responses for binary evaluations). In Section D.3 in the Appendix, we show that there can be no mutual information measure satisfying information monotonicity whose unbiased estimate can be constructed from two agents’ responses to a single evaluation task, even in the homogeneous, binary responses setting. Thus, the KS mechanism design framework fails to incentivize truthful behavior in a single evaluation even in the homogeneous responses setting.

Via the design of SRA, we effectively show that if one leverages distributional information obtained from multiple tasks in the homogeneous responses setting, then there is a mutual information measure (the square root agreement measure) and a deviation from the KS mechanism design framework that utilizes the learned distributional information along with agent responses to compute payments, such that the mechanism strongly incentivizes truthful behavior in even a single evaluation.

Based on SRA’s design, one may wonder if it is possible to obtain a generic adaptation of the KS mechanism to the multi-task, homogeneous responses setting, which utilizes distributional information (obtained from many tasks) to incentivize truthful single responses under *any* mutual information measure. In Section D.4 in the Appendix, we show that this is not true by considering the Shannon mutual information (Cover and Thomas 2012): we show that a mechanism along the lines of SRA that leverages the Shannon mutual information instead of the square root agreement measure is not truthful in general for homogenous responses. This underscores the importance of the discovery of SRAM and shows that it is arguably the “right” agreement measure for the purpose of strongly incentivizing objective evaluations.

## 2.2. Other prominent multi-task mechanisms

We now discuss multi-task mechanisms that are not known to be strongly truthful in general.

1. **Witkowski and Parkes (2013)**. The mechanism proposed by Witkowski and Parkes (2013) is minimal and requires each agent to perform only one task; however, it requires that the responses are binary, and hence it is not uniformly applicable to the homogeneous responses setting. Additionally, no equilibrium dominance properties are known for this mechanism.

2. **Correlated Agreement mechanism (CA) (Shnayder et al. 2016)**. CA is a multi-task mechanism that incentivizes truthful behavior in the general heterogeneous responses setting. CA operates on a pair of agents, and it requires at least two evaluations per agent. As originally

<sup>5</sup> Schoenebeck and Yu (2021) recently proposed a sample-efficient approach to directly learn an appropriate scoring-rule for scoring the agents’ reports that implements the mutual-information mechanism of Kong and Schoenebeck (2019), rather than learning the joint distribution and then estimating the mutual information.

described, it assumes that certain information about the joint distribution of the agents’ responses is known to the principal. However, this knowledge assumption can be relaxed in the multi-task setting with homogeneous responses, since this distribution can be estimated from the response data obtained from a large number of tasks (appealing to the self-fulfilling prophecy of truthful behavior). Additionally, the requirement of two evaluations per agent can also be relaxed: in Section D.1 in the Appendix, we describe an adaptation of CA to our setting that only requires one evaluation per agent. We call this mechanism CA for homogenous responses (CA-HR).

In the setting in which CA is originally defined, it is assumed that there is no extraneous identifying information for the evaluation tasks and the task allocation is randomized across agents, in effect eliminating the need to consider task-contingent reporting strategies of the kind allowed in our setting. In this setting, CA satisfies the dominance property of *informed truthfulness*, which is weaker than strong truthfulness. Under informed truthfulness, the truthful equilibrium yields the highest equilibrium payoff to each agent, which is higher than an agent’s payoff in any equilibrium where her reporting strategy is fully *uninformative*, i.e., her reports are independent of her observation. However, in contrast to strong truthfulness, there could be other equilibrium strategy profiles that are not fully informative, which yield the same payoff as the truthful equilibrium. In particular, although the informed truthfulness property precludes fully uninformative equilibria where all observations map to a single response, the CA mechanism remains vulnerable to equilibria where agents map smaller sets of responses to a single response (e.g. if the responses are  $\{1, 2, 3, 4\}$  then  $\{1, 2\}$  map to 1 and  $\{3, 4\}$  map to 4), which a strongly truthful mechanism precludes (i.e., strictly payoff-dominates). Such types of equilibria are payoff-equivalent to the truthful equilibrium under CA if the joint distribution of observations of a pair of agents for a fixed evaluation task is “clustered,” as defined in Definition 10 (from Shnayder et al. (2016)) in the Appendix. In our practically motivated experimental setup in Section 6, we show that instances with clustered observations are encountered with a very high frequency; see Remark 4.

In our setting, since reporting strategies contingent on task identities are allowed, CA and CA-HR are not informed truthful; we present an example Section D.2 in the Appendix illustrating this fact. But they can be shown to be (asymptotically) informed truthful *across symmetric equilibria*. On the other hand, in Section 5.3, we show that if strategies contingent on task identities are disallowed in our setting and tasks are randomly allocated across agents, the SRA is (asymptotically) *strongly truthful* across all equilibria. Thus, SRA satisfies a stronger equilibrium dominance property compared to CA or CA-HR for homogeneous responses.

**3. Multi-task extension of the peer-prediction method (Miller et al. 2005).** The peer-prediction method is a minimal mechanism and incentivizes truthful responses with only one evaluation per agent in the homogeneous responses setting. As originally described, it assumes that the joint distribution of the agents' responses is commonly known to the principal and the agents. However, this assumption can be relaxed in the multi-task setting with homogeneous responses since this distribution can be estimated from the response data obtained from a large number of tasks. As we discussed in Section 1, this mechanism achieves truthfulness by rewarding an agent's predicted probability distribution of her peer's answer, as implied by her own answer, using a Proper Scoring Rule (PSR).

This mechanism, however, does not satisfy any equilibrium dominance properties: in particular, there exist uninformative equilibria that yield the highest possible payoff to each agent, irrespective of the PSR utilized.<sup>6</sup> For example, for any PSR, the highest possible utility under this mechanism (in the many tasks limit) is achieved when all agents simply report the same answer irrespective of their observation. To see this, note that under such a symmetric strategy, in the many tasks limit, the joint distribution of responses estimated by the principal puts a unit probability mass on the event of the two peer agents reporting the fixed answer. Hence, the conditional distribution on the peer's response implied by an agent's response perfectly predicts the peer's response and thus, achieves the highest score. In contrast, we show that such equilibrium leads to a strictly lower payoff than the truthful equilibrium under SRA when the number of tasks is large enough (see Remark 2).

**4. Radanovic and Faltings (2015).** Radanovic and Faltings (2015) describe a mechanism for the homogeneous responses setting that only utilizes one evaluation per agent. This mechanism is a multi-task extension of the peer-prediction method that relaxes the common prior assumption: they show how an unbiased estimate of the payoff under the peer-prediction method utilizing the quadratic scoring rule (a PSR) can be constructed in the homogeneous responses setting with response data from a random but almost surely finite number of tasks. However, the mechanism inherits the chief concern regarding the peer-prediction method that we discuss above, in that the strategy profile where everyone reports a fixed answer irrespective of their observations results in the highest possible score and thus a higher payoff than the truthful equilibrium.

<sup>6</sup> This observation has already been made for the original non-detail-free version of the peer-prediction method by (Jurca and Faltings 2005).

### 3. Model

We now describe the details of our model.

**Operational details.** We consider a setting with  $N$  evaluation tasks denoted by the set  $\mathcal{N}$  and labeled as  $i = 1, \dots, N$ . Let  $\mathcal{M}$  denote the population of  $M$  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 possible observations and the set of possible answers in each evaluation task is assumed to be the same finite set, denoted as  $\mathcal{Y}$ . A generic element of  $\mathcal{Y}$  will be denoted as  $y$ .

**Statistical assumptions.** The distribution of observations of agents performing task  $i$  is specified by an unobservable type of the entity being evaluated in task  $i$ , denoted as the random variable  $X_i$ , which takes values in the finite set  $\mathcal{X}$ . A generic element of  $\mathcal{X}$  will be denoted as  $x$ . This set of possible types  $\mathcal{X}$  is common across all tasks. The distribution of the observations of the agents for any task, as a function of the type  $x \in \mathcal{X}$  of the entity being evaluated in that task, is denoted as  $\mathbf{p}(x) = (p_y(x); y \in \mathcal{Y})$ . The observation of an agent  $j$  in  $\mathcal{M}_i$  is denoted as the random variable  $Y_j^i$ , which is independently drawn from  $\mathbf{p}(X_i)$  for each such agent. This implies that the observations of different agents for a single task  $i$  are independent conditioned on  $X_i$ , but may be dependent otherwise.<sup>7</sup> Further, we assume that the types of entities being evaluated in the different tasks are independently sampled from a common distribution,  $P_X$ . We refer to this property as the tasks being *statistically similar*.

Finally, we assume that from the perspective of any agent  $j$ , there are no other observable features of the evaluation task  $i$  except the observation  $Y_j^i$ . The probability distribution over types,  $P_X$ , and the function  $\mathbf{p}$  together form a *probability generating model* (henceforth referred to as the generating model) of the agent observations, denoted as the pair  $(P_X, \mathbf{p})$ . In particular, this pair fully specifies a joint distribution on the underlying types of the different entities being evaluated and the different agents' observations across tasks. The following example illustrates this model.

**EXAMPLE 1.** Consider a situation where a labor platform wants to obtain feedback on punctuality of plumbers that operate on the platform. Suppose that each plumber could be of 2 possible unobservable types  $\mathcal{X} = \{\text{Punctual}, \text{Not Punctual}\}$ , with  $P_X(\text{Punctual}) = P_X(\text{Not Punctual}) = 0.5$ . Each plumber's type is independently sampled from the distribution  $P_X$ . The question is "Did the plumber show up within 5 minutes of his/her appointed time?". The two possible observations/answers are  $\mathcal{Y} =$

<sup>7</sup> Note that  $\{Y_j^i; j \in \mathcal{M}_i\}$  is a set of exchangeable random variables by the virtue of the fact that they are i.i.d. conditioned on the unobservable entity type  $X_i$ . Instead of explicitly assuming the existence of such a type, we can assume that the set of observations made by any and potentially an infinite number of agents for an entity  $i$  are exchangeable random variables, which a property that is expected to hold in practice for objective evaluations. De Finetti's theorem (Aldous 1985) would then imply the existence of a type for each entity such that the observations of the agents are conditionally i.i.d.

$\{\text{Yes}, \text{No}\}$ . And the distributions of true answers as a function of type are  $\mathbf{p}(\text{Punctual}) = (0.95, 0.05)$  and  $\mathbf{p}(\text{Not Punctual}) = (0.5, 0.5)$ .

We want to remark that the above example is merely illustrative: the type of plumber need not have any semantic interpretation. It may lie in some abstract space. More importantly, we emphasize that this type is unobservable, both, to the agents and the principal.

**The payment mechanism.** Our goal is to design a payment mechanism that elicits observations from the agents. For any  $j \in \mathcal{M}_i$ , let  $r_j^i$  denote agent  $j$ 's reported answer for task  $i$ . We define a payment mechanism as follows.

**DEFINITION 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).

**Agents' strategies.** 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 observation being  $y$ . We emphasize that agents are allowed to choose different reporting strategies for the different tasks, i.e., their reporting strategy can be task-contingent. We however restrict ourselves to considering task-contingent reporting strategies in which the reported answer of an evaluation only depends on the observation for *that* evaluation instead of potentially depending on the observations for all the other evaluations that the agent performs. This restriction is simply for the ease of exposition. All the incentive properties continue to hold for our proposed mechanism even if such strategies are allowed. This is simply because the payment in our mechanism to any agent is additive over the tasks that she performs. Hence, by the expectation operator's additivity, only the marginal distributions of the responses for the different tasks matter in determining the total expected payoff to an agent. Thus, choosing a reporting strategy for a given task that depends on the observations for other tasks is equivalent in terms of expected payoff to choosing the reporting strategy based on the output of some randomization device that produces values that are identically distributed to the observations for these other tasks. Such a reporting strategy is already included in the space of strategies for each agent.

**Equilibrium notion.** We define the following notion of a Bayes-Nash equilibrium (Myerson 2013) in the game induced by a payment mechanism.

**DEFINITION 2 (BAYES-NASH EQUILIBRIUM).** 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  $\{\bar{\mathbf{q}}^{ij} : i \in \mathcal{N}_j\} \neq \{\mathbf{q}^{ij} : i \in \mathcal{N}_j\}$ , 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 3 (BAYES-NASH INCENTIVE COMPATIBILITY).** 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 in the game induced by the mechanism. If this equilibrium is strict, we say that the mechanism is strictly Bayes-Nash incentive compatible.

**Informational assumptions.** We make the following informational assumptions.

1. The principal is *not* assumed to know  $P_X$  or  $\mathbf{p}$ . Hence,  $(P_X, \mathbf{p})$  is not an input to the payment mechanism.
2. We assume that the structure of the underlying generating model, i.e., the existence of some prior distribution  $P_X$  from which the type of any evaluated entity is drawn, and the function  $\mathbf{p}$  that captures the conditional distribution of observations given the type, that is common across entities being evaluated, is common-knowledge across all agents. In particular, this means that all the agents commonly know that all the tasks are statistically similar, and the observations of agents performing each evaluation are statistically homogeneous. Additionally, we will also assume that it is commonly known to all the agents that the generating model  $(P_X, \mathbf{p})$  satisfies a certain separation property, which we define in Section 4.2 (Definition 5), where we also discuss simple interpretations of this property. Finally, we assume that the payment mechanism, once fixed by the principal, is publicly announced and is commonly known to all agents.

Note that the definitions of Bayes-Nash equilibrium and Bayes-Nash incentive compatibility (Definitions 2 and 3) assume that the generating model is commonly known to the agents. However, we show that SRA is Bayes-Nash incentive-compatible with respect to any commonly known generating model that satisfies the separation property (discussed in Sections 4.2; Definition 5). Consequently, it is only necessary for the agents to commonly know that this separation property is satisfied by the generating model to obtain truthful behavior. Similarly, we show that the other properties we discuss concerning the Bayes-Nash equilibria in the game induced by SRA hold for *any* commonly known generating model that satisfies the separation property. Hence, to obtain these properties, we only need to assume that it is common knowledge amongst the agents that the generating model satisfies this property.

#### 4. The square root agreement rule (SRA)

Our main proposal, the square root agreement rule (SRA), is formally defined in Mechanism 1. Informally, the mechanism can be described as follows. Consider an agent  $j$  who has performed evaluation task  $i$ . Suppose that she submits an answer  $y \in \mathcal{Y}$ . Then, she receives payment for this answer only if another agent  $j'$ , who has performed the same task  $i$ , and who is chosen to be her peer, also reports the same answer  $y$ . This payment denoted as  $e_j(y)$  is inversely proportional to the square root of the empirical frequency at which arbitrarily chosen agents agree on answer  $y$  across all tasks that  $j$  has not performed; see Equation 3. This empirical frequency is denoted by  $\bar{f}_j(y)$  and is computed in Equation 2. To ensure that  $\bar{f}_j(y) > 0$  and the inverse is well defined, we use a smoothed version of empirical frequency, i.e., we add 1 to the total number of agreements on each answer before dividing by the number of tasks  $j$  has not performed. The following example illustrates the mechanism.

**MECHANISM 1: THE SQUARE ROOT AGREEMENT RULE (SRA).** ASSUMES  $|\mathcal{M}_i| \geq 2$  FOR ALL  $i \in \mathcal{N}$ .

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

$$\bar{f}_j(y) = \frac{1}{N - |\mathcal{N}_j|} \left( 1 + \sum_{i \in \mathcal{N} \setminus \mathcal{N}_j} \mathbf{1}_{\{r_{j_1(i)}^i = y\}} \mathbf{1}_{\{r_{j_2(i)}^i = y\}} \right). \quad (2)$$

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

$$e_j(y) = \frac{K}{\sqrt{\bar{f}_j(y)}}. \quad (3)$$

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  $i \in \mathcal{N}_j$ , 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.

**EXAMPLE 2 (SRA IN ACTION).** Consider the labor platform presented in Example 1. During an evaluation period, the platform solicits answers to the question “Did the plumber arrive within 5 minutes of his/her appointed time?” from all the customers who have hired a plumber from a set of a priori similar plumbers operating on the platform (e.g., new plumbers who have recently joined

the platform) during this period. Suppose a customer, Susan, reports an answer “Yes,” meaning that she reports that the plumber that she hired, Tim, did arrive within 5 minutes of his appointed time. Then, Susan gets a reward only if a randomly chosen customer who has also hired Tim in the same period also reports the answer “Yes.” The reward is inversely proportional to a popularity index of the answer “Yes” across the customer population computed from the response data. SRA defines this popularity index as the square root of the (smoothed) empirical frequency at which two customers who hire the same plumber both report the answer “Yes” across all the plumbers that Susan hasn’t evaluated. The payment procedure is similar if Susan reports “No” instead.

REMARK 1. Note that in SRA, a separate set of popularity indices for the different answers is computed for each agent based on population responses for tasks that this agent hasn’t performed. In practice, however, these indices are expected to be almost identical across agents when the total number of tasks is large relative to the number of tasks each agent performs; hence one can potentially calculate a single set of popularity indices of the answers and use them for all agents with negligible impact on incentives. Our theoretical results, however, pertain to SRA as it is formally defined.

#### 4.1. The main ideas behind SRA

SRA is a *biased* output agreement mechanism, i.e., an agent gets paid for her evaluation only if her answer matches the answer reported by her peer agent who has made the same evaluation, where the payment depends on the answer. The simplest description of the core idea of the mechanism is obtained in the *hypothetical* scenario where the generating model  $(P_X, \mathbf{p})$  is known to the principal and is commonly known to the agents (this is the setting considered by Miller et al. (2005)). In this setting, a biased output agreement scheme is defined as follows.

1. Each agent  $j$  is paired with another randomly chosen peer agent  $j'$ , and their responses are compared.
2. If their responses don’t match, then  $j$  gets no reward.
3. If their responses match and this common response is  $y \in \mathcal{Y}$ , then  $j$  gets a positive reward  $e(y)$ .

The agreement rewards  $e(y)$  for  $y \in \mathcal{Y}$  are defined as a function of the generating model of responses. The main innovation in SRA is how these agreement rewards are defined. To motivate their design, we first discuss why the naïve output agreement mechanism fails to incentivize truthful behavior.

**Failure of naïve output agreement.** In the naïve output agreement mechanism,  $e(y)$  is defined to be a constant  $K > 0$  for all  $y \in \mathcal{Y}$ . This mechanism tries to exploit an intuitive property that one may naïvely expect to hold in many scenarios, which is that the peer agent  $j'$  has the highest conditional likelihood of observing the *same* answer as that observed by an agent  $j$ . However, this property is not always true. For instance, it is violated if, irrespective of the observation made by

an agent, one “popular” answer has an overwhelming conditional likelihood of being observed by the peer agent. This feature can be observed in Example 1.

EXAMPLE 3. Consider the setting in Example 1. The question is “Did the plumber show up within 5 minutes of his/her appointed time?” with the possible observations/answers being  $\mathcal{Y} = \{\text{Yes}, \text{No}\}$ . Suppose an agent observes that the plumber did not arrive within 5 minutes of her appointed time, i.e., her observation was “No.” Then the conditional probability of the peer agent, assumed to be truthful, replying “Yes” can be computed to be  $0.5409$ , which is higher than the conditional probability of her replying “No,” which can be computed to be  $0.45909$ . Hence, replying “Yes,” i.e., lying, results in a higher expected payoff than being truthful and replying “No.”

In the example above, plumbers are a priori overwhelmingly likely to turn up on time, to the extent that even if an agent observes that a plumber was late, she will still find it more likely that the same plumber will be observed to be on time by her peer agent. Thus, assuming that the peer agent is truthful, it is better to lie and say that the plumber was on time. Summarily, truthful behavior is not an equilibrium in this case.

To overcome this shortcoming of the naïve output agreement scheme, the key obstacle that one must tackle is this tendency of regressing to the conditionally most popular answer. Formally, if the observation of an agent  $j$  is  $Y_j = y$  for some  $y \in \mathcal{Y}$ , her tendency is to report  $\arg\max_{y' \in \mathcal{Y}} P(Y_{j'} = y' | Y_j = y)$  so as to maximize the probability of agreement. As we saw in the example above, this optimal answer need not necessarily be  $y$ , i.e., the inequality

$$P(Y_{j'} = y | Y_j = y) \geq P(Y_{j'} = y' | Y_j = y), \quad (4)$$

doesn’t necessarily hold for all  $y, y' \in \mathcal{Y}$ , even in the binary responses setting where  $|\mathcal{Y}| = 2$ . Biased output agreement schemes can tackle this issue by scaling the rewards for agreement depending on the answer, essentially lowering rewards for answers that are expected to be more (conditionally) likely and increasing rewards for answers that are less (conditionally) likely. The intuition is that if an agent observes an answer that is less likely to also be observed by her peer, she is still incentivized to report that answer since the matching reward on that answer is higher. Conversely, if she observes an answer that is more likely to also be observed by her peer, she prefers to report that answer despite the low reward for agreement relative to other answers since the probability of agreement is higher. The challenge is to design the reward scalings for different answers so that these incentives for truthful reporting for an agent are satisfied for *each* answer (assuming every other agent is truthful).

A straightforward way of making an agent indifferent between different reports (and hence weakly incentivize truthful behavior) is by defining the agreement reward for answer  $y' \in \mathcal{Y}$  to be proportional to  $1/P(Y_{j'} = y' | Y_j = y)$ , where  $y$  is the answer observed by the agent. This approach is

rendered infeasible by observing that  $y$  is not known to the mechanism – indeed, the entire exercise is meant for the purpose of eliciting  $y$ . With this observation in perspective, we discuss two approaches of defining the rewards before we present our approach in SRA.

**Approach 1: Peer Truth Serum.** The PTS mechanism (Jurca and Faltings 2011, Faltings et al. 2017) scales the agreement reward for answer  $y' \in \mathcal{Y}$  by  $1/P(Y_{j'} = y')$ , i.e., reward for agreement on an answer is inversely proportional to the probability that an evaluating agent makes that observation. We thus obtain truthful behavior from agent  $j$  under PTS if for all  $y \in \mathcal{Y}$  and  $y' \neq y$ ,

$$\frac{P(Y_{j'} = y | Y_j = y)}{P(Y_{j'} = y)} \geq \frac{P(Y_{j'} = y' | Y_j = y)}{P(Y_{j'} = y')}, \text{ i.e., if,} \quad (5)$$

$$P(Y_{j'} = y | Y_j = y) \geq P(Y_{j'} = y | Y_j = y'). \quad (6)$$

In the literature, this is referred to as the ‘self-predicting responses’ condition (that we discussed earlier in Section 2), which is not satisfied in general for homogeneous responses, except when  $|\mathcal{Y}| = 2$ . Intuitively, defining the popularity index of an answer to be equal to the *marginal* probability of a single agent making that observation does not capture the fact that the agent evaluates the *conditional* probabilities of agreement for the different answers, which depend on the *joint probabilities* of a pair of agents making various observations.

**Approach 2: A “conditional” peer truth serum (CPTS).** With the goal of incorporating the conditional distribution of the observations in the definition of the agreement rewards, another proxy for the ideal scaling can be defined to be  $1/P(Y_{j'} = y' | Y_j = y')$ , i.e.,  $e(y')$  is defined to be proportional to  $1/P(Y_{j'} = y' | Y_j = y')$ . In this case, we obtain truthful behavior from agent  $j$  if for all  $y \in \mathcal{Y}$  and  $y' \neq y$ ,

$$\frac{P(Y_{j'} = y | Y_j = y)}{P(Y_{j'} = y | Y_j = y)} = 1 \geq \frac{P(Y_{j'} = y' | Y_j = y)}{P(Y_{j'} = y' | Y_j = y')}, \text{ i.e., if,} \quad (7)$$

$$P(Y_{j'} = y' | Y_j = y') \geq P(Y_{j'} = y' | Y_j = y). \quad (8)$$

This condition is exactly the self-prediction condition from Equation 6. Thus, the conditions that are necessary for the PTS and CPTS mechanisms to induce truthful behavior are the same. This shows that a naïve incorporation of conditional probabilities in the reward scalings need not provide any advantage over PTS, at least in terms of the conditions required for truthfulness.

**SRA’s approach.** SRA defines the rewards for agreement as  $e(y') = K/\sqrt{P(Y_j = Y_{j'} = y')}$  for each  $y' \in \mathcal{Y}$ , for some  $K > 0$ . To put it simply, *the reward for an agreement is inversely proportional to the square root of the probability of that agreement*. Thus a more probable agreement earns a

lower reward than an agreement that is relatively less probable. To see the relation to PTS and CPTS, observe that SRA simply replaces the scaling  $1/P(Y_{j'} = y')$  of PTS and the scaling  $1/P(Y_{j'} = y' | Y_j = y')$  of CPTS by the product of the square root of the two scalings  $1/\sqrt{P(Y_{j'} = y')} \times 1/\sqrt{P(Y_{j'} = y' | Y_j = y')} = 1/\sqrt{P(Y_j = Y_{j'} = y')}$ . By making this change, SRA explicitly incorporates the joint probabilities of agreements in the rewards.

It turns out that the scaling under SRA successfully overcomes the shortcoming of the naïve output agreement scheme in homogeneous responses settings. In particular, truthful behavior is a Bayes-Nash equilibrium in this mechanism. To see this, 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,

$$K \frac{P(Y_{j'} = y | Y_j = y)}{\sqrt{P(Y_j = Y_{j'} = y)}} = K \frac{\sqrt{P(Y_{j'} = Y_j = y)}}{P(Y_j = y)}. \quad (9)$$

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

$$K \frac{P(Y_{j'} = y' | Y_j = y)}{\sqrt{P(Y_j = Y_{j'} = y')}} = K \frac{P(Y_{j'} = y', Y_j = y)}{P(Y_j = y) \sqrt{P(Y_j = Y_{j'} = y')}}. \quad (10)$$

Thus being truthful gives a higher reward if the quantity in Equation 9 is larger than the quantity in Equation 10, which is, if,

$$\sqrt{P(Y_{j'} = Y_j = y)} \sqrt{P(Y_j = Y_{j'} = y')} \geq P(Y_{j'} = y', Y_j = y). \quad (11)$$

This inequality resembles the well-known Cauchy-Schwarz inequality that relates second-order moments of two random variables  $X$  and  $Y$  as,  $E[XY] \leq \sqrt{E(X^2)E(Y^2)}$ . Hence, if the joint distribution of responses of the two agents satisfies this inequality for each  $y, y' \in \mathcal{Y}$ , we say that the distribution satisfies the ‘‘Cauchy-Schwarz’’ (CS) property. And in this case, truthful equilibrium is a Bayes-Nash equilibrium under SRA.

Now the key interesting fact is that this property is *always* satisfied for homogeneous responses, i.e., for objective evaluations. To see this, the inequality in Equation 11 can be expressed as follows in the homogeneous responses setting.

$$\sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_y(x)^2} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y'}(x)^2} \geq \sum_{x \in \mathcal{X}} P_X(x) p_y(x) p_{y'}(x). \quad (12)$$

But this is precisely the Cauchy-Schwarz inequality that always holds. Hence, truthful behavior is incentivized as a Bayes-Nash equilibrium under SRA. For truthful behavior to be a *strict* Bayes-Nash equilibrium, we need the above inequality to be strict. We show that this requirement is satisfied if the generating model is ‘separated’: a condition we discuss in Section 4.2.

Now implementing the mechanism above in our setting is infeasible since the principal does not know the generating model  $(P_X, \mathbf{p})$ . Moreover, the generating model is not assumed to be commonly

known to the agents. SRA addresses these issues by replacing the required agreement probabilities with consistent statistical estimates computed from reports obtained across multiple tasks. Observe that if everyone except agent  $j$  is truthful, then in the definition of SRA,  $E[\bar{f}_j(y)] = P(Y_{j_1(i)}^i = Y_{j_2(i)}^i = y) + 1/(N - |\mathcal{N}_j|)$ . 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, i.e.,  $\bar{f}_j(y)$  is an asymptotically consistent estimate of  $P(Y_{j_1(i)}^i = Y_{j_2(i)}^i = y)$ . For a large enough  $N$ , we can show that the estimate's quality is sufficiently high to ensure that truthfulness is recovered as a strict Bayes-Nash equilibrium for any 'separated' generating model (see our main result in Section 4.3). Thus, the agents only need to commonly know the generating model's structure and that it is separated to obtain truthful behavior.

## 4.2. Obtaining strictness

An important goal for any reward mechanism is to *strictly* incentivize truthfulness, i.e., in the truthful equilibrium, each agent gets a strictly higher reward by being truthful than by adopting any other strategy. Without this property, trivial mechanisms like the one that gives a fixed payment to each agent regardless of her report, in principle, weakly incentivize truthfulness. For truthfulness to be a strict equilibrium under SRA, we need the Cauchy-Schwarz inequality in Equation 12 to be strict for every pair  $y, y' \in \mathcal{Y}$ . It will be useful to define the following notion of the ‘‘inequality gap.’’

DEFINITION 4 (CAUCHY-SCHWARZ INEQUALITY GAP). 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. Consider the following definition, which will be useful for our forthcoming discussion.

DEFINITION 5 (SEPARATION). We say that a generating model  $(P_X, \mathbf{p})$  is separated if  $\delta(P_X, \mathbf{p}) > 0$ . We say that it is  $\alpha$ -separated for any  $\alpha > 0$  if  $\delta(P_X, \mathbf{p}) \geq \alpha$ .

To understand whether separation is a reasonable assumption on the generating model, a little demystification of this condition is in order. For any answer  $y \in \mathcal{Y}$ , define the vector

$$\mathbf{v}(y) \triangleq (\sqrt{P_X(x)} p_y(x); x \in \mathcal{X}). \quad (13)$$

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

$$|\mathbf{v}(y) \cdot \mathbf{v}(y')| \leq \|\mathbf{v}(y)\| \|\mathbf{v}(y')\|. \quad (14)$$

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

$$\theta(\mathbf{u}, \mathbf{v}) \triangleq \arccos \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}, \quad (15)$$

when both  $\mathbf{u}$  and  $\mathbf{v}$  are non-zero and as 0 when either of them is a zero vector. We can then show that the inequality in Equation 14 is strict if and only if the angle between the vectors  $\mathbf{v}(y)$  and  $\mathbf{v}(y')$  is positive. The following proposition gives a precise statement.

**Proposition 4.1** *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  $(P_X, \mathbf{p})$  is  $\alpha$ -separated.*
2. *There is some  $\gamma > 0$  such that  $\theta(\mathbf{v}(y), \mathbf{v}(y')) \geq \gamma$  for all  $y, y' \in \mathcal{Y}$  such that  $y \neq y'$ .*

Thus, separation is equivalent to assuming that the angle between  $\mathbf{v}(y)$  and  $\mathbf{v}(y')$  is positive 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  $p_y(x) = Cp_{y'}(x)$  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, separation is also equivalent to the stochastic relevance condition that is imposed to obtain strictness in several works in this domain, starting from Miller et al. (2005). 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  $\theta(\mathbf{v}(y), \mathbf{v}(y')) > 0$ , which is equivalent to separation by Proposition 4.1. Showing that separation implies stochastic relevance is less straightforward and we show it in Proposition E.2 in the Appendix.

### 4.3. Main result

The following result presents the main incentive property satisfied by SRA.

**THEOREM 1.** *Consider an  $\alpha$ -separated generating model  $(P_X, \mathbf{p})$  that is commonly known to the agents. Further, suppose that  $N_j \leq n$  for all  $j \in \mathcal{M}$  and  $|\mathcal{M}_i| \geq 2$  for all  $i \in \mathcal{N}$ . 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*

- SRA 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.

Note that the bounds we derive in the proof of this result can be used to explicitly calculate an  $N_0$  as a function of  $\omega$ ,  $\alpha$ ,  $|\mathcal{Y}|$ ,  $n$ , and  $K$ . Also note that, although the theorem assumes that the generating model is commonly known to the agents, the mentioned properties hold for any  $\alpha$ -separated generating model. In particular, the dependence of  $N_0$  on the generating model is only through  $\alpha$ . Hence, this result implies that even if only the fact that  $(P_X, \mathbf{p})$  is  $\alpha$ -separated is common knowledge amongst agents, then irrespective of their individual beliefs about the specifics of the generating model, truthful reporting is strictly incentivized in the game induced by the mechanism for a large enough  $N$ .

The second claim in the theorem is crucial too: it says that in the truthful 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. Such strict incentives allow the principal to account for any costs that the agents may incur for their evaluation effort by appropriately scaling the mechanism's rewards. Note that it is not possible to ensure that the difference between the expected payoff under truthful behavior and that under *any* other strategy is bounded away from zero, since one can choose a randomized strategy that chooses a non-truthful response with an arbitrarily small probability. However, our goal here is to deter agents who report an arbitrary answer without investing effort into making an observation. Our result ensures that any perceived cost for such effort can be absorbed in the difference in the payoff under truthful reporting and under any reporting strategy that ignores the observation, by appropriately scaling the rewards.<sup>8</sup>

Finally, we note that although the common knowledge assumption above is necessary to obtain the theoretical properties of our mechanism, our numerical experiments in Section 6 test SRA under more practical considerations.

## 5. Equilibrium selection

In this section, we address the issue of multiplicity of equilibria in the game induced by SRA. First, observe that if truthful behavior is an equilibrium, then so is any *symmetric fully informative strategy profile* in which 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

<sup>8</sup> Liu and Chen (2016) show how to learn such a scaling when there is heterogeneity in the cost for effort in classical output agreement mechanisms. We believe a similar approach can be adopted in our case.

for the agents to play a fully informative equilibrium other than truthful behavior 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 a weak dual to this result: under a 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 define.

### 5.1. Truthfulness vs. symmetric equilibria as $N \rightarrow \infty$

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 9),

$$\sum_{y \in \mathcal{Y}} P(Y_j^i = y) K \frac{\sqrt{P(Y_j^i = Y_j^i = y)}}{P(Y_j^i = y)} = K \sum_{y \in \mathcal{Y}} \sqrt{P(Y_j^i = Y_j^i = y)} = K \sum_{y \in \mathcal{Y}} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_y(x)^2}. \quad (16)$$

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

$$K \sum_{y \in \mathcal{Y}} \sqrt{P(Z_j^i = Z_j^i = y)} = K \sum_{y \in \mathcal{Y}} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) \left( \sum_{y' \in \mathcal{Y}} p_{y'}(x) q_y(y') \right)^2}. \quad (17)$$

It turns out that the quantity in Equation 17 is, in general, lower than the quantity in Equation 16. How much lower depends on the ‘uninformativeness’ of the strategy  $\mathbf{q}$ : the more uninformative the strategy  $\mathbf{q}$ , the higher is the difference. We will describe this phenomenon more generally since we believe it has applications beyond this work (see Appendix B.1 for a discussion). We first define the following notion of a *square root agreement measure*.

**DEFINITION 6 (SQUARE ROOT AGREEMENT MEASURE (SRAM)).** 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 square root agreement measure between  $Y_1$  and  $Y_2$  is defined as

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

From the definition of SRAM, it is clear that under any 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 SRAM between the reported responses (see Equation 17). Some important properties of the SRAM are presented in Appendix B.

Next, we define a new notion of *uninformativeness* of a reporting strategy. 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 7 (AN UNINFORMATIVENESS MEASURE).** The uninformativeness of a reporting strategy  $\mathbf{q}$  is defined as

$$\Omega(\mathbf{q}) = \frac{1}{|\mathcal{Y}|(|\mathcal{Y}| - 1)} \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}; y' \neq y''} \sqrt{q_y(y')q_y(y'')}. \quad (18)$$

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

Certain important properties of the uninformativeness measure are presented in Appendix C. In particular,  $\Omega(\mathbf{q}) = 0$  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, and  $\Omega(\mathbf{q})$  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 independently of the true answer.

We finally present the following information monotonicity property, which is key to our results.

**Proposition 5.1 (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,*

$$\Gamma(Z_1, Z_2) \leq \Gamma(Y_1, Y_2) - \frac{\delta(P_X, \mathbf{p})\Omega(\mathbf{q})^2(|\mathcal{Y}| - 1)}{2\sqrt{|\mathcal{Y}|}}. \quad (19)$$

To see how this property helps us, recall from Equations 16 and 17 that the expected payoff under any symmetric fully informative strategy profile is  $K\Gamma(Y_1, Y_2)$ , and that under any other symmetric strategy profile  $\mathbf{q}$  is  $K\Gamma(Z_1, Z_2)$ , where  $Z_1$  and  $Z_2$  is obtained by applying a reporting strategy  $\mathbf{q}$  independently to  $Y_1$  and  $Y_2$ . 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 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. In other words, SRA is asymptotically strongly truthful across symmetric equilibria. In the next section, we use Proposition 5.1 to address the case where  $N$  is large but finite.

## 5.2. Equilibrium selection in the finite $N$ regime

Now we turn to the finite  $N$  setting. In this case, the expected payoffs under the fully informative strategy and 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 compared to any other symmetric equilibrium.

**THEOREM 2.** *Consider an  $\alpha$ -separated generating model  $(P_X, \mathbf{p})$  that is commonly known to the agents. Further, suppose that  $\mathcal{N}_j \leq n$  for all  $j \in \mathcal{M}$  and  $|\mathcal{M}_i| \geq 2$  for all  $i \in \mathcal{N}$ . Then for any  $\omega > 0$ , there is a positive integer  $N_0$  that depends on  $\alpha, \omega, n, K$ , and  $|\mathcal{Y}|$ , such that for any  $N > N_0$ , under SRA,*

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

Note once again that the bounds we derive in the proof of this result can be used to explicitly calculate an  $N_0$  as a function of  $\alpha, \omega, n, K$ , and  $|\mathcal{Y}|$ . Also, note that similar to Theorem 1, although Theorem 2 assumes that the generating model is commonly known to the agents, the result holds for any  $\alpha$ -separated generating model (since the dependence of  $N_0$  on  $(P_X, \mathbf{p})$  is only through  $\alpha$ ). Hence, we conclude that these properties presented in Theorem 2 hold even in the game where it is only commonly known to the agents that the generating model is  $\alpha$ -separated.

Finally, note that in the statement of Theorem 2,  $N > N_0$  suffices to ensure that the truthful equilibrium gives a payoff at most  $\omega$  lower than that under *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.

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 5.1 should result in an upper bound on the uninformativeness 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, where the symmetric strategy profile 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 3.** *Consider an  $\alpha$ -separated generating model  $(P_X, \mathbf{p})$ . Further, suppose that  $N_j \leq n$  for all  $j \in \mathcal{M}$  and  $|\mathcal{M}_i| \geq 2$  for all  $i \in \mathcal{N}$ . Then for any  $\omega > 0$  and  $\eta > 0$ , there is a positive integer  $N_0$  that depends on  $\omega, \alpha, \eta, n, K$ , and  $|\mathcal{Y}|$ , such that for any  $N > N_0$ , under SRA, any symmetric strategy profile in which the probability of reporting any answer  $y \in \mathcal{Y}$  is either 0 or at least  $\eta$ , and that gives a higher expected payoff to an agent than the truthful strategy profile, is at most  $\omega$ -uninformative.*

**REMARK 2.** Theorem 3 implies that for a large enough  $N$ , the truthful equilibrium yields a strictly higher payoff to each agent than any equilibrium where all agents report a fixed answer irrespective of the observation. To see this, note that the latter strategy profile is 1-uninformative (see properties of the uninformativeness measure presented in Appendix C); hence, we can choose  $\omega \in (0, 1)$ . Moreover,  $\eta$  can be chosen to be any number in  $(0, 1)$  since mapping all observations to a single response implies that the probability of any response is either 0 or 1. Choosing such  $\eta$  and  $\omega$ , we see that the above property will hold for any  $N$  larger than the corresponding  $N_0$ . This property is in contrast to the multi-task extension of the peer-prediction method (Miller et al. 2005) and the mechanism of Radanovic and Faltings (2015), in which each agent reporting a fixed answer is an equilibrium that yields the highest possible payoff. A similar argument shows that for a large enough  $N$ , the truthful equilibrium yields a strictly higher payoff to each agent than any equilibrium where all agents map smaller sets of responses to a single response, e.g., if the responses are  $\{a, b, c, d\}$  then  $\{a, b\}$  map to  $b$  and  $\{c, d\}$  map to  $d$ . This contrasts with the informed truthfulness property of CA, under which the truthful equilibrium doesn't necessarily strictly payoff-dominate equilibria of this type (particularly in scenarios where the observations are clustered; see Definition 10 in the Appendix).

### 5.3. Strong truthfulness vs. Strong truthfulness across symmetric equilibria

In this section, we show that the property of strong truthfulness across symmetric equilibria of SRA can be strengthened to strong truthfulness (across all equilibria) under certain assumptions commonly made in the peer-prediction literature. To do so, we first discuss what kind of *asymmetric* equilibria may arise under SRA that yield a higher reward to the agents than the truthful equilibrium.

**EXAMPLE 4.** *Consider the labor platform setting considered in Example 1. In this example, the observation that is expected to be least frequently agreed upon is “No” (i.e., the plumber did not arrive within 5 mins of his/her appointed time). Let John be a plumber operating on this platform, and*

consider the strategy profile where the customers who evaluate John always report “No” for him, while all customers report truthfully for every other plumber. This is an asymmetric strategy profile: the agents who evaluate John follow a different reporting strategy from those who don’t. In the limit where the number of plumbers on the platform is large, under the assumptions of Theorem 1, we claim that this strategy profile constitutes an equilibrium. The key idea is that the popularity indices are not expected to be impacted by John’s reports in this limit. Hence, assuming everyone else adheres to this strategy profile, being truthful is optimal on all tasks other than evaluating John. Moreover, for John’s evaluation, reporting “No” is optimal since it is the only way to obtain a positive payment. Thus this strategy profile constitutes an equilibrium. Moreover, this equilibrium gives a strictly higher reward to agents who have evaluated John than the truthful equilibrium, since the answer “No” is expected to have the lowest popularity index and hence the highest reward for agreement.

The ability to construct such equilibria relies on agents coordinating their behavior on a task or a small subset of tasks. Our model assumes that task-identifying information is available and the task allocation is exogenously specified, so we cannot preclude this possibility. However, when tasks are randomly allocated to the agents and agents do not make their reporting strategy contingent on the task identity, such equilibria are not expected to arise. Formally, consider the following assumptions.

**ASSUMPTION 1 (Randomized Task Allocation).** *Suppose that there are  $N$  tasks and  $M$  agents. Each task is to be performed by  $m$  agents, where we assume that  $m \geq 2$ . Also, suppose that no agent should perform more than  $n$  tasks on average. To ensure that this is feasible, we consider a system regime in which  $N$  and  $M$  simultaneously grow such that  $M > mN/n$ . For each task,  $m$  distinct agents are uniformly sampled and assigned to that task. Note that each agent gets picked for a task with probability  $m/M$  and thus performs  $mN/M$  tasks on average, which is less than  $n$ .*

Note that the average number of tasks performed by an agent,  $n$ , can be chosen to be as small as required (at the cost of requiring a large  $M$ ) to ensure that each agent doesn’t perform more than one evaluation with a high probability.

**ASSUMPTION 2 (Task-independent strategies).** *Assume that each agent  $j$  picks a reporting strategy  $\mathbf{q}^j$  before the allocation of evaluation tasks, with the assumption that the same reporting strategy will be applied to every task that the agent performs.*

Such assumptions are often justifiable in crowdsourcing applications and hence are commonly made in the multi-task peer-prediction literature as we discussed in Section 2. We can argue that under such assumptions, the fact that SRA is asymptotically strongly truthful across symmetric equilibria implies that it is, in fact, asymptotically strongly truthful. To state our result, we need the notion of the population average reporting strategy given a strategy profile, defined as  $\bar{\mathbf{q}}(y) = \frac{1}{M} \sum_{j \in \mathcal{M}} \mathbf{q}^j(y)$  for all  $y \in \mathcal{Y}$ . We have the following result.

THEOREM 4. Consider an  $\alpha$ -separated generating model  $(P_X, \mathbf{p})$  that is commonly known to the agents. Suppose that Assumptions 1 and 2 are satisfied. Then for any  $\omega > 0$  and  $\eta > 0$ , there is a positive integer  $N_0$  that depends on  $\alpha, \eta, \omega, m, n, K$ , and  $|\mathcal{Y}|$ , such that for any  $N > N_0$ , under SRA, the following holds.

1. Any symmetric fully informative strategy profile is a strict Bayes-Nash equilibrium.
2. For any Bayes-Nash equilibrium strategy profile such that probability of reporting any answer  $y \in \mathcal{Y}$  is either 0 or at least  $\eta$  under the population average reporting strategy, the expected payoff of any agent is at most  $\omega$  higher than that under any symmetric fully informative strategy profile.
3. Consider any Bayes-Nash equilibrium strategy profile such that probability of reporting any answer  $y \in \mathcal{Y}$  is either 0 or at least  $\eta$  under the population average reporting strategy, and that yields a weakly higher expected payoff to any agent than that under any symmetric fully informative strategy profile. Then the population average reporting strategy under this strategy profile is at most  $\omega$ -uninformative.

Note that for any  $\omega > 0$ , we are not able to achieve  $\omega$ -domination of the fully informative equilibrium over *all* equilibria for some finite but large enough  $N$ , but only those equilibria in which the population average reporting probabilities are either 0 or bounded away from 0. This is because if the average reporting probability for an answer across the population becomes arbitrarily small in the number of tasks  $N$  as  $N$  grows, then it is not possible to obtain decaying concentration bounds on the agreement rewards; refer to the similar discussion in the context of Theorem 3 above. However, for any fixed  $\omega > 0$ , in the limit as  $N \rightarrow \infty$ , *all* equilibria are  $\omega$ -dominated by the truthful equilibrium. Moreover, for any fixed  $\omega > 0$ , in the limit as  $N \rightarrow \infty$ , *any* equilibrium that yields a weakly higher payoff than the truthful equilibrium is at most  $\omega$ -uninformative. Since  $\omega$  is arbitrary, this shows that the mechanism is *asymptotically* strongly truthful.<sup>9</sup>

Informally, the idea of the proof of Theorem 4 is the following. The random allocation of tasks, as well as non-contingency of reporting strategies on task identities, imply that from the perspective of each agent, no task is special, i.e., they are expected to be paired with a generic peer from the population of agents irrespective of the identity of the task. The assumption that the number of tasks performed by each agent is bounded on average implies that one agent's reports are expected to have a vanishing impact on the population indices. Because of these assumptions, from the perspective of *each* agent, the remainder of the population with an asymmetric strategy profile can be replaced by a population where each agent utilizes the *population average* reporting strategy  $\bar{\mathbf{q}}$ ,

<sup>9</sup> We remark that agents are expected to use simple strategies in practice, and thus strategies in which the reporting probabilities of the answers become arbitrarily close to zero are unlikely to arise in finite but large  $N$  settings under an  $\alpha$ -separated generating model.

while (approximately) preserving the payoff structure of the game. Additionally, the fact that *each* of the different strategies utilized by the agents is near-optimal against this *average* reporting strategy of the population implies that the average strategy  $\bar{\mathbf{q}}$  is *also* near-optimal against  $\bar{\mathbf{q}}$ , because the payoff of an agent is linear in her reporting strategy. Hence, the agents' payoff is approximately the same as one in the strategy profile where everyone follows  $\bar{\mathbf{q}}$ , which constitutes a symmetric strategy profile (with the approximation error in each of the above arguments vanishing in the large tasks limit). The result then essentially follows from the fact that SRA is asymptotically strongly truthful across symmetric equilibria.

We argue that Assumptions 1 and 2 are generally not appropriate in the context of feedback elicitation on online platforms. In this context, tasks are not randomly allocated but are chosen by the agents themselves. Also, entities being evaluated have a multitude of extraneous identifying features, e.g., the name of a plumber, the race or ethnicity of the Airbnb property owner, etc. Given this information, one cannot preclude the possibility of the agents choosing their reporting strategy based on such extraneous features and achieving coordination in their reports, potentially achieving higher payoffs compared to truthful reporting of their observations.

However, such extraneous coordination could be more or less likely depending on the context. For instance, in the example above, the possibility that agents coordinate their behavior with respect to this one plumber, John, out of potentially hundreds seems unlikely. However, such coordination may not be unreasonable in reviewing a popular neighborhood restaurant amongst people residing in that neighborhood.

To ensure that practitioners concerned with our target applications are not misled by SRA's properties reported in our work, we choose not to make restricting assumptions that disallow such coordination and content ourselves with the weaker property of strong truthfulness *across symmetric equilibria*. As we argued in Section 2, in the absence of such assumptions, the dominance properties of other mechanisms are also restricted.

## 6. Numerical Evaluation

In this section, we numerically evaluate the performance of SRA, both, in the case of objective evaluations using synthetic data (Section 6.1) and in the case of subjective evaluations using real data from online platforms (Section 6.2). In the latter case, we also compare SRA's performance to other related mechanisms.

### 6.1. Performance on objective evaluations

In this section, we test SRA's robustness in inducing truthful behavior in the finite  $N$  regime, in settings where population homogeneity may not necessarily hold exactly for objective evaluations due to observation or reporting biases. To do so, we assume the perspective of a single agent

operating in a platform environment where SRA is deployed, and examine her incentives for truthful reporting given her beliefs about the other agents' generating model of the observations and their reporting behavior. This exercise also serves as a practical alternative to arguing about whether the common knowledge assumptions that are required for our results hold in practice; we demonstrate numerically that each agent has strong practical incentives to be truthful under SRA if she believes that agents are largely unbiased and they report their observations faithfully. We begin by defining our experimental setup and the performance measures that we consider.

**Setting.** We consider the setting of an online service platform such as Thumbtack,<sup>10</sup> on which a large number of moving companies offer local moving services. The platform seeks to collect information about how punctual the different moving companies are in adhering to the committed time to start the move. Poor scheduling and information collection practices can result in large variability in the start and finish times of different moves in a day, potentially leading to disgruntled customers.

In our simulation setup, we assume that a moving company's delay in showing up is an exponentially distributed random variable with a certain mean. We assume that there are  $|\mathcal{X}| = 5$  types of moving companies, where

$$\mathcal{X} = \{1 \text{ (Mostly Punctual)}, 2 \text{ (Somewhat Punctual)}, 3 \text{ (Mostly Tardy)}, 4 \text{ (Tardy)}, 5 \text{ (Very Tardy)}\}.$$

Each of these types is associated with a mean for the distribution of delay. We denote these means as  $\boldsymbol{\mu} = \{\mu_1, \mu_2, \mu_3, \mu_4, \mu_5\}$ , where we assume that  $\mu_1 < \mu_2 < \mu_3 < \mu_4 < \mu_5$ . The distribution of these five types across the population is denoted as  $P_X$ . The platform asks the customers the following question: "How long (in minutes) after the scheduled time of appointment did the movers show up?" We consider two settings that differ in the possible answers to this question.

1. The number of answers is  $|\mathcal{Y}| = 3$  where,  $\mathcal{Y} = \{0 \text{ to } 30, 30 \text{ to } 60, 60 \text{ to } \infty\}$ .
2. The number of answers is  $|\mathcal{Y}| = 5$  where,  $\mathcal{Y} = \{0 \text{ to } 15, 15 \text{ to } 30, 30 \text{ to } 45, 45 \text{ to } 60, 60 \text{ to } \infty\}$ .

Given the type of a mover  $x \in \mathcal{X}$ , the delay is exponentially distributed with mean  $\mu_x$ , and hence the conditional probability of making an observation  $y = l$  to  $u$  is given by:

$$p_y(x) = \exp\left(-\frac{l}{\mu_x}\right) - \exp\left(-\frac{u}{\mu_x}\right). \quad (20)$$

**Instances.** An instance is specified by the probability distribution of the types  $P_X$ , and the mean delays associated with the types  $\boldsymbol{\mu}$ . We generate 10000 instances. In each instance,  $P_X$  is determined by sampling 5 numbers independently and uniformly in  $[0, 1]$  and dividing them by their total to

<sup>10</sup> Thumbtack is an online service platform that matches customers with local professionals; see <https://www.thumbtack.com/>

obtain a probability distribution.  $\boldsymbol{\mu}$  is obtained by sampling 5 numbers independently and uniformly in  $[0, 60]$  (hence the mean delay of a moving company can be at most 60 min) and sorting them in an increasing order to satisfy the requirement that that  $\mu_1 < \mu_2 < \mu_3 < \mu_4 < \mu_5$ . Once  $P_X$  and  $\boldsymbol{\mu}$  are thus specified, the conditional distributions over observations,  $\mathbf{p}(x)$  for each  $x \in \mathcal{X}$ , get specified as well according to Equation 20, for both the settings of  $|\mathcal{Y}| = 3$  and  $|\mathcal{Y}| = 5$ .

REMARK 3. We note that for  $|\mathcal{Y}| = 3$ , 9787 out of the 10000 instances thus generated failed to satisfy the self-predicting responses condition (see item 2 in Section 2.1) and 9995 instances failed to satisfy the categorical responses condition (see item 1 in Section 2.1). For  $|\mathcal{Y}| = 5$ , none of the 10000 instances satisfied either of the two conditions. Hence, the PTSC mechanism (Radanovic et al. 2016) and the mechanism of Dasgupta and Ghosh (2013) are inapplicable with a high frequency in this setting. This finding continues to hold if  $P_X$  is chosen to be more “regular.” We sampled 10000 instances in two additional settings when  $|\mathcal{Y}| = 5$ : (a) one where  $P_X(x)$  is decreasing in  $x$ , i.e., types with higher mean delay are more rare, and (b)  $P_X(x)$  is increasing in  $x$ , i.e., types with higher mean delay are more frequent. In both these settings, none of the instances satisfied the self-predicting responses or categorical responses conditions.

REMARK 4. The CA-HR mechanism described in Section D.1 in the Appendix is informed truthful across symmetric equilibria in general. However, if the instance is such that the joint distribution of observations of a pair of agents for a common evaluation task is not “clustered,” as defined in Definition 10 in the Appendix, then the mechanism is strongly truthful across symmetric equilibria for that instance. We find that for  $|\mathcal{Y}| = 3$ , 9995 out of the 10000 instances thus generated had clustered observations. For  $|\mathcal{Y}| = 5$ , 9940 out of the 10000 instances had clustered observations. Thus, in this setting, with a very high frequency, SRA is the only applicable mechanism requiring one task per agent that is also strongly truthful across symmetric equilibria.

**Agent beliefs.** We focus on a single agent (a customer on the platform), whom we refer to as agent  $j$ , and examine her incentives for being truthful under SRA. Each instance that we define above represents a belief that agent  $j$  has about the generating model for the movers’ true delays. Additionally, we allow the agent to account for potential biases in the observation-making process of other agents. In particular, agent  $j$  believes that, while she can make perfect observations, other agents do not observe the delay perfectly, but rather observe a biased version. Owing to this, she believes that given a mover type  $x$ , the conditional distribution of observations made by a generic agent in the population is not  $\mathbf{p}(x)$ , but  $\mathbf{p}'(x)$ , which is (slightly) different. We assume that

$$\mathbf{p}'(x) = (1 - \epsilon)\mathbf{p}(x) + \epsilon\mathbf{q}'(x), \quad (21)$$

where  $\mathbf{q}'$  represents agent  $j$ ’s belief of the bias in the population, and  $\epsilon$  represents the magnitude of the bias. In each of the 10000 instances, for both the settings of  $|\mathcal{Y}| = 3$  and  $|\mathcal{Y}| = 5$  answers, we

independently sample an associated bias  $\mathbf{q}'(x)$  for each  $x \in \mathcal{X}$ , again by generating 3 and 5 values uniformly in  $[0, 1]$  and dividing each by their sum to obtain a distribution. In our analysis, we will independently consider different values of  $\epsilon \in \{0, 0.1, 0.2\}$ , where  $\epsilon = 0$  is the case where agent  $j$  believes that everyone else makes perfect observations.

**Performance measures.** We define two performance measures that capture the relative attractiveness of non-truthful behavior compared to being truthful from the perspective of agent  $j$ , assuming every other agent truthfully reports her (potentially biased) observation. Before we define these measures, we first compute the expected reward for agreement on each answer  $y \in \mathcal{Y}$  under the mechanism from the perspective of agent  $j$ , i.e., compute  $E(e_j(y))$ . For each instance, assuming  $K = 1$  and denoting  $N - |\mathcal{N}_j| = N'$ , we denote  $r(y, N', \epsilon) \triangleq E(e_j(y))$  to be the expected payment that agent  $j$  receives if she and her peer both give a matching response  $y$ , for each  $y \in \mathcal{Y}$ . We consider values of  $N'$  in the set  $\{200, 400, 600, 800, 1000\}$  and  $\epsilon \in \{0, 0.1, 0.2\}$ . Note that  $r(y, N', \epsilon)$  can be computed exactly given the generating model.<sup>11</sup> Also note that under SRA,  $e_j(y)$  is computed based on answers of agents other than  $j$ ; hence, from the perspective of agent  $j$ ,  $e_j(y)$  incorporates the bias of the agents in making their observations. This bias is reflected in the computation of  $E(e_j(y)) = r(y, N', \epsilon)$ .

Based on our calculation of  $r(y, N', \epsilon)$ , we next define three quantities that will be utilized to define our performance measures. In defining all of these quantities, we assume that all agents other than  $j$  truthfully report their (potentially biased) observations.

1. **Truthful reward.** First, for each instance, we define the expected reward of agent  $j$  under truthful behavior:

$$\mathbf{truthful-reward}(N', \epsilon) \triangleq \sum_{y \in \mathcal{Y}} P(Y_j = Y_{j'} = y) r(y, N', \epsilon). \quad (22)$$

Note that in defining this reward, we account for the fact that agent  $j$  believes that she makes perfect observations while her peer  $j'$  is potentially biased. In particular, we have,

$$P(Y_{j'} = Y_j = y) = \sum_{x \in \mathcal{X}} P_X(x) p_y(x) p'_y(x),$$

where  $\mathbf{p}(x)$  is defined in Equation 20 and  $\mathbf{p}'(x)$  is defined in Equation 21 (capturing the fact that the peer agent is biased).

<sup>11</sup> In SRA,  $\bar{f}_j(y)$  is a discrete random variable taking  $N' + 1$  possible values in the set  $\{1/N', 2/N', \dots, (N' + 1)/N'\}$ . In particular,  $\bar{f}_j(y) = 1/N' + Z(y)/N'$ , where  $Z(y)$  is a binomial random variable arising from  $N'$  trials, with probability of success equalling the probability of agreement on  $y$  (between two (potentially) biased agents). Thus the expectation of  $e_j(y)$ , which is also a discrete random variable and a function of  $\bar{f}_j(y)$ , can be exactly computed.

2. **Optimal Reward.** Next, we define the expected reward of agent  $j$  from the optimal reporting strategy that maximizes her expected reward (which could potentially entail lying). In doing so, we address an important consideration. Although the question simply asks for the interval in which the true delay of the mover lies, agent  $j$  can actually observe the true delay. Thus the optimal report must be determined conditioned not on the true answer, but the true delay (from which the true answer can be determined). We assume that agent  $j$  can accurately observe the delay to within a minute and she stops observing if the delay is larger than 180 min (i.e., 3 hours).<sup>12</sup> Formally, we assume that the observations of the delay lie in the finite set  $\bar{\mathcal{Y}} = \{a \text{ to } a + 1; \text{ for } a \in \{0, 1, \dots, 179\}\} \cup \{180 \text{ to } \infty\}$ . We denote the observed delay of agent  $j$  by the random element  $\bar{Y}_j \in \bar{\mathcal{Y}}$ . As before, since the delay is exponentially distributed with mean  $\mu_x$  given the type  $x \in \mathcal{X}$ , the conditional probability of observing the delay  $y = l$  to  $u \in \bar{\mathcal{Y}}$  is given by:

$$\bar{p}_y(x) = \exp\left(-\frac{l}{\mu_x}\right) - \exp\left(-\frac{u}{\mu_x}\right). \quad (23)$$

Accounting for this consideration, we finally define the expected reward from the optimal reporting strategy as:

$$\text{optimal-reward}(N', \epsilon) \triangleq \sum_{y \in \bar{\mathcal{Y}}} P(\bar{Y}_j = y) \max_{y' \in \mathcal{Y}} P(Y_{j'} = y' | \bar{Y}_j = y) r(y', N', \epsilon) \quad (24)$$

$$= \sum_{y \in \bar{\mathcal{Y}}} \max_{y' \in \mathcal{Y}} P(Y_{j'} = y', \bar{Y}_j = y) r(y', N', \epsilon). \quad (25)$$

Here we have,

$$P(Y_{j'} = y', \bar{Y}_j = y) = \sum_{x \in \mathcal{X}} P_X(x) \bar{p}_y(x) p'_{y'}(x),$$

where  $\mathbf{p}'(x)$  is defined in Equation 21 (capturing the fact that the peer agent is biased) and  $\bar{\mathbf{p}}(x)$  is defined in Equation 23 (capturing the fact that the agent observes the true delay in the set  $\bar{\mathcal{Y}}$ ).

3. **Effortless reward.** Finally, as a baseline, we define the reward obtained by agent  $j$  by choosing an answer  $y \in \mathcal{Y}$  uniformly at random without making any observation:

$$\text{effortless-reward}(N', \epsilon) \triangleq \sum_{y \in \mathcal{Y}} \frac{1}{|\mathcal{Y}|} P(Y_{j'} = y) r(y, N', \epsilon), \quad (26)$$

where we have,

$$P(Y_{j'} = y) = \sum_{x \in \mathcal{X}} P_X(x) p'_y(x),$$

where  $\mathbf{p}'(x)$  is defined in Equation 21.

<sup>12</sup> This assumption is for simplicity of computation of the optimal reporting strategy as a function of the observed delay. Our findings are not expected to change significantly if the optimal reports are computed conditioned on finer feedback information.

We finally define our two main performance measures.

1. **Lying gain.** The first measure we define captures the percentage gain in expected reward of agent  $j$  by reporting optimally rather than simply being truthful. It is defined as,

$$\mathbf{lying-gain}(N', \epsilon) \triangleq \frac{\mathbf{optimal-reward}(N', \epsilon) - \mathbf{truthful-reward}(N', \epsilon)}{\mathbf{truthful-reward}(N', \epsilon)} \times 100\%. \quad (27)$$

Ideally, we would like this gain to be small.

2. **Truthful coverage.** Although the measure that we define above is a natural one to consider, it provides at best a partial picture of the incentives generated by the mechanism. In particular, one way of ensuring a small lying-gain is to simply add a very large fixed reward to the reward under the mechanism so that the denominator in Equation 27 becomes large. By scaling this fixed reward, one can ensure that the lying-gain is as small as one desires without changing the incentive properties of the mechanism. By doing so, even mechanisms with poor incentives for truthful behavior can result in small lying-gain. In other words, although the lying-gain is invariant to multiplicative scaling of the rewards under the mechanism, it is not invariant to additive shifts of the rewards. To address this concern, we define the following relative performance measure, which is invariant to both additive shifts as well as multiplicative scaling of the rewards.

$$\mathbf{truthful-coverage}(N', \epsilon) \triangleq \frac{\mathbf{truthful-reward}(N', \epsilon) - \mathbf{effortless-reward}(N', \epsilon)}{\mathbf{optimal-reward}(N', \epsilon) - \mathbf{effortless-reward}(N', \epsilon)} \times 100\%. \quad (28)$$

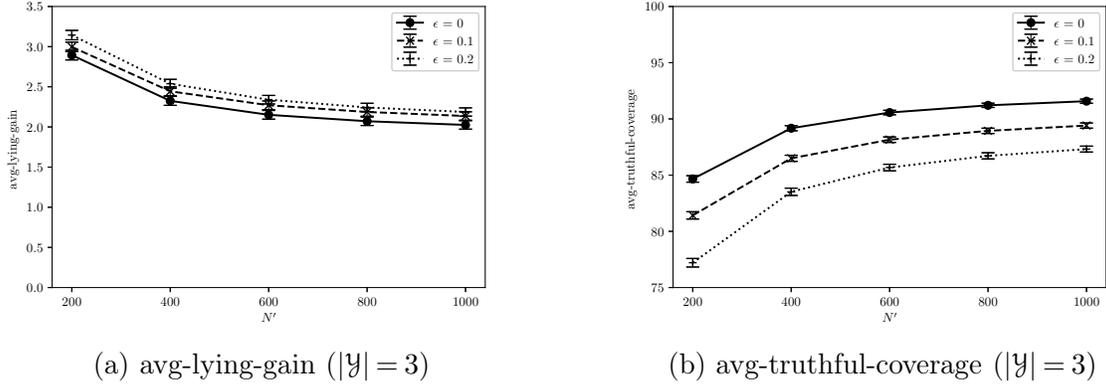
This quantity measures the fraction of the incremental gain resulting from optimal reporting as compared to reporting randomly, that can be attained by truthful reporting. This quantity should ideally be large.

We note that due to the bias in the population, and given that agent  $j$  evaluates the incentives for lying conditioned on the observed delay as opposed to the true answer, being truthful is not guaranteed to be optimal for agent  $j'$  under SRA even in the  $N' \rightarrow \infty$  limit. Hence, it is expected that  $\mathbf{lying-gain}(N', \epsilon) > 0\%$  and  $\mathbf{truthful-coverage}(N', \epsilon) < 100\%$ .

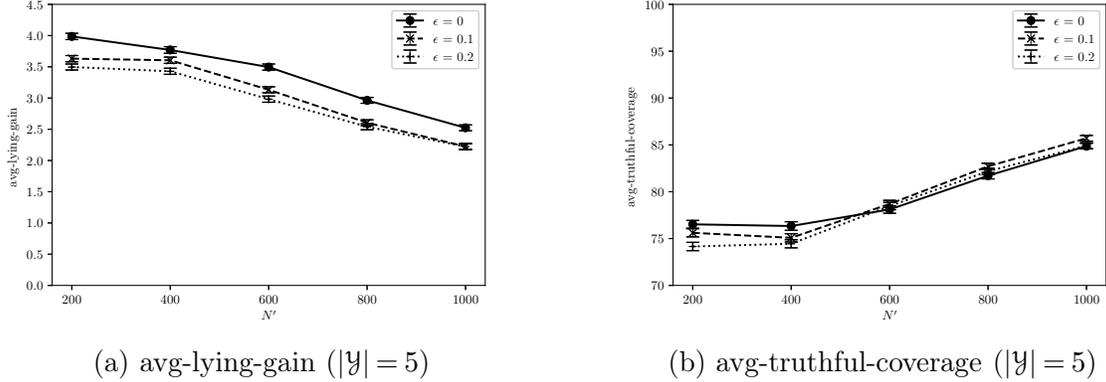
We define  $\mathbf{avg-lying-gain}(N', \epsilon)$  to be the average across 10000 instances of the  $\mathbf{lying-gain}(N', \epsilon)$ . Similarly, we define  $\mathbf{avg-truthful-coverage}(N', \epsilon)$  to be the average across 10000 instances of the  $\mathbf{truthful-coverage}(N', \epsilon)$ . When the context is clear, for notational convenience, we will refer to these aggregate quantities as avg-lying-gain and avg-truthful-coverage respectively.

**Results.** The aggregate performance measures and standard errors<sup>13</sup> for the different values of  $N'$  and  $\epsilon$  are presented for the  $|\mathcal{Y}| = 3$  setting in Figure 1 and for the  $|\mathcal{Y}| = 5$  setting in Figure 2. We note two main observations.

<sup>13</sup> The standard errors of the avg-lying-gain and avg-truthful-coverage are the standard errors of these mean quantities, defined as 1.96 times the empirical standard deviation of the lying-gains or truthful-coverages divided by the square root of the sample size (10000).



**Figure 1** The avg-lying-gains along with standard errors under SRA for  $N' \in \{100, 200, 300, 400, 500\}$  (X-axis) for different values of  $\epsilon$  and for the  $|\mathcal{Y}| = 3$  setting.



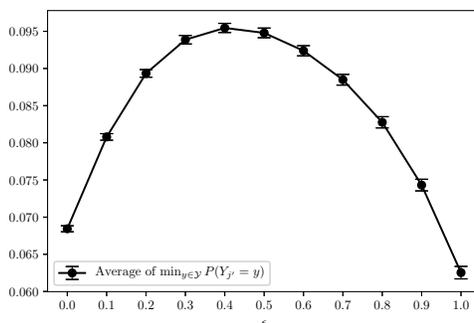
**Figure 2** The avg-lying-gains (Y-axis) along with standard errors under SRA for  $N' \in \{100, 200, 300, 400, 500\}$  (X-axis) for different values of  $\epsilon$ , for the  $|\mathcal{Y}| = 5$  setting.

First, SRA displays reasonably good performance in incentivizing truthful behavior despite the limited number of tasks, even when agent  $j$  believes that the rest of the population is biased. For instance, for  $|\mathcal{Y}| = 3$  and  $N' = 1000$ , the avg-lying-gain is at most about 2.5% across all settings of biases. Correspondingly, for  $|\mathcal{Y}| = 5$  and  $N' = 1000$ , the avg-lying-gain is at most about 2.7%. Additionally, SRA displays reasonably good truthful-coverage, e.g., for  $|\mathcal{Y}| = 3$  and  $N' = 1000$ , truthful behavior attains about 86% of the maximal potential gain over random reporting on average. For  $|\mathcal{Y}| = 5$  and  $N' = 1000$ , truthful behavior attains about 84% of the maximal potential gain over random reporting on average.

Second, we note an interesting phenomenon in the  $|\mathcal{Y}| = 5$  setting: when agent  $j$ 's belief about the bias in the population increases, i.e.,  $\epsilon$  increases from 0 to 0.2, her incentive to lie seems to *decrease* under almost all performance measures, especially for  $N' \in \{600, 800, 1000\}$ . In particular, the lying-gain is smaller and truthful-coverage is larger in aggregate for  $\epsilon = 0.1$  as compared to  $\epsilon = 0$ .

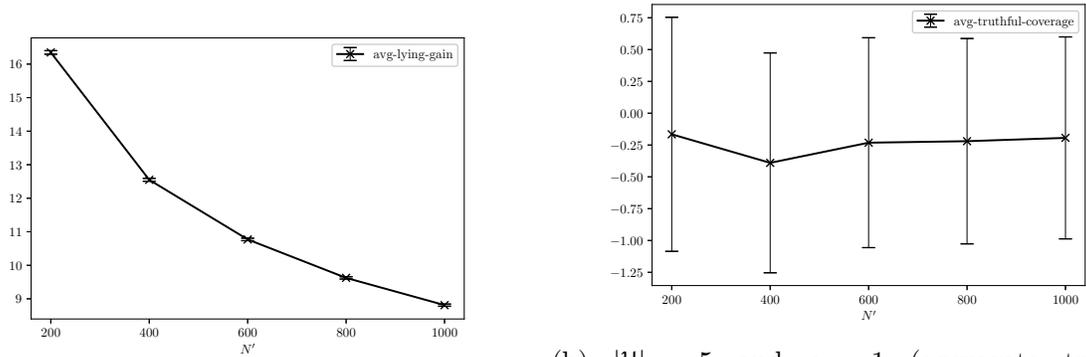
This is quite unlike the  $|\mathcal{Y}| = 3$  setting, in which the incentive to lie *increases* with  $\epsilon$ , as expected, under all measures.

This observation is surprising, given that population homogeneity is crucial to the properties of SRA. It can be explained in light of the fact that adding a small random bias typically increases the probability of rare answers in the population. Informally, adding a small observation noise makes the distribution of answers “better mixed,” since the reporting probability of answers with a low probability of occurrence improves simply because of the noise. This results in a higher inequality gap in the Cauchy-Schwarz inequality, as shown in Proposition 4.1. In other words, the biased generating model is better separated on average than the unbiased model. In turn, this fact results in faster convergence of the popularity indices of the answers to their stable values. In Figure 3, we show the average across the 10000 instances of the smallest probability of observation (in  $\mathcal{Y}$ ) made by a generic biased agent in the population, where we observe that this average indeed increases with  $\epsilon$  when  $\epsilon$  is small, before decreasing (the mixing effect is highest around  $\epsilon \approx 0.4$ ). Thus, in finite  $N'$  settings, the incentive to lie can potentially be higher when the population is assumed to be unbiased since the population indices are expected to be farther from their asymptotic stable values, as compared to the case where the population is expected to be mildly biased.



**Figure 3** For the  $|\mathcal{Y}| = 5$  setting, the average across 10000 instances of  $\min_{y \in \mathcal{Y}} P(Y_{j'} = y)$  (Y axis) for some generic agent  $j'$  in the population as  $\epsilon$  (X axis), i.e., the magnitude of the bias in the population, increases.

However, as the bias in the population becomes large, i.e.,  $\epsilon$  increases, this effect is overpowered by the increased incentive to lie, since the generating model of the peer agent’s observations starts to look starkly different from agent  $j$ ’s model, i.e., the response homogeneity assumption is violated to a higher degree. We indeed verify this to be the case. In Figure 4, we plot the various aggregate measures in the  $|\mathcal{Y}| = 5$  setting for  $\epsilon = 1$ , i.e., when the distributions of observations conditioned on the type are completely uncorrelated for agent  $j$  and her randomly chosen peer agent. As expected, the incentive to lie is significantly higher across all values of  $N'$  in this case compared to settings with smaller values of  $\epsilon$ .



(a)  $|\mathcal{Y}| = 5$  and  $\epsilon = 1$  (aggregate lying-gains)      (b)  $|\mathcal{Y}| = 5$  and  $\epsilon = 1$  (aggregate truthful-coverages)

**Figure 4** For  $|\mathcal{Y}| = 5$ , the different aggregate measures (Y-axis) along with standard errors under SRA for  $N' \in \{100, 200, 300, 400, 500\}$  (X-axis) in the fully biased setting, i.e.,  $\epsilon = 1$  (the distributions of observations of agent  $j$  and her randomly chosen peer are independent).

Overall, these observations suggest that the belief that there exists a mild observation bias in the population may in fact improve incentives for truthful behavior in finite  $N'$  settings, as long as this bias is small enough.

## 6.2. Performance on subjective evaluations

As we have discussed in Section 2, strongly truthful or informed truthful mechanisms for eliciting subjective evaluations from heterogeneous agents require multiple evaluations from each agent. Moreover, these mechanisms require that each agent uses the same strategy for each evaluation. These constraints may hinder the practical applicability of these mechanisms in many platform environments. In the face of these drawbacks, mechanisms tailored to homogeneous response settings that require a single task per agent could be a practical alternative. In this section, we hence evaluate the performance of SRA, PTSC, and CA-HR for incentivizing truthful responses to subjective evaluations in real settings.

**Datasets.** We test these mechanisms using publicly available rating datasets from different online platforms. In particular, we consider the following three datasets.

1. **Goodreads.** We consider book rating data from Goodreads, which is a popular book review platform.<sup>14</sup> We restrict our attention to books belonging to two largest and similarly-sized categories: (a) romance and (b) fantasy and paranormal. We assume books in each of these categories to be a priori statistically similar and we test the performance of different mechanisms for these two categories independently.

<sup>14</sup>The data is publicly available at <https://sites.google.com/eng.ucsd.edu/ucsdbookgraph/home>. The source requires us to cite Wan and McAuley (2018) and Wan et al. (2019).

2. **Amazon.** We next consider product rating data from the e-commerce platform, Amazon.<sup>15</sup> We restrict ourselves to the “Clothing, Shoes, and Jewelery” (CSJ) category, which is by far the largest product category other than books.<sup>16</sup>
3. **Netflix.** We finally consider movie rating data from the streaming platform, Netflix, that was released as part of the Netflix Prize challenge.<sup>17</sup>

In all of the above cases, the ratings are integers on a scale from 1 to 5. Moreover, the ratings are expected to have a strongly subjective influence, especially so in the case of books and movies. Table 2 provides some basic information about these datasets.

**Table 2** Properties of datasets. The rating strength represents the highest lower bound on the number of ratings given by the top 1000 high-contributing users.

	No. of entities	No. of users	No. of ratings	Rating strength
Goodreads: romance	334957 books	198141	3565378	472
Goodreads: fantasy/paranormal	258212 books	256088	3424641	278
Amazon: CSJ	2681297 products	12483678	32292099	104
Netflix	17770 movies	480189	100480507	2087

**Testing procedure.** In each of the above cases, we focus on the top 1000 users who have rated the most number of entities (books, movies, or products). Assuming that all ratings in the dataset are truthful, we estimate the reporting behavior of these users and investigate their incentives for lying under the various mechanisms. Formally, let  $\mathcal{H}$  denote the set of high-contributing users and consider a user  $i \in \mathcal{H}$ . Based on  $i$ 's ratings across the books they have rated and the ratings of randomly chosen peers for these books, we estimate the joint distribution of the rating of  $i$  and that of a randomly chosen peer agent for a random book that they rate. Let us denote this estimate as  $(Q_i(y, y'))_{y, y' \in \mathcal{Y}}$ . Similarly, we estimate the joint distribution of the ratings of two randomly chosen agents for a random book by sampling two agents at random for each book in the data set and computing the empirical distribution of the resulting answers. Let us denote this estimate as  $(\bar{Q}(y, y'))_{y, y' \in \mathcal{Y}}$ . The estimates  $Q_i$  and  $\bar{Q}$  are expected to be different, in line with the expectation that agent  $i$ 's responses are statistically different from a randomly chosen agent in the population due to the subjectivity of responses.

Based on the estimates  $Q_i$  and  $\bar{Q}$ , we can estimate the truthful-coverage of each mechanism for each agent  $i \in \mathcal{H}$ , assuming that (a)  $i$ 's belief about the joint distribution of her rating and that of

<sup>15</sup> The data is publicly available at <https://nijianmo.github.io/amazon/index.html>. The source requires us to cite Ni et al. (2019).

<sup>16</sup> The Goodreads data is richer than Amazon's rating data for books.

<sup>17</sup> The data is publicly available at <https://www.kaggle.com/netflix-inc/netflix-prize-data>. Information about the Netflix prize is available at [https://en.wikipedia.org/wiki/Netflix\\_prize](https://en.wikipedia.org/wiki/Netflix_prize).

a random peer for a random book is identical to  $Q_i$ , and (b) her belief about the joint distribution of the ratings of two randomly chosen agents for a random book is identical to  $\bar{Q}$ . Note that we focus on truthful-coverage since, as we discussed earlier, the lying-gain is not invariant to additive shifts in the rewards.

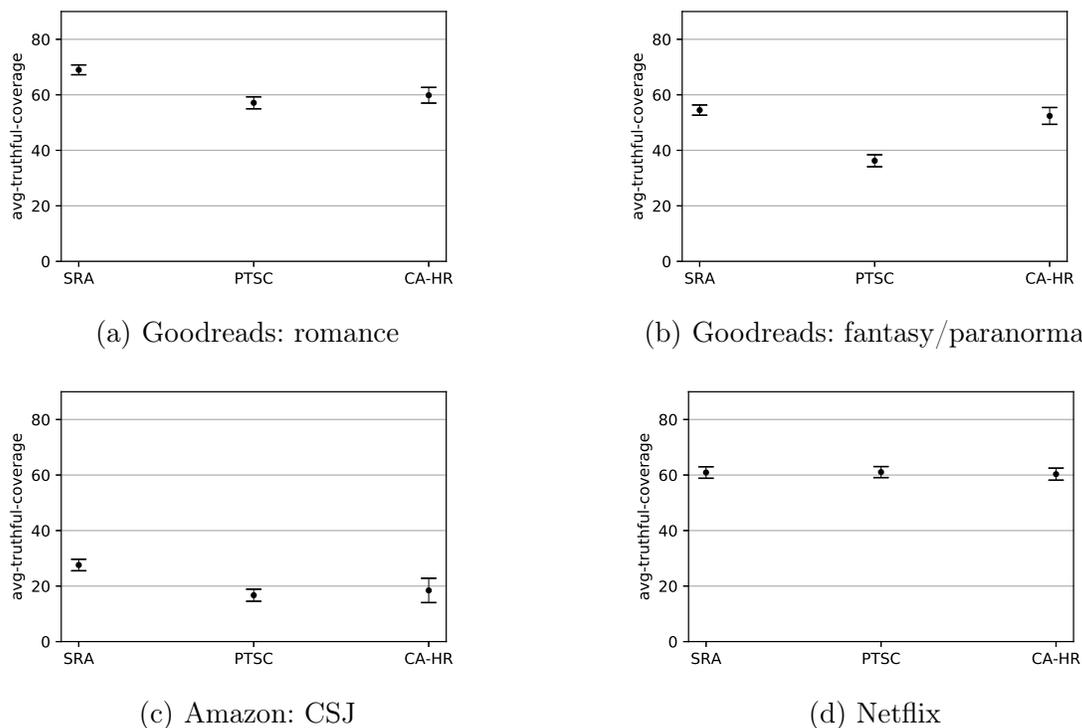
**Results.** The results are shown in Figure 5. First, we observe that SRA achieves an average truthful-coverage of at least 50% across all settings except for Amazon. In this latter case, the performance of all mechanisms is relatively poor. The reason for this may be that the definition of the CSJ category is quite broad, and significant diversity is expected across the products in this category; hence, the assumption of the products being a priori statistically similar likely doesn't hold.

Next, we observe that SRA outperforms PTSC in all settings except for the case of Netflix, where their performance is statistically similar. SRA outperforms CA-HR in the Goodreads setting for the romance category and in the case of Amazon, while their performance is statistically similar in the other two cases. To investigate the difference in SRA and PTSC, we consider the estimate  $\bar{Q}$  of the joint distribution of two agents' responses to a common evaluation. The incentive to lie for an agent  $i$  stems from two sources: (a)  $Q_i$  may be different from  $\bar{Q}$ , and (b)  $\bar{Q}$  may not satisfy the conditions necessary for inducing truthful behavior even when all agents are identical. We find that in all settings,  $\bar{Q}$  satisfies the Cauchy-Schwarz property required for SRA to be truthful (Equation 11), while the self-prediction property that is needed for PTSC to be truthful (see item 2 in Section 2.1) is not satisfied in any of the settings; see Table 3.

Moreover, we find that the response distribution  $\bar{Q}$  is 'clustered' in all settings, as defined in Definition 10 (from Shnayder et al. (2016)) in the Appendix. This implies that under the CA-HR mechanism, the agents can achieve the same payoff by merging their responses. For example, considering the  $\bar{Q}$  from the Netflix setting, we find that agents need not distinguish between the ratings 4 and 5 under the CA-HR mechanism (see Table 3). This points to the importance of the distinction between strong and informed truthfulness in these settings.

**Table 3** Properties of population average joint distribution of responses  $\bar{Q}$ .

	Satisfies CS property	Satisfies self-prediction	Clustered ratings
Goodreads: romance	Yes	No	(1, 2)
Goodreads: fantasy/paranormal	Yes	No	(1, 2, 3)
Amazon: CSJ	Yes	No	(2, 3)
Netflix	Yes	No	(4, 5)



**Figure 5** The avg-truthful-coverages along with standard errors under the different mechanisms in different settings.

## 7. Discussion and Conclusion

In the paper, we focus on the practical setting of reputation systems in online platforms where objective evaluations must be strongly incentivized; ideally, without imposing any constraints on the number of evaluations performed by each agent. Our results show that SRA is the first mechanism that achieves this goal.

While there are other mechanisms, such as those of Kong and Schoenebeck (2019), Kong (2020), or CA (Shnayder et al. 2016), that incentivize truthful behavior despite response heterogeneity across agents, these mechanisms incur a high operational cost of requiring multiple evaluations from each agent, which could be prohibitive in many scenarios, including in online platforms. On the other hand, our numerical evaluations show that the truthfulness property of SRA is robust to mild degrees of heterogeneity and subjectivity in the population. This observation overall suggests that SRA can be a simpler alternative to these more complex mechanisms in settings where response homogeneity is a reasonable approximation to the mild degree of heterogeneity and subjectivity expected in the evaluations. Eliciting objective evaluations in online platforms is one such setting that we focused on in the paper. Additionally, our tests on real data show that SRA generates strong incentives for truthful behavior even when evaluations are expected to be highly subjective.

At the same time, we acknowledge that there are settings where other mechanisms could be preferable over SRA. Moreover, metrics such as lying-gain or truthful-coverage may not adequately

inform the practical utility of mechanisms in these settings and other operational considerations may take precedence. For example, in applications such as crowdsourcing and peer-grading,<sup>18</sup> agents typically perform several evaluations in a short span of time. Moreover, subjectivity in evaluations could be a major concern in settings like peer-grading for courses in the arts and the humanities. In this case, Kong (2020)'s mechanism would provide strong truthfulness guarantees without requiring the homogeneity assumption and hence could be preferable over SRA, even if, hypothetically, it turns out to be the case that SRA achieves better truthful-coverage or lying-gain than Kong's mechanism in homogeneous settings with comparable data. As another example, if the responses can be validated to be self-predicting, then PTSC may be preferable owing to the simpler description of the agreement rewards.

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, behavioral, and strategic concerns (Hu et al. 2017, Filippas et al. 2018, Nosko and Tadelis 2015). We believe that appropriate incentive mechanisms that are simple and intuitive can go a long way in addressing some of these concerns, and hence 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.

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. It is important to understand the fundamental tradeoffs faced by dynamic mechanisms that balance incentives with different statistical accuracy objectives in such situations. We are optimistic that our framework and insights can be used as building blocks in this pursuit.

## References

Aldous, D. J. (1985). Exchangeability and related topics. In *École d'Été de Probabilités de Saint-Flour XIII-1983*, pages 1–198. Springer.

<sup>18</sup> Peer-grading is the idea of having students grade each others' assignments and examinations. This idea is key in obtaining a scalable solution to the problem of grading in massive open online courses (MOOC). Incentivizing students to grade truthfully is an important concern in such settings.

- 
- Brier, G. W. (1950). Verification of forecasts expressed in terms of probability. *Monthly weather review*, 78(1):1–3.
- Cover, T. and Thomas, J. (2012). *Elements of Information Theory*. Wiley.
- Dasgupta, A. and Ghosh, A. (2013). Crowdsourced judgement elicitation with endogenous proficiency. In *Proceedings of the 22nd International Conference on World Wide Web*, pages 319–330.
- Faltings, B., Jurca, R., and Radanovic, G. (2017). Peer truth serum: incentives for crowdsourcing measurements and opinions. *arXiv preprint arXiv:1704.05269*.
- Filippas, A., Horton, J. J., and Golden, J. (2018). Reputation inflation. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 483–484.
- Gao, X. A., Wright, J. R., and Leyton-Brown, K. (2019). Incentivizing evaluation with peer prediction and limited access to ground truth. *Artificial Intelligence*, 275:618–638.
- Gneiting, T. and Raftery, A. E. (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378.
- Harsanyi, J. C., Selten, R., et al. (1988). A general theory of equilibrium selection in games. *MIT Press Books*, 1.
- Hu, N., Pavlou, P. A., and Zhang, J. J. (2017). On self-selection biases in online product reviews. *MIS Q.*, 41(2):449–471.
- Jøsang, A., Ismail, R., and Boyd, C. (2007). A survey of trust and reputation systems for online service provision. *Decision support systems*, 43(2):618–644.
- Jurca, R. and Faltings, B. (2005). Enforcing truthful strategies in incentive compatible reputation mechanisms. In *International Workshop on Internet and Network Economics*, pages 268–277. Springer.
- Jurca, R. and Faltings, B. (2008). Incentives for expressing opinions in online polls. In *Proceedings of the 9th ACM Conference on Electronic Commerce*, pages 119–128.
- Jurca, R. and Faltings, B. (2011). Incentives for answering hypothetical questions. In *Workshop on Social Computing and User Generated Content, EC-11*.
- Kong, Y. (2020). Dominantly truthful multi-task peer prediction with a constant number of tasks. In *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 2398–2411. SIAM.
- Kong, Y. and Schoenebeck, G. (2019). An information theoretic framework for designing information elicitation mechanisms that reward truth-telling. *ACM Transactions on Economics and Computation (TEAC)*, 7(1):1–33.
- Li, L., Tadelis, S., and Zhou, X. (2020). Buying reputation as a signal of quality: Evidence from an online marketplace. *The RAND Journal of Economics*, 51(4):965–988.

- Liu, Y. and Chen, Y. (2016). Learning to incentivize: eliciting effort via output agreement. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, pages 3782–3788.
- Luca, M. (2017). Designing online marketplaces: Trust and reputation mechanisms. *Innovation Policy and the Economy*, 17(1):77–93.
- Miller, N., Resnick, P., and Zeckhauser, R. (2005). Eliciting informative feedback: The peer-prediction method. *Management Science*, 51(9):1359–1373.
- Myerson, R. B. (2013). *Game theory*. Harvard University Press.
- Ni, J., Li, J., and McAuley, J. (2019). Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 188–197.
- Nosko, C. and Tadelis, S. (2015). The limits of reputation in platform markets: An empirical analysis and field experiment. Technical report, National Bureau of Economic Research.
- Prelec, D. (2004). A Bayesian truth serum for subjective data. *Science*, 306(5695):462–466.
- Radanovic, G. and Faltings, B. (2013). A robust bayesian truth serum for non-binary signals. In *Proceedings of the Twenty-Seventh AAAI Conference on Artificial Intelligence*, pages 833–839.
- Radanovic, G. and Faltings, B. (2015). Incentives for subjective evaluations with private beliefs. In *Proceedings of the 29th AAAI Conference on Artificial Intelligence (AAAI’15)*.
- Radanovic, G., Faltings, B., and Jurca, R. (2016). Incentives for effort in crowdsourcing using the peer truth serum. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 7(4):48.
- Resnick, P., Kuwabara, K., Zeckhauser, R., and Friedman, E. (2000). Reputation systems. *Communications of the ACM*, 43(12):45–48.
- Savage, L. J. (1971). Elicitation of personal probabilities and expectations. *Journal of the American Statistical Association*, 66(336):783–801.
- Schoenebeck, G. and Yu, F.-Y. (2020). Two strongly truthful mechanisms for three heterogeneous agents answering one question. In *International Conference on Web and Internet Economics*, pages 119–132. Springer.
- Schoenebeck, G. and Yu, F.-Y. (2021). Learning and strongly truthful multi-task peer prediction: A variational approach. In *12th Innovations in Theoretical Computer Science Conference (ITCS 2021)*.
- Shnayder, V., Agarwal, A., Frongillo, R., and Parkes, D. C. (2016). Informed truthfulness in multi-task peer prediction. In *Proceedings of the 2016 ACM Conference on Economics and Computation*, pages 179–196. ACM.
- Tadelis, S. (2016). Reputation and feedback systems in online platform markets. *Annual Review of Economics*, 8:321–340.

- 
- Van Damme, E. (2002). Strategic equilibrium. *Handbook of Game Theory with Economic Applications*, 3:1521–1596.
- Von Ahn, L. and Dabbish, L. (2004). Labeling images with a computer game. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 319–326. ACM.
- Von Ahn, L. and Dabbish, L. (2008). Designing games with a purpose. *Communications of the ACM*, 51(8):58–67.
- Wan, M. and McAuley, J. (2018). Item recommendation on monotonic behavior chains. In *Proceedings of the 12th ACM Conference on Recommender Systems*, pages 86–94.
- Wan, M., Misra, R., Nakashole, N., and McAuley, J. (2019). Fine-grained spoiler detection from large-scale review corpora. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2605–2610.
- Witkowski, J. and Parkes, D. C. (2012). A robust bayesian truth serum for small populations. In *Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence*, pages 1492–1498.
- Witkowski, J. and Parkes, D. C. (2013). Learning the prior in minimal peer prediction. In *3rd Workshop on Social Computing and User Generated Content at the ACM Conference on Electronic Commerce*.

## Appendix

### A. Proofs

*Proof of Proposition 4.1* Let  $m = \min_{y \in \mathcal{Y}} \|\mathbf{v}(y)\|$ . To see that 2 implies 1, note that  $\theta(\mathbf{v}(y), \mathbf{v}(y')) \geq \gamma$  for all  $y, y' \in \mathcal{Y}$  such that  $y \neq y'$ , implies that  $m > 0$  and that

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

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

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

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 both  $\|\mathbf{v}(y)\|$  and  $\|\mathbf{v}(y')\|$  are non-zero. Then dividing on both sides, we get:

$$\begin{aligned} \frac{|\mathbf{v}(y) \cdot \mathbf{v}(y')|}{\|\mathbf{v}(y')\| \|\mathbf{v}(y)\|} &\leq 1 - \frac{\alpha}{\|\mathbf{v}(y')\| \|\mathbf{v}(y)\|} \\ &\leq 1 - \alpha \end{aligned}$$

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

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

This implies that  $\theta(\mathbf{v}(y), \mathbf{v}(y')) \geq \arccos(1 - \alpha)$ . Note that  $\alpha > 0$  so that  $\arccos(1 - \alpha) > 0$ .  $\square$

*Proof of Theorem 1* 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,

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

In the proof of Proposition 4.1, 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

$$e_j(y) = \frac{K}{\sqrt{\bar{f}_j(y)}}.$$

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

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

$$\begin{aligned}
&\geq \frac{K}{\sqrt{g(y)(1+\epsilon)+1/N'}} - 2\exp(-\epsilon^2\alpha^4N')\frac{K}{\alpha\sqrt{1+\epsilon}} \\
&\geq \frac{K}{\sqrt{g(y)(1+\epsilon+1/(g(y)N'))}} - 2\exp(-\epsilon^2\alpha^4N')\frac{K}{\alpha\sqrt{1+\epsilon}} \\
&\geq \frac{K}{\sqrt{g(y)(1+\epsilon+1/(\alpha^2N'))}} - 2\exp(-\epsilon^2\alpha^4N')\frac{K}{\alpha\sqrt{1+\epsilon}} \\
&\stackrel{(b)}{\geq} \frac{K}{\sqrt{g(y)}}(1-\epsilon-1/(\alpha^2N')) - 2\exp(-\epsilon^2\alpha^4N')\frac{K}{\alpha} \quad (\text{for large enough } N') \\
&\geq \frac{K}{\sqrt{g(y)}}(1-\epsilon-1/(\alpha^2(N-n))) - 2\exp(-\epsilon^2\alpha^4(N-n))\frac{K}{\alpha} \quad (\text{for large enough } N). \quad (29)
\end{aligned}$$

Here (a) follows from Hoeffding's inequality, and (b) is because  $\frac{1}{\sqrt{1+a}} \geq 1-a$  for every  $a \in (0, 1)$ . 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:

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

Next, we also have,

$$\begin{aligned}
\mathbb{E}(e_j(y)) &\leq P(\bar{f}_j(y) - 1/N' \in [g(y)(1-\epsilon), g(y)(1+\epsilon)])\frac{K}{\sqrt{g(y)(1-\epsilon)}} \\
&\quad + \mathbb{E}\left(\mathbf{1}_{\{\bar{f}_j(y) - 1/N' \notin [g(y)(1-\epsilon), g(y)(1+\epsilon)]\}} \frac{K}{\sqrt{f_j(y)}}\right) \\
&\stackrel{(a)}{\leq} \frac{K}{\sqrt{g(y)(1-\epsilon)}} + P(\mathbf{1}_{\bar{f}_j(y) - 1/N' \notin [g(y)(1-\epsilon), g(y)(1+\epsilon)]}) K\sqrt{N'} \\
&\stackrel{(b)}{\leq} \frac{K}{\sqrt{g(y)(1-\epsilon)}} + 2K\sqrt{N'}\exp(-\epsilon^2g(y)^2N') \\
&\leq \frac{K}{\sqrt{g(y)(1-\epsilon)}} + 2K\sqrt{N'}\exp(-\epsilon^2\alpha^4N') \\
&\stackrel{(c)}{\leq} \frac{K}{\sqrt{g(y)}}\left(1 + \frac{\epsilon}{2} + w(\epsilon)\right) + 2K\sqrt{N'}\exp(-\epsilon^2\alpha^4N') \\
&\leq \frac{K}{\sqrt{g(y)}} + \frac{\epsilon K}{2\alpha} + \frac{|w(\epsilon)|K}{\alpha} + 2K\sqrt{N'}\exp(-\epsilon^2\alpha^4N') \\
&\leq \frac{K}{\sqrt{g(y)}} + \frac{\epsilon K}{2\alpha} + \frac{|w(\epsilon)|K}{\alpha} + 2K\sqrt{N}\exp(-\epsilon^2\alpha^4(N-n)). \quad (30)
\end{aligned}$$

Here, (a) results from the fact that  $\bar{f}_j(y) \geq 1/N'$  (because of the smoothing). (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:

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

Thus, we finally have  $|\mathbb{E}(e_j(y)) - \frac{K}{\sqrt{g(y)}}| \leq \sigma(N) = o(N)$ , where  $\sigma(N) \geq 0$  is a function of  $N$  that depends only on  $\alpha$ ,  $n$  and  $K$  and not on  $y$  (note that our bounds explicitly define this function: we have  $w(\epsilon) < \epsilon/2$  for  $\epsilon < 1/2$  and thus  $w(\epsilon)$  can be replaced by  $\epsilon/2$  for  $N-n \geq 2^4 = 16$ ).

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,

$$R(\mathbf{q}^{ij}) \triangleq \sum_{y \in \mathcal{Y}} P(Y_{j'}^i = y, Y_j^i = y) \mathbb{E}(r_j(y)) = \sum_{y \in \mathcal{Y}} \mathbb{E}(r_j(y)) \sum_{x \in \mathcal{X}} P_X(x) p_y(x) \sum_{y' \in \mathcal{Y}} p_{y'}(x) q_y^{ij}(y'). \quad (31)$$

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,

$$\begin{aligned}
R(\mathbf{q}^{ij}) &= \sum_{y \in \mathcal{Y}} \mathbb{E}(e_j(y)) \sum_{y' \in \mathcal{Y}} q_y^{ij}(y') \sum_{x \in \mathcal{X}} P_X(x) p_y(x) p_{y'}(x) \\
&\stackrel{(a)}{\leq} \sum_{y \in \mathcal{Y}} \mathbb{E}(e_j(y)) \sum_{y' \in \mathcal{Y}} q_y^{ij}(y') \left( \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_y(x)^2} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y'}(x)^2} - \alpha \mathbf{1}_{\{y \neq y'\}} \right) \\
&\leq \sum_{y \in \mathcal{Y}} \left( \frac{K}{\sqrt{g(y)}} + \sigma(N) \right) \sum_{y' \in \mathcal{Y}} q_y^{ij}(y') \left( \sqrt{g(y)g(y')} - \alpha \mathbf{1}_{\{y \neq y'\}} \right) \\
&\leq K \sum_{y' \in \mathcal{Y}} \sqrt{g(y')} - \alpha K \sum_{y' \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} \mathbf{1}_{\{y \neq y'\}} q_y^{ij}(y') + \sigma(N) \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}} q_y^{ij}(y') \sqrt{g(y)g(y')} \quad (32) \\
&\stackrel{(b)}{\leq} K \sum_{y' \in \mathcal{Y}} \sqrt{g(y')} - \alpha K \mathbf{1}_{\{\mathbf{q}^{ij} \neq \mathbf{t}\}} + |\mathcal{Y}| \sigma(N). \quad (33)
\end{aligned}$$

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 (32),

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

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.

*Proof of Proposition 5.1* We have,

$$\begin{aligned}
\Gamma(Z_1, Z_2) &= \sum_{y \in \mathcal{Y}} \sqrt{\sum_{x \in \mathcal{X}, y_1 \in \mathcal{Y}, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) q_y(y_1) q_y(y_2)} \\
&= \sum_{y \in \mathcal{Y}} \sqrt{\sum_{y_1 \in \mathcal{Y}, y_2 \in \mathcal{Y}} q_y(y_1) q_y(y_2) \sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x) p_{y_2}(x)} \\
&\stackrel{(a)}{\leq} \sum_{y \in \mathcal{Y}} \sqrt{\sum_{y_1 \in \mathcal{Y}, y_2 \in \mathcal{Y}} q_y(y_1) q_y(y_2) \left( \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x)^2} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_2}(x)^2} - \delta(P_X, \mathbf{p}) \mathbf{1}_{\{y_1 \neq y_2\}} \right)} \\
&= \sum_{y \in \mathcal{Y}} \sqrt{\left( \sum_{y_1 \in \mathcal{Y}} q_y(y_1) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x)^2} \right)^2 - \delta(P_X, \mathbf{p}) \sum_{y_1 \in \mathcal{Y}, y_2 \in \mathcal{Y}} q_y(y_1) q_y(y_2) \mathbf{1}_{\{y_1 \neq y_2\}}} \\
&\stackrel{(b)}{\leq} \sum_{y \in \mathcal{Y}, y_1 \in \mathcal{Y}} q_y(y_1) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x)^2} - \frac{\delta(P_X, \mathbf{p})}{2} \sum_{y \in \mathcal{Y}} \frac{\left( \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} q_y(y') q_y(y'') \mathbf{1}_{\{y' \neq y''\}} \right)}{\sum_{y_1 \in \mathcal{Y}} q_y(y_1) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x)^2}} \\
&\stackrel{(c)}{\leq} \Gamma(Y_1, Y_2) - \frac{\delta(P_X, \mathbf{p})}{2} \sum_{y \in \mathcal{Y}} \frac{\left( \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} q_y(y') q_y(y'') \mathbf{1}_{\{y' \neq y''\}} \right)}{\Gamma(Y_1, Y_2)} \\
&\stackrel{(d)}{\leq} \Gamma(Y_1, Y_2) - \frac{\delta(P_X, \mathbf{p})}{2\sqrt{|\mathcal{Y}|}} \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} q_y(y') q_y(y'') \mathbf{1}_{\{y' \neq y''\}} \\
&\stackrel{(e)}{\leq} \Gamma(Y_1, Y_2) - \frac{\delta(P_X, \mathbf{p}) \Omega(\mathbf{q})^2 (|\mathcal{Y}| - 1)}{2\sqrt{|\mathcal{Y}|}}.
\end{aligned}$$

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 2* The first statement follows from Theorem 1: 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 1 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  $\alpha$ ,  $n$  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 (30)

in the proof of Theorem 1, 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  $|\sigma'(N)| = o(1)$ , and it depends on  $\alpha$ ,  $K$ ,  $n$  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  $b(y) < \delta(P_X, \mathbf{p})^2/|\mathcal{Y}|$ ,  $q'_y(y') = 0$  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 4.1. 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,

$$\begin{aligned} G(\mathbf{q}) &= G(\mathbf{q}', \mathbf{q}^{-j}) \\ &\leq \sum_{y \in \mathcal{Y}; b(y) \geq \delta(P_X, \mathbf{p})^2/|\mathcal{Y}|} \mathbb{E}(e_j^s(y)) \sum_{x \in \mathcal{X}} P_X(x) \left[ \sum_{y_1 \in \mathcal{Y}} p_{y_1}(x) q'_y(y_1) \right] \left[ \sum_{y_2 \in \mathcal{Y}} p_{y_2}(x) q_y(y_2) \right] \\ &\stackrel{(a)}{\leq} \sum_{y \in \mathcal{Y}; b(y) \geq \delta(P_X, \mathbf{p})^2/|\mathcal{Y}|} \mathbb{E}(e_j^s(y)) \sqrt{s(y) s'(y)} \\ &\leq \sum_{y \in \mathcal{Y}; b(y) \geq \delta(P_X, \mathbf{p})^2/|\mathcal{Y}|} \left( \frac{K}{\sqrt{s(y)}} + \sigma'(N) \right) \sqrt{s(y) s'(y)} \\ &\leq \sum_{y \in \mathcal{Y}; b(y) \geq \delta(P_X, \mathbf{p})^2/|\mathcal{Y}|} K \sqrt{s'(y)} + |\mathcal{Y}| \sigma'(N) \\ &\stackrel{(b)}{=} K \sum_{y \in \mathcal{Y}} \sqrt{s'(y)} + |\mathcal{Y}| \sigma'(N). \end{aligned} \tag{34}$$

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,

$$\begin{aligned} G(\mathbf{t}) &= \sum_{y \in \mathcal{Y}} \mathbb{E}(e_j^t(y)) g(y) \\ &\geq \sum_{y \in \mathcal{Y}} K \sqrt{g(y)} - \sum_{y \in \mathcal{Y}} \sigma(N) g(y) \\ &\geq K \Gamma(Y_j^i, Y_{j'}^i) - |\mathcal{Y}| \sigma(N) \\ &\stackrel{(a)}{\geq} K \sum_{y \in \mathcal{Y}} \sqrt{s'(y)} - |\mathcal{Y}| \sigma(N). \end{aligned} \tag{35}$$

Here, (a) follows from Proposition 5.1. Finally, (35) and (34) 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 3* 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,  $b(y) \geq \eta$  if  $b(y) \neq 0$ , and since  $s(y) \geq b(y)^2$ , we have  $s(y) \geq \eta^2$  if  $b(y) \neq 0$ . Then using arguments similar to the ones leading up to (30) in the proof of Theorem 1, 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 for all  $y \in \mathcal{Y}$  such that  $b(y) \neq 0$ ,  $|\mathbb{E}(e_j^s(y)) - \frac{K}{\sqrt{s(y)}}| \leq \sigma'(N) = o(1)$ , where  $\sigma(N) \geq 0$  is some function of  $N$  that depends only on  $\alpha, n$  and  $K$ , and  $\sigma'(N) \geq 0$  is some function of  $N$  that depends only on  $\alpha, \eta, n$  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  $\mathbf{q}$ , respectively. 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,

$$\begin{aligned} G(\mathbf{q}) &= \sum_{y \in \mathcal{Y}} \mathbb{E}(e_j^s(y)) s(y) \\ &\stackrel{(a)}{\leq} \sum_{y \in \mathcal{Y}} K \sqrt{s(y)} + \sum_{y \in \mathcal{Y}} s(y) \sigma'(N) \\ &\leq K\Gamma(Z_j^i, Z_{j'}^i) + |\mathcal{Y}| \sigma'(N). \end{aligned} \quad (36)$$

Here, (a) follows from the fact that if  $b(y) = 0$ , then  $s(y) = 0$  and moreover, for any  $y$  such that  $b(y) \neq 0$ , we have  $|\mathbb{E}(e_j^s(y)) - \frac{K}{\sqrt{s(y)}}| \leq \sigma'(N)$  from above. Similarly, we can show that

$$\begin{aligned} G(\mathbf{t}) &= \sum_{y \in \mathcal{Y}} \mathbb{E}(e_i^s(y)) g(y) \\ &\geq \sum_{y \in \mathcal{Y}} K \sqrt{g(y)} - \sum_{y \in \mathcal{Y}} g(y) \sigma(N) \\ &\geq K\Gamma(Y_j^i, Y_{j'}^i) - |\mathcal{Y}| \sigma(N) \\ &\geq K\Gamma(Z_j^i, Z_{j'}^i) + \frac{K\delta(P_X, \mathbf{p})\Omega(\mathbf{q})^2(|\mathcal{Y}| - 1)}{2\sqrt{|\mathcal{Y}|}} - |\mathcal{Y}| \sigma(N). \end{aligned} \quad (37)$$

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

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

which implies that

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

or that,

$$\Omega(\mathbf{q}) \leq \sqrt{\frac{2|\mathcal{Y}|^{3/2}(\sigma(N) + \sigma'(N))}{K\delta(P_X, \mathbf{p})(|\mathcal{Y}| - 1)}} < \sqrt{\frac{2|\mathcal{Y}|^{3/2}(\sigma(N) + \sigma'(N))}{K\alpha(|\mathcal{Y}| - 1)}}. \quad (38)$$

Now the quantity on the right is  $o(1)$  (depending only on  $\alpha, \eta, n, |\mathcal{Y}|$ , and  $K$ ). Thus for any  $\omega > 0$  and  $\eta > 0$ , there exists some  $N_0$  such that for any  $N \geq N_0$ , any symmetric strategy profile in which the probability of

reporting any answer  $y \in \mathcal{Y}$  is either 0 or at least  $\eta$ , and 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 Theorem 4.* We will use the following notion for the proof.

DEFINITION 8. For any strategy profile  $(\mathbf{q}^j)_{j \in \mathcal{M}}$  across agents, the average reporting strategy excluding the set of agents  $\mathcal{J}$  is defined as

$$\bar{\mathbf{q}}^{-\mathcal{J}}(y) = \frac{1}{M - |\mathcal{J}|} \sum_{j' \in \mathcal{M} \setminus \mathcal{J}} \mathbf{q}^{j'}(y).$$

Manipulating this definition, we have,

$$\bar{\mathbf{q}}^{-\mathcal{J}}(y) = \frac{M}{M - |\mathcal{J}|} \left( \bar{\mathbf{q}}(y) - \frac{1}{M} \sum_{j' \in \mathcal{J}} \mathbf{q}^{j'}(y) \right). \quad (39)$$

We then directly have that

$$\bar{\mathbf{q}}^{-\mathcal{J}}(y) \geq \bar{\mathbf{q}}(y) - \frac{|\mathcal{J}|}{M}.$$

Next, since  $1/(1 - \frac{|\mathcal{J}|}{M})$  is  $1 + \frac{|\mathcal{J}|}{M} + o(\frac{|\mathcal{J}|}{M})$  as  $M \rightarrow \infty$  (from the Taylor series expansion), we can conclude that there exists some  $\kappa > 1$  such that for any  $M$  large enough, we have

$$\bar{\mathbf{q}}^{-\mathcal{J}}(y) \leq \bar{\mathbf{q}}(y) + \frac{\kappa|\mathcal{J}|}{M}. \quad (40)$$

To summarize, for some  $\kappa > 1$  and any  $M$  large enough, we thus have

$$\bar{\mathbf{q}}(y) - \frac{\kappa|\mathcal{J}|}{M} \leq \bar{\mathbf{q}}^{-\mathcal{J}}(y) \leq \bar{\mathbf{q}}(y) + \frac{\kappa|\mathcal{J}|}{M}. \quad (41)$$

We now present the proof of Theorem 4. Let  $j$  be a fixed agent evaluating a fixed task  $i$ . Let  $J$  be her (random) peer on task  $i$ . Now upon observing  $y$ , her expected reward on reporting any  $y'$  such that the probability of reporting  $y'$  is 0 under the population average strategy is 0 (since there is no hope of matching  $y'$  with any peer). We thus focus on only those  $y' \in \mathcal{Y}$  such that their reporting probability is at least  $\eta$ . For any such  $y'$ ,  $j$ 's expected reward on reporting  $y'$  when she observed  $y$  is given by,

$$G_i^j(y, y') = \mathbb{E} \left[ \mathbb{1}_{r_j^i = y'} \frac{K \sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + \sum_{i' \in \mathcal{N} \setminus \mathcal{W}_j} \mathbb{1}_{r_{J_1(i')}^{i'} = y'} \mathbb{1}_{r_{J_2(i')}^{i'} = y'}}} \mid Y_j^i = y \right] \quad (42)$$

$$\stackrel{(a)}{=} \underbrace{\mathbb{E} \left( \frac{K \sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + \sum_{i' \in \mathcal{N} \setminus \mathcal{W}_j} \mathbb{1}_{r_{J_1(i')}^{i'} = y'} \mathbb{1}_{r_{J_2(i')}^{i'} = y'}}} \right)}_A \underbrace{\mathbb{E} \left[ \mathbb{1}_{r_j^i = y'} \mid Y_j^i = y \right]}_B \text{ (probability of matching } y'). \quad (43)$$

Here,  $J_1(i')$  and  $J_2(i')$  are the (random) agents who have evaluated task  $i'$ , chosen to compute the agreement rewards for  $j$ . (a) results from the fact that the agreement rewards are independent of  $Y_j^i$  and  $Y_j^i$ : the former because of the fact that the agreement rewards only depend on the tasks that  $j$  does not perform, and the latter because of the random task allocation policy ( $Y_j^i$  may contain information about  $J$ , but that doesn't give any information about agents who will be utilized in computing the agreement rewards since the agent allocation to each task is i.i.d.).

We first focus on term A in Equation 43, which is the expected reward for matching on the answer  $y'$ . Note that the random variables  $\mathbb{1}_{r_{J_1(i')}=y'} \mathbb{1}_{r_{J_2(i')}=y'}$  across  $i'$  are i.i.d. owing to our random task allocation policy with  $\mathbb{E}(\mathbb{1}_{r_{J_1(i')}=y'} \mathbb{1}_{r_{J_2(i')}=y'})$  defined as follows (for notational simplicity we drop the dependence on  $i'$ ).

$$\mathbb{E}(\mathbb{1}_{r_{J_1}=y'} \mathbb{1}_{r_{J_2}=y'}) = \mathbb{E} \left[ \sum_{x \in \mathcal{X}, y_1, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) q_{y'}^{J_1}(y_1) q_{y'}^{J_2}(y_2) \right]. \quad (44)$$

Here, the latter expectation is over the random choice of  $J_1$  and  $J_2$ . Now we have for a large enough  $M$ ,

$$\begin{aligned} & \mathbb{E} \left[ \sum_{x \in \mathcal{X}, y_1, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) q_{y'}^{J_1}(y_1) q_{y'}^{J_2}(y_2) \right] \\ & \stackrel{(a)}{=} \mathbb{E} \left[ \sum_{x \in \mathcal{X}, y_1, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) q_{y'}^{J_1}(y_1) \bar{q}_{y'}^{-\{J_1, j\}}(y_2) \right] \\ & \stackrel{(b)}{\leq} \mathbb{E} \left[ \sum_{x \in \mathcal{X}, y_1, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) q_{y'}^{J_1}(y_1) (\bar{q}_{y'}(y_2) + 2\kappa/M) \right] \\ & \stackrel{(c)}{\leq} \mathbb{E} \left[ \sum_{x \in \mathcal{X}, y_1, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) q_{y'}^{J_1}(y_1) \bar{q}_{y'}(y_2) \right] + 2\kappa/M \\ & \stackrel{(d)}{\leq} \sum_{x \in \mathcal{X}, y_1, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) \bar{q}_{y'}^{-j}(y_1) \bar{q}_{y'}(y_2) + 2\kappa/M \\ & \stackrel{(e)}{\leq} \sum_{x \in \mathcal{X}, y_1, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) (\bar{q}_{y'}(y_1) + \kappa/M) \bar{q}_{y'}(y_2) + 2\kappa/M \\ & \leq \sum_{x \in \mathcal{X}, y_1, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) \bar{q}_{y'}(y_1) \bar{q}_{y'}(y_2) + 3\kappa/M. \end{aligned} \quad (45)$$

Here, (a) follows from the fact that,  $J_2$  is equally likely to be any of the remaining agents other than  $J_1$  and  $j$ , again by the random task allocation policy. (b) follows from Equation 41. (c) follows from the fact that the coefficient of  $2\kappa/M$  after the expansion is at most 1. (d) follows from taking expectation over  $J_1$ , who is equally likely to be any agent other than  $j$ . Finally, (e) again follows from Equation 41. Similarly, we have

$$\mathbb{E} \left[ \sum_{x \in \mathcal{X}, y_1, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) q_{y'}^{J_1}(y_1) q_{y'}^{J_2}(y_2) \right] \geq \sum_{x \in \mathcal{X}, y_1, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) \bar{q}_{y'}(y_1) \bar{q}_{y'}(y_2) - 3\kappa/M \quad (46)$$

for any large enough  $M$ . Let  $\mathbb{E}(\mathbb{1}_{r_{J_1(i')}=y} \mathbb{1}_{r_{J_2(i')}=y})$  be denoted as  $h(y)$ , and define

$$s(y) = \sum_{x \in \mathcal{X}, y_1, y_2 \in \mathcal{Y}} P_X(x) p_{y_1}(x) p_{y_2}(x) \bar{q}_y(y_1) \bar{q}_y(y_2).$$

We have then concluded that  $|h(y') - s(y')| \leq 3\kappa/M$  for each  $y' \in \mathcal{Y}$ . We also have that  $s(y) \stackrel{(a)}{\geq} \left( \sum_{x \in \mathcal{X}, y' \in \mathcal{Y}} P_X(x) p_{y'}(x) \bar{q}_y \right)^2 \geq \eta^2$ , where (a) follows from Jensen's inequality applied to the function  $f(x) = x^2$ . Now by the multiplicative Hoeffding's inequality, for any  $\epsilon > 0$ , we have,

$$P \left( \sum_{i' \in \mathcal{N} \setminus \mathcal{W}_j} \mathbb{1}_{r_{J_1(i')}=y'} \mathbb{1}_{r_{J_2(i')}=y'} \geq (N - |\mathcal{W}_j|) h(y') (1 + \epsilon) \right) \leq \exp(-\epsilon^2 h(y')/3), \text{ and,} \quad (47)$$

$$P \left( \sum_{i' \in \mathcal{N} \setminus \mathcal{W}_j} \mathbb{1}_{r_{J_1(i')}=y'} \mathbb{1}_{r_{J_2(i')}=y'} \leq (N - |\mathcal{W}_j|) h(y') (1 - \epsilon) \right) \leq \exp(-\epsilon^2 h(y')/3). \quad (48)$$

Thus, for any  $\epsilon > 0$ , and  $N$  large enough, we have,

$$\begin{aligned}
& \mathbb{E} \left( \frac{\sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + \sum_{i' \in \mathcal{N} \setminus \mathcal{W}_j} \mathbb{1}_{r_{J_1(i')} = y'} \mathbb{1}_{r_{J_2(i')} = y'}}}} \right) \\
& \stackrel{(a)}{\leq} \mathbb{E} \left( \frac{\sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + (N - |\mathcal{W}_j|)h(y')(1 - \epsilon)}} + \exp(-\epsilon^2 h(y')(N - |\mathcal{W}_j|)/3) \sqrt{N - |\mathcal{W}_j|} \right) \\
& \stackrel{(b)}{\leq} \mathbb{E} \left( \frac{\sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + (N - |\mathcal{W}_j|)(s(y') - \frac{3\kappa}{M})(1 - \epsilon)}} + \exp(-\frac{\epsilon^2 (s(y') - \frac{3\kappa}{M})(N - |\mathcal{W}_j|)}{3}) \sqrt{N - |\mathcal{W}_j|} \right) \\
& \stackrel{(c)}{\leq} \mathbb{E} \left( \frac{\sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + (N - |\mathcal{W}_j|)(s(y') - \frac{3\kappa n}{mN})(1 - \epsilon)}} \right) + \mathbb{E} \left( \exp(-\frac{\epsilon^2 (\eta^2 - \frac{3\kappa n}{mN})(N - |\mathcal{W}_j|)}{3}) \sqrt{N - |\mathcal{W}_j|} \right) \\
& \stackrel{(d)}{\leq} \frac{1}{\sqrt{(s(y') - \frac{3\kappa n}{mN})(1 - \epsilon)}} + \mathbb{E} \left( \exp(-\frac{\epsilon^2 (\eta^2 - \frac{3\kappa n}{mN})(N - |\mathcal{W}_j|)}{3}) \right) \sqrt{N} \\
& = \frac{1}{\sqrt{(s(y') - \frac{3\kappa n}{mN})(1 - \epsilon)}} + \sqrt{N} \exp(-\frac{\epsilon^2 (\eta^2 - \frac{3\kappa n}{mN})N}{3}) \mathbb{E} \left( \exp(\frac{\epsilon^2 (\eta^2 - \frac{3\kappa n}{mN})|\mathcal{W}_j|}{3}) \right). \tag{49}
\end{aligned}$$

Here, (a) results from Equation 48 and the fact that in the worst case, the left hand side is at most  $\sqrt{N - |\mathcal{W}_j|}$ . (b) results from the fact that  $|h(y') - s(y')| \leq 3\kappa/M$ , and (c) results from the fact that  $s(y) \geq \eta^2$  and  $M > mN/n$ . All the expectations are with respect to the randomness in  $|\mathcal{W}_j|$ . (d) follows from (i) noting that  $s(y') \geq \eta^2 > \frac{3\kappa n}{mN}$  for  $N$  large enough, (ii) ignoring the constant 1 in the denominator of the first term, and (iii) ignoring  $|\mathcal{W}_j| \geq 0$  in the second term.

Now, due to the randomized task allocation policy,  $|\mathcal{W}_j|$  is distributed as Binomial( $N, m/N$ ). Since the moment generating function of a Binomially distributed random variable  $X$  with parameters  $(n, p)$  is  $\mathbb{E}(\exp(Xt)) = (1 - p + pe^t)^n$ , we have that

$$\mathbb{E} \left( \exp(\frac{\epsilon^2 (\eta^2 - \frac{3\kappa n}{mN})|\mathcal{W}_j|}{3}) \right) = \left( 1 - \frac{m}{N} + \frac{m}{N} \exp(\frac{\epsilon^2 (\eta^2 - \frac{3\kappa n}{mN})}{3}) \right)^N. \tag{50}$$

Choosing  $\epsilon = N^{-1/4}$ , we have that

$$\lim_{N \rightarrow \infty} \left( 1 - \frac{m}{N} + \frac{m}{N} \exp(\frac{\eta^2 - \frac{3\kappa n}{mN}}{3\sqrt{N}}) \right)^N = \lim_{N \rightarrow \infty} \left( 1 + \frac{m}{N} (\exp(\frac{\eta^2 - \frac{3\kappa n}{mN}}{3\sqrt{N}}) - 1) \right)^N \leq \lim_{N \rightarrow \infty} \left( 1 + \frac{m}{N} \right)^N = \exp(m). \tag{51}$$

Thus, choosing  $\epsilon = N^{-1/4}$  in Equation 49, and combining Equation 51 with the fact that

$$\lim_{N \rightarrow \infty} \sqrt{N} \exp(-\frac{(\eta^2 - \frac{3\kappa n}{mN})\sqrt{N}}{3}) = 0,$$

we have that

$$\lim_{N \rightarrow \infty} \mathbb{E} \left( \frac{\sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + \sum_{i' \in \mathcal{N} \setminus \mathcal{W}_j} \mathbb{1}_{r_{J_1(i')} = y'} \mathbb{1}_{r_{J_2(i')} = y'}}}} \right) \leq \frac{1}{\sqrt{s(y')}}. \tag{52}$$

Next, we also have that, for a large enough  $N$ ,

$$\mathbb{E} \left( \frac{\sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + \sum_{i' \in \mathcal{N} \setminus \mathcal{W}_j} \mathbb{1}_{r_{J_1(i')} = y'} \mathbb{1}_{r_{J_2(i')} = y'}}}} \right)$$

$$\begin{aligned}
&\stackrel{(a)}{\geq} \mathbb{E} \left( \frac{\sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + (N - |\mathcal{W}_j|)h(y')(1 + \epsilon)}} (1 - \exp(-\epsilon^2 h(y')(N - |\mathcal{W}_j|)/3)) \right) \\
&\stackrel{(b)}{\geq} \mathbb{E} \left( \frac{\sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + N(s(y') + \frac{3\kappa}{M})(1 + \epsilon)}} - \exp\left(-\frac{\epsilon^2(s(y') - 3\kappa/M)(N - |\mathcal{W}_j|)}{3}\right) \frac{\sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + N(s(y') - \frac{3\kappa}{M})(1 + \epsilon)}} \right) \\
&\stackrel{(c)}{\geq} \mathbb{E} \left( \frac{\sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + N(s(y') + \frac{3\kappa n}{mN})(1 + \epsilon)}} \right) - \mathbb{E} \left( \exp\left(-\frac{\epsilon^2(\eta^2 - \frac{3\kappa n}{mN})(N - |\mathcal{W}_j|)}{3}\right) \frac{\sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + N(\eta^2 - \frac{3\kappa n}{mN})(1 + \epsilon)}} \right) \\
&\geq \frac{\mathbb{E} \left( \sqrt{1 - \frac{|\mathcal{W}_j|}{N}} \right)}{\sqrt{1/N + (s(y') + \frac{3\kappa n}{mN})(1 + \epsilon)}} - \mathbb{E} \left( \exp\left(\frac{\epsilon^2(\eta^2 - \frac{3\kappa n}{mN})(|\mathcal{W}_j|)}{3}\right) \right) \exp\left(-\frac{\epsilon^2(\eta^2 - \frac{3\kappa n}{mN})N}{3}\right) \frac{\sqrt{N}}{\sqrt{1 + N(\eta^2 - \frac{3\kappa n}{mN})(1 + \epsilon)}} \\
&\stackrel{(d)}{\geq} \frac{\mathbb{E} \left( \sqrt{1 - \frac{|\mathcal{W}_j|}{N}} \right)}{\sqrt{1/N + (s(y') + \frac{3\kappa n}{mN})(1 + \epsilon)}} - \left(1 - \frac{m}{N} + \frac{m}{N} \exp\left(\frac{\epsilon^2(\eta^2 - \frac{3\kappa n}{mN})}{3}\right)\right)^N \exp\left(-\frac{\epsilon^2(\eta^2 - \frac{3\kappa n}{mN})N}{3}\right) \frac{\sqrt{N}}{\sqrt{1 + N(\eta^2 - \frac{3\kappa n}{mN})(1 + \epsilon)}}. \tag{53}
\end{aligned}$$

Here, (a) follows from Equation 47, and the fact that the agreement rewards are always positive. (b) results from the fact that  $|h(y') - s(y')| \leq 3\kappa/M$  and by ignoring the  $|\mathcal{W}_j|$  term in the denominator, and (c) results from the fact that  $s(y) \geq \eta^2$  and  $M > mN/n$ . (d) follows from Equation 50. We once again choose  $\epsilon = N^{-1/4}$ . Then, by Equation 51, the second term in Equation 53 converges to 0 as  $N \rightarrow \infty$ . We now focus on the first term. The denominator of this term clearly converges to  $\sqrt{s(y')}$ . It is now easy to show that the numerator converges to 1. This is because  $|\mathcal{W}_j|$  is distributed as Binomial( $N, m/N$ ), and thus  $|\mathcal{W}_j|/N$  converges in distribution to the constant 0. Since  $f(x) = \sqrt{1-x}$  is a bounded, continuous function on the domain  $[0, 1]$ , it follows (by the Portmanteau's theorem on the equivalence of definitions of convergence in distribution) that  $\mathbb{E} \left( \sqrt{1 - |\mathcal{W}_j|/N} \right)$  converges to 1 as  $N \rightarrow \infty$ . Thus, we finally have,

$$\lim_{N \rightarrow \infty} \mathbb{E} \left( \frac{\sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + \sum_{i' \in \mathcal{N} \setminus \mathcal{W}_j} \mathbb{1}_{r_{J_1(i')} = y'} \mathbb{1}_{r_{J_2(i')} = y'}}} \right) \geq \frac{1}{\sqrt{s(y')}}. \tag{54}$$

Thus, from Equations 52 and 54, we finally have,

$$\left| \mathbb{E} \left( \frac{K \sqrt{N - |\mathcal{W}_j|}}{\sqrt{1 + \sum_{i' \in \mathcal{N} \setminus \mathcal{W}_j} \mathbb{1}_{r_{J_1(i')} = y'} \mathbb{1}_{r_{J_2(i')} = y'}}} \right) - \frac{K}{\sqrt{s(y')}} \right| \leq \sigma(N). \tag{55}$$

where  $\sigma(N) = o(1)$ . Now before we proceed, note that the convergence of the expected matching reward for answer  $y$  to  $K/\sqrt{s(y)}$  for each  $y \in \mathcal{Y}$  is all that is required for strict truthfulness to follow for a large enough  $N$ , as we show in the proof of Theorem 1. For the truthful strategy profile, by the  $\alpha$ -separation assumption, we have that  $s(y) \geq \alpha^2$  for all  $y \in \mathcal{Y}$ . Thus, by replacing  $\eta$  with  $\alpha$  in the arguments leading up to Equation 55, we can conclude the convergence of the matching rewards to  $K/\sqrt{s(y)}$  for each  $y \in \mathcal{Y}$ . Thus, there exists  $N_1$  such that for  $N \geq N_1$ , the truthful strategy profile is a Bayes-Nash equilibrium. We will not repeat the proof here for conciseness. The first statement of the theorem thus follows and we hence focus on proving the second and third statement.

To that effect, we now proceed to focus on term B in Equation 43. We have

$$\begin{aligned}
P(r_J^i = y', Y_j^i = y) &\stackrel{(a)}{=} \sum_{x \in \mathcal{X}, y_2 \in \mathcal{Y}} P_X(x) p_y(x) p_{y_2}(x) \bar{q}_{y'}^{\{j\}}(y_2) \\
&\stackrel{(b)}{\leq} \sum_{x \in \mathcal{X}, y_2 \in \mathcal{Y}} P_X(x) p_y(x) p_{y_2}(x) (\bar{q}_{y'}(y_2) + \kappa/M) \\
&\leq \sum_{x \in \mathcal{X}, y_2 \in \mathcal{Y}} P_X(x) p_y(x) p_{y_2}(x) \bar{q}_{y'}(y_2) + \kappa/M.
\end{aligned} \tag{56}$$

Here, (a) follows from the fact that  $J$  is equally likely to be any agent other than  $j$ , by the random task allocation policy. (b) again follows from Equation 41. Similarly, we have

$$P(r_J^i = y', Y_j^i = y) \geq \sum_{x \in \mathcal{X}, y_2 \in \mathcal{Y}} P_X(x) p_y(x) p_{y_2}(x) \bar{q}_{y'}(y_2) - \kappa/M. \tag{57}$$

Thus, we have, for all  $y'$  such that the population average probability of reporting is at least  $\eta$ , we have

$$\begin{aligned}
P(Y_j^i = y) G_i^j(y, y') &= \left( \frac{K}{\sqrt{s(y')}} + o(1) \right) \left( \sum_{x \in \mathcal{X}, y_2 \in \mathcal{Y}} P_X(x) p_y(x) p_{y_2}(x) \bar{q}_{y'}(y_2) + o(1) \right) \\
&= \frac{K \left( \sum_{x \in \mathcal{X}, y_2 \in \mathcal{Y}} P_X(x) p_y(x) p_{y_2}(x) \bar{q}_{y'}(y_2) \right)}{\sqrt{s(y')}} + o(1).
\end{aligned} \tag{58}$$

Here, we use the fact that  $s(y) \geq \eta^2$ . For every other  $y'$  such that the population average probability of reporting is 0, we have

$$P(Y_j^i = y) G_i^j(y, y') = 0. \tag{59}$$

Let  $\mathcal{Y}'$  denote the set of responses such that the population average probability of reporting the response is at least  $\eta$ . Then the expected payoff of agent  $j$  on task  $i$  under policy  $\mathbf{q}^j$  (fixing everyone else's policy) is:

$$\mathbf{G}_i^j(\mathbf{q}^j) = \sum_{y, y' \in \mathcal{Y}} P(Y_j^i = y) G_i^j(y, y') q_{y'}^j(y) \tag{60}$$

$$\stackrel{(a)}{=} \sum_{y \in \mathcal{Y}, y' \in \mathcal{Y}'} \frac{K \left( \sum_{x \in \mathcal{X}, y_2 \in \mathcal{Y}} P_X(x) p_y(x) p_{y_2}(x) \bar{q}_{y'}(y_2) \right)}{\sqrt{s(y')}} q_{y'}^j(y) + o(1) \tag{61}$$

$$\stackrel{(b)}{=} \max_{\mathbf{q}} \sum_{y, y' \in \mathcal{Y}} P(Y_j^i = y) G_i^j(y, y') q_{y'}^j(y) \tag{62}$$

$$\stackrel{(c)}{=} \max_{\mathbf{q}} \sum_{y \in \mathcal{Y}, y' \in \mathcal{Y}'} \frac{K \left( \sum_{x \in \mathcal{X}, y_2 \in \mathcal{Y}} P_X(x) p_y(x) p_{y_2}(x) \bar{q}_{y'}(y_2) \right)}{\sqrt{s(y')}} q_{y'}^j(y) + o(1). \tag{63}$$

Here, (a) follows from Equation 58, (b) follows from the fact that  $\mathbf{q}^j$  is a best-response strategy, and (c) again follows from Equation 58. Note that the final right hand side neither depends on the identity of agent  $j$  nor does it depend on the identity of task  $i$ . It only depends on the population average strategy  $\bar{\mathbf{q}}$ . Since, each policy  $\mathbf{q}^{j'}$  optimizes  $\mathbf{G}_i^{j'}(\mathbf{q})$ , we have that

$$\mathbf{G}_i^j(\bar{\mathbf{q}}) = \sum_{y, y' \in \mathcal{Y}} P(Y_j^i = y) G_i^j(y, y') \bar{q}_{y'}(y) \tag{64}$$

$$\stackrel{(a)}{=} \sum_{y \in \mathcal{Y}, y' \in \mathcal{Y}'} \frac{K \left( \sum_{x \in \mathcal{X}, y_2 \in \mathcal{Y}} P_X(x) p_y(x) p_{y_2}(x) \bar{q}_{y'}(y_2) \right)}{\sqrt{s(y')}} \bar{q}_{y'}(y) + o(1) \tag{65}$$

$$\stackrel{(b)}{=} \max_{\mathbf{q}} \sum_{y \in \mathcal{Y}, y' \in \mathcal{Y}'} \frac{K \left( \sum_{x \in \mathcal{X}, y_2 \in \mathcal{Y}} P_X(x) p_y(x) p_{y_2}(x) \bar{q}_{y'}(y_2) \right)}{\sqrt{s(y')}} q_{y'}(y) + o(1) \tag{66}$$

$$= \mathbf{G}_i^j(\mathbf{q}^j) + o(1). \tag{67}$$

Here, (a) results from Equation 58. (b) results from averaging the expression in Equation 61 and the expression in Equation 63 across all agents and realizing that the expression in Equation 63 is identical across the agents. Hence, we have that  $|\mathbf{G}_i^j(\mathbf{q}^j) - \mathbf{G}_i^j(\bar{\mathbf{q}})| = o(1)$ . Hence, we finally have,

$$\mathbf{G}_i^j(\mathbf{q}^j) = \mathbf{G}_i^j(\bar{\mathbf{q}}) + o(1) \quad (68)$$

$$= \sum_{y \in \mathcal{Y}, y' \in \mathcal{Y}'} \frac{K(\sum_{x \in \mathcal{X}, y_2 \in \mathcal{Y}} P_X(x) p_y(x) p_{y_2}(x) \bar{q}_{y'}(y_2))}{\sqrt{s(y')}} \bar{q}_{y'}(y) + o(1) \quad (69)$$

$$= \sum_{y' \in \mathcal{Y}'} \frac{K \sum_{x \in \mathcal{X}, y_2 \in \mathcal{Y}, y \in \mathcal{Y}} P_X(x) p_y(x) p_{y_2}(x) \bar{q}_{y'}(y_2) \bar{q}_{y'}(y)}{\sqrt{s(y')}} + o(1) \quad (70)$$

$$= \sum_{y' \in \mathcal{Y}'} \frac{K s(y')}{\sqrt{s(y')}} + o(1) \quad (71)$$

$$\stackrel{(a)}{=} \sum_{y' \in \mathcal{Y}} K \sqrt{s(y')} + o(1) \quad (72)$$

$$\stackrel{(b)}{\leq} K\Gamma(Y_j^i, Y_{j'}^i) - \frac{K\delta(P_X, \mathbf{p})\Omega(\bar{\mathbf{q}})^2(|\mathcal{Y}| - 1)}{2\sqrt{|\mathcal{Y}|}} + o(1) \quad (73)$$

$$\stackrel{(c)}{=} \mathbf{G}(\mathbf{t}) - \frac{K\delta(P_X, \mathbf{p})\Omega(\bar{\mathbf{q}})^2(|\mathcal{Y}| - 1)}{2\sqrt{|\mathcal{Y}|}} + o(1) \quad (74)$$

where  $\mathbf{G}(\mathbf{t})$  is the expected reward to each agent for an evaluation task under the truthful strategy profile. Here, (a) follows from the fact that  $s(y) = 0$  for all  $y \in \mathcal{Y} \setminus \mathcal{Y}'$ . (b) follows from Proposition 5.1, and (c) follows from the fact that  $\mathbf{G}(\mathbf{t}) = K\Gamma(Y_j^i, Y_{j'}^i) + o(1)$ . This latter conclusion results from the fact that under the truthful strategy profile, because of  $\alpha$ -separation of the generating model, the probability of reporting any answer  $y \in \mathcal{Y}$  is at least  $\alpha^2$  (see proof of Proposition 4.1). We can thus use the same arguments as that used for deriving the expression in Equation 72 as the expected payoff of each agent, while replacing  $\eta$  with  $\alpha$ .

Now the second statement of the theorem immediately follows, since for any  $\omega > 0$ , there is an  $N_2$  such that for any  $N > N_2$  we have that  $\mathbf{G}_i^j(\mathbf{q}^j) \leq \mathbf{G}(\mathbf{t}) + \omega$ . Moreover, we have that if  $\mathbf{G}_i^j(\mathbf{q}^j) > \mathbf{G}(\mathbf{t})$ , then, Equation 74 allows us to conclude that

$$\frac{K\delta(P_X, \mathbf{p})\Omega(\bar{\mathbf{q}})^2(|\mathcal{Y}| - 1)}{2\sqrt{|\mathcal{Y}|}} \leq o(1), \text{ or,} \quad (75)$$

$$\Omega(\bar{\mathbf{q}}) \leq \sqrt{\frac{2\sqrt{|\mathcal{Y}|}o(1)}{K\delta(P_X, \mathbf{p})(|\mathcal{Y}| - 1)}} \leq \sqrt{\frac{2\sqrt{|\mathcal{Y}|}o(1)}{K\alpha(|\mathcal{Y}| - 1)}}. \quad (76)$$

Thus for any  $\omega > 0$ , there is an  $N_3$  such that for any  $N > N_3$ , we have that  $\Omega(\bar{\mathbf{q}}) \leq \omega$  for any population strategy profile where (a) the average probability of reporting any answer  $y$  is either 0 or at least  $\eta$ , and (b) there exists an agent whose expected payoff is larger than the expected payoff under the truthful strategy profile. Thus, all the statements of the theorem hold for any  $N$  larger than  $N_0 = \max(N_1, N_2, N_3)$ .  $\square$

## B. Properties of the square-root agreement measure

The SRAM 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  $y \in \mathcal{Y}$ .

We also prove the following inequality satisfied by the SRAM, which generalizes Proposition 5.1 without the characterizing the inequality gap.

**Proposition B.1 (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,*

$$\sum_{y \in \mathcal{Y}} \sqrt{P(Z_1 = Z_2 = y)} \leq \Gamma(Y_1, Y_2). \quad (77)$$

Moreover, if  $\delta(P_X, \mathbf{p}) > 0$ , then the above inequality is an equality if and only if  $\mathbf{q} = \mathbf{q}'$  and  $\Omega(\mathbf{q}) = 0$ , i.e., if and only if the two reporting strategies are identical and fully informative.

*Proof of Proposition B.1* 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.

Now suppose that  $\delta(P_X, \mathbf{p}) > 0$ . Then (a) is an equality if and only if  $q_y(y_1)q'_y(y_2) = 0$  for every  $y$  and every  $y_1 \neq y_2$ . Further, (b) is an equality, i.e., arithmetic mean equals geometric mean, if and only if all the terms are equal. This means for all  $y \in \mathcal{Y}$ ,

$$\sum_{y_1 \in \mathcal{Y}} q_y(y_1) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_1}(x)^2} = \sum_{y_2 \in \mathcal{Y}} q'_y(y_2) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y_2}(x)^2},$$

i.e., if

$$\sum_{y' \in \mathcal{Y}} (q_y(y') - q'_y(y')) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y'}(x)^2} = 0. \quad (78)$$

Squaring both sides, we obtain, for all  $y \in \mathcal{Y}$ ,

$$\begin{aligned} & \sum_{y' \in \mathcal{Y}} (q_y(y') - q'_y(y'))^2 \left( \sum_{x \in \mathcal{X}} P_X(x) p_{y'}(x)^2 \right) \\ & + \sum_{\substack{y' \neq y'' \\ y' \in \mathcal{Y}}} (q_y(y') - q'_y(y')) (q_y(y'') - q'_y(y'')) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y'}(x)^2} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y''}(x)^2} = 0. \end{aligned} \quad (79)$$

Substituting  $q_y(y')q'_y(y'') = 0$  for all  $y' \neq y''$ , we obtain,

$$\begin{aligned} & \sum_{y' \in \mathcal{Y}} (q_y(y') - q'_y(y'))^2 \left( \sum_{x \in \mathcal{X}} P_X(x) p_{y'}(x)^2 \right) \\ & + \sum_{\substack{y' \neq y'' \\ y' \in \mathcal{Y}}} (q_y(y')q_y(y'') + q'_y(y')q'_y(y'')) \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y'}(x)^2} \sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y''}(x)^2} = 0. \end{aligned} \quad (80)$$

But if  $\delta(P_X, \mathbf{p}) > 0$ , then we know from Proposition 4.1 that  $\sqrt{\sum_{x \in \mathcal{X}} P_X(x) p_{y'}(x)^2} \geq \sum_{x \in \mathcal{X}} P_X(x) p_{y'}(x) > 0$  for all  $y' \in \mathcal{Y}$ . Hence, we conclude that if  $\delta(P_X, \mathbf{p}) > 0$ , then  $\sum_{y \in \mathcal{Y}} \sqrt{P(Z_1 = Z_2 = y)} = \Gamma(Y_1, Y_2)$  holds, if and only if all the terms in Equation 80 are 0, i.e., if and only if

1.  $q_y(y') = q'_y(y')$  for all  $y, y' \in \mathcal{Y}$ , i.e.,  $q$  and  $q'$  are identical, and,
2.  $q_y(y')q_y(y'') = 0$  for all  $y \in \mathcal{Y}$  and  $y' \neq y''$ , i.e.,  $\Omega(q) = 0$ .

This finishes the proof.  $\square$

### B.1. Utility of the square-root agreement measure beyond our work

Definition 6 essentially defines an agreement measure between any two random variables that are independent and identically distributed conditioned on some latent random variable. But we could just as well define an agreement measure between any two random variables that take values in some common finite set.

DEFINITION 9. Consider two random variables  $X$  and  $X'$ , which take values in a finite set  $\mathcal{S}$ . Then the square-root agreement measure between  $X$  and  $X'$  is defined as

$$\Gamma(X, X') = \sum_{s \in \mathcal{S}} \sqrt{P(X = X' = s)}.$$

Proposition B.1 implies that if  $X \rightarrow X' \rightarrow Y$  form a Markov chain, i.e.,  $X$  is conditionally independent of  $Y$  given  $X'$ , and if, conditioned on some latent random variable  $U$ ,  $X$  and  $X'$  are independent and identically distributed random variables, then,

$$\Gamma(X, Y) \leq \Gamma(X, X').$$

Inequalities of this form are called *data processing* inequalities and they have several applications in information theory, statistics, causal inference, and related fields. For example, such inequalities provide testable hypotheses to determine the validity of conditional independence assumptions across variables from data. Several *mutual information* measures between two random variables are known to satisfy such inequality. These measures are typically constructed from two classes of divergences or distance notions between probability distributions, called f-divergences and Bregman divergences; see Kong and Schoenebeck (2019) and references therein. It is interesting to note that our SRAM does not result from such a construction, and to the best of our knowledge, the resulting data-processing inequality was not known in the literature. Moreover, typical mutual information measures depend on the entire joint distribution of two variables, i.e., to estimate these measures from data, one typically needs to learn  $|\mathcal{S}|^2$  probability values where  $\mathcal{S}$  is the support set of each variable. On the other hand, the SRAM only depends on the diagonal values of the joint probability distribution, i.e., only the probabilities of agreement matter. Hence, to estimate the SRAM from data, one only needs to learn  $|\mathcal{S}|$  probability values. It is important to note that for the data processing inequality to hold for the SRAM,  $X$  and  $X'$  need to be conditionally independent and identically distributed (conditioned on some latent random variable). There is typically no such requirement for other measures. To show that this condition is necessary, consider the following counterexample. Suppose that  $X$  is uniformly distributed on the discrete set  $\{-1, +1\}$ , and  $X' = -X$ . Thus  $\Gamma(X, X') = 0$ . Whereas if  $Y = -X'$ , then it is true that  $X \rightarrow X' \rightarrow Y$  forms a Markov chain, and  $\Gamma(X, Y) = \Gamma(X, X) = \sqrt{2}$ . Hence,  $\Gamma(X, X') < \Gamma(X, Y)$ .

### C. Properties of the uninformativeness measure

The uninformativeness measure has the following properties.

1. 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.
2.  $\Omega(\mathbf{q})$  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 independently of the true answer. To see this, observe that,

$$\begin{aligned}
& \frac{1}{|\mathcal{Y}|(|\mathcal{Y}| - 1)} \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} \sqrt{q_y(y') q_y(y'') \mathbf{1}_{\{y' \neq y''\}}} \\
& \stackrel{(a)}{\leq} \frac{1}{|\mathcal{Y}|(|\mathcal{Y}| - 1)} \sum_{y \in \mathcal{Y}} \sqrt{\left( \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} q_y(y') \mathbf{1}_{\{y' \neq y''\}} \right) \left( \sum_{y' \in \mathcal{Y}, y'' \in \mathcal{Y}} q_y(y'') \mathbf{1}_{\{y' \neq y''\}} \right)} \\
& = \frac{1}{|\mathcal{Y}|(|\mathcal{Y}| - 1)} \sum_{y \in \mathcal{Y}} \sqrt{(|\mathcal{Y}| - 1)^2 \left( \sum_{y' \in \mathcal{Y}} q_y(y') \right)^2} \\
& = \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}} q_y(y') \\
& = 1.
\end{aligned} \tag{81}$$

Here, (a) follows from the Cauchy-Schwarz inequality.

## D. Miscellaneous remarks on existing mechanisms

### D.1. An adaptation of the Correlated Agreement (CA) mechanism to the homogeneous responses setting

In this section, we present an adaptation of CA to the homogeneous responses setting, which induces truthful behavior while only requiring one evaluation per agent. In order to define this mechanism, we first present the original CA mechanism.

**Original CA mechanism.** CA operates on a pair of agents. Both agents perform a common “bonus” evaluation task, say A, and individually perform one independent “penalty” evaluation task that the other agent doesn’t perform; say agent 1 performs B, and agent 2 performs C. In keeping with our notation, let  $r_j^i$  for  $j = 1, 2$ , and  $i \in \{A, B, C\}$ , be the response of agent  $j$  in task  $i$ , where the responses are taken to be the null  $\phi$  if an agent doesn’t perform the corresponding task (hence  $r_1^C = r_2^B = \phi$ ). CA defines an intermediate scoring function that maps two responses to a real number, which, informally, is monotonically increasing in the expected correlation between the responses, i.e.,  $S$  is higher if the two responses are expected to frequently occur together. Formally, denoting  $\Delta_{ab} = P(Y_j^i = a, Y_{j'}^i = b) - P(Y_j^i = a)P(Y_{j'}^i = b)$ , for any two responses  $a, b$  (where it is assumed that agents  $j$  and  $j'$  have both performed task  $i$ ), the intermediate score for these responses is defined to be:

$$S(a, b) = \text{sgn}(\Delta_{ab}),$$

where  $\text{sgn}$  is the sign function. In the multi-task, homogeneous responses setting, this scoring function can be estimated from the response data obtained from the large number of other participants (i.e., excluding the two agents under consideration) operating on the platform, and relying on the “self-fulfilling prophecy of truthfulness” to assume truthful behavior.<sup>19</sup> The final score/payment to an agent  $i$  is then defined to be

$$S(r_1^A, r_2^A) - S(r_1^B, r_2^C),$$

i.e., the final score/payment is the difference between the bonus score and the penalty score. The intuition is that the payment scheme rewards *incremental* correlation in the responses to the bonus task over what is expected anyway from the responses to two independent evaluation tasks.

**An adaptation of CA for homogeneous responses requiring one evaluation per agent.** Consider the following adaptation of CA. Consider an agent, say 1, who has performed evaluation task A, and whose payment needs to be determined. Let 2 be another agent who has performed task A. Let 3 be a third agent who has performed some task B that 1 hasn’t performed. Then the payment to agent 1 is defined to be:

$$S(r_1^A, r_2^A) - S(r_1^A, r_3^B).$$

Here, it is assumed that these scoring functions are calculated as in the original CA mechanism based on the responses data from all tasks other than A and B. In a natural practical implementation of this mechanism, to calculate the payment of an agent  $j$ , the platform would randomly pick a peer agent who has performed the

<sup>19</sup> In the general setting with non-homogeneous responses, this function can be estimated by having the two agents perform a large number of overlapping and disjoint tasks.

same evaluation to calculate the bonus score and similarly, randomly pick another agent who has performed some task that the agent  $j$  hasn't performed to calculate the penalty score. These scoring functions are calculated from the response data of all tasks that  $j$  hasn't performed. We call this mechanism CA for homogeneous responses (CA-HR).

It is clear that this mechanism only requires each agent to perform one evaluation as long as (a) each task is performed by at least two agents, and (b) there is a large number of tasks while each agent performs only a small number (so that the scores can be estimated accurately, independently of the agent's reports). These assumptions are almost the same as that required by SRA for its properties, and they are easily satisfied on most platforms.

It is easy to argue that truthful behavior is an equilibrium under CA-HR (in the large tasks regime where the scoring function estimates are reasonably accurate) in the homogeneous responses setting. This is because, due to the statistical indistinguishability of agent 2 and 3's responses to an arbitrary task assuming that they are truthful, replacing agent 2 with agent 3 in the calculation of the penalty score is inconsequential from the perspective of agent 1. All that matters from the perspective of agent 1 (in terms of aligning with the incentives generated by the original CA mechanism) is that this penalty score is computed on the basis of some agent's response to a task that 1 hasn't performed. Thus, the fact that truthful behavior is a best response under the original CA mechanism implies that it is a best response under this modification as well.

## D.2. Remarks on the properties of CA/CA-HR in our setting.

Although CA is informed truthful in the setting in which it is originally defined, neither CA nor CA-HR are informed truthful in our setting. This is because in our setting, task allocations are exogenously specified and agents can choose task-contingent reporting strategies based on task identities. We present an example below that shows this for CA-HR.

EXAMPLE 5. Consider the setting in Example 1 again. For the sake of the present discussion, suppose that the plumbers are numbered  $i = 1, \dots, N$  (just as the tasks are numbered in our formal model). If everyone is truthful, the accurate scoring function is  $S(a, b) = \mathbf{1}_{\{a=b\}} - \mathbf{1}_{\{a \neq b\}}$ . Suppose that  $j$  has evaluated plumber  $A$ .  $j'$  is her randomly chosen peer who also has also evaluated  $A$ . Let  $j''$  be another randomly chosen peer who has evaluated plumber  $B$ , whom  $j$  hasn't evaluated. Then the (random) payment of agent  $j$  under CA-HR is

$$\mathbf{1}_{\{Y_j^A = Y_{j'}^A\}} - \mathbf{1}_{\{Y_j^A \neq Y_{j'}^A\}} - \mathbf{1}_{\{Y_j^A = Y_{j''}^B\}} + \mathbf{1}_{\{Y_j^A \neq Y_{j''}^B\}}.$$

Thus, the expected payment of agent  $j$  can be determined to be

$$2(P(Y_j^A = Y_{j'}^A) - P(Y_j^A = Y_{j''}^B))$$

This can be computed to be 0.2025, given the generating model. On the other hand, consider the following strategy profile. For all even tasks, agents report 'Yes,' and for all odd tasks, agents report 'No.' We first argue that this strategy profile is an equilibrium in the many tasks regime. Note that under this strategy profile, the scoring function that will be estimated by the platform is  $\bar{S}(a, b) = \mathbf{1}_{\{a=b\}} - \mathbf{1}_{\{a \neq b\}}$ , same as that under truthful behavior. Thus, the expected payment of an agent for reporting 'No' on an even task (or reporting 'Yes' on an odd task) is  $-1$ , whereas the expected payment from following the prescribed strategy

is 1. This argument shows both, that (a) this strategy profile constitutes an equilibrium and (b) the expected payoff to any agent under this strategy profile (which is 1) is strictly higher than the expected payment under the truthful equilibrium (which is 0.2025). Thus CA-HR is not informed truthful.

The same construction of a non-truthful strategy profile also shows that CA is also not informed truthful in our setting.

Moreover, unlike the mechanism of Kong and Schoenebeck (2019) (see Section 3), in our setting, neither CA nor CA-HR are informed truthful across all equilibria where agents choose the same reporting strategy for each task they perform. This is because task allocations are exogenously specified: in the example above, it could very well be the case that every agent performs exactly one task under CA-HR. In this case, the non-truthful equilibrium strategy profile constructed above respects the constraint that each agent chooses the same reporting strategy for each task they perform, simply because each agent performs only one task. A similar argument shows this for CA by considering a situation in which each agent performs either even tasks only or odd tasks only.

Although CA and CA-HR are not informed truthful in our setting, Lemma 5.12 in Shnayder et al. (2016) implies that these mechanisms are informed truthful across symmetric equilibria, i.e., they are informed truthful when restricted to symmetric equilibria, in the many tasks limit.

Next, we discuss why CA and CA-HR are *not* (asymptotically) strongly truthful across symmetric equilibria in general for homogeneous responses, i.e., there could be symmetric strategy profiles that are *not* fully informative, that asymptotically yield the same payoff as the truthful equilibrium. The existence of such strategy profiles is related to the following notion of “clustered observations” as defined in Shnayder et al. (2016).

**DEFINITION 10.** (Shnayder et al. 2016) A distribution of two agents’ observations for a common evaluation is said to be clustered if there exist at least two identical rows in the matrix  $[\text{sgn}(\Delta_{yy'})]_{y \in \mathcal{Y}, y' \in \mathcal{Y}}$ . (Note that  $[\text{sgn}(\Delta_{yy'})]_{y \in \mathcal{Y}, y' \in \mathcal{Y}}$  is a symmetric matrix under homogeneous responses)

In the presence of clustered observations, there are symmetric equilibrium strategy profiles that are not fully informative, that yield the same payoff asymptotically as the truthful equilibrium under CA/CA-HR. To see this for CA-HR, suppose that  $y$  and  $\bar{y}$  are two observations for which the corresponding rows  $(\text{sgn}(\Delta_{yy'}); y' \in \mathcal{Y})$  and  $(\text{sgn}(\Delta_{\bar{y}y'}); y' \in \mathcal{Y})$  are identical. Then, if all agents report a fixed observation, e.g.,  $y$ , irrespective of whether they observe  $y$  or  $\bar{y}$ , the scoring function estimated by the platform under CA-HR is the same as that under truthful behavior, except with the answer  $\bar{y}$  eliminated as a possible report. However, if everyone else was truthful, the bonus and penalty scores obtained by an agent would have anyway been identical irrespective of whether any of the three agents involved in computing a payment report  $y$  or  $\bar{y}$ . Thus the payments to all agents remain the same if everyone reports  $y$  irrespective of whether they observe  $y$  or  $\bar{y}$ . It thus follows that this strategy profile is an equilibrium under CA-HR, which yields the same expected payoff to any agent as the truthful equilibrium. This strategy profile is not a fully informative strategy profile, and hence, CA-HR is not asymptotically strongly truthful across symmetric equilibria. The mechanism is essentially incapable of identifying the difference between  $y$  and  $y'$  since it depends only on the sign structure of the  $\Delta$  matrix and not the values themselves.

If an instance does not possess clustered observations, CA and CA-HR are strongly truthful across symmetric equilibria. In our practically motivated experimental setup, however, we find that clustered observations are encountered with a high frequency; see Remark 4.

### D.3. Insufficiency of a single evaluation per agent with homogenous responses in the Kong and Schoenebeck (2019) (KS) mechanism design framework

In this section, we show that it is impossible to design a mechanism within the KS mechanism design framework in the homogeneous responses setting, that incentivizes truthfulness with one evaluation per agent. The KS framework operates on a pair of agents and the payment of each agent is defined to be some scaling of an unbiased estimate of some mutual information measure constructed from their responses to a common set of tasks. The sufficiency of a single response per agent within this framework implies that the payment must be decided based only on the pair of agents' responses to a single task. We argue that such a payment scheme cannot strictly incentivize truthful behavior even in the homogenous, binary response setting. This result is not new; it has been shown in the general homogeneous responses setting in Jurca and Faltings (2011) (Theorem 1). We present a proof of the simpler binary responses case below for completeness. This result implies that there cannot be any mutual information measure satisfying information monotonicity, whose unbiased estimate can be constructed based on two agents' responses to a single evaluation task.

**Proposition D.1** (*Jurca and Faltings 2011*) *In any truthful mechanism in the homogenous, binary responses setting that calculates the payment of an agent only as a function of the responses of the agent and her peer to a single evaluation task, the payment to the agent does not depend on her own responses.*

*Proof.* Consider an evaluation task with only two responses:  $\mathcal{Y} = \{\text{Yes}, \text{No}\}$ . The payment scheme that depends on the responses of an agent and her peer to a common task is a specification of payment to the agent for every possible pair of responses. One of these payments can be 0 without loss of generality since additive shifts of payments across all possibilities do not change the incentive structure of the game. Let us suppose that the payments are as shown in Table 4, where it is assumed that the agent is the row player.

	Yes	No
Yes	a	b
No	c	0

**Table 4** The payments to the row agent corresponding to the pair of responses for the common evaluation task.

Let the generating model have two possible types  $\mathcal{X} = \{A, B\}$ , with  $P_X = (1/2, 1/2)$ ,  $\mathbf{p}(A) = (p, 1 - p)$ , and  $\mathbf{p}(B) = (q, 1 - q)$ . The expected payment of the agent if she reports 'Yes' on observing 'Yes' can be determined to be:

$$\frac{a(p^2/2 + q^2/2) + b(p(1-p)/2 + q(1-q)/2)}{p/2 + q/2}. \quad (82)$$

The expected payment of the agent if she reports 'No' on observing 'Yes' can be determined to be:

$$\frac{c(p^2/2 + q^2/2)}{p/2 + q/2}. \quad (83)$$

Thus reporting ‘Yes’ on observing ‘Yes’ yields a higher expected payment if

$$(c - a)(p^2 + q^2) \leq b(p(1 - p) + q(1 - q)). \quad (84)$$

The expected payment of the agent if she reports ‘No’ on observing ‘No’ can be determined to be:

$$\frac{c(p(1 - p)/2 + q(1 - q)/2)}{(1 - p)/2 + (1 - q)/2}. \quad (85)$$

The expected payment of the agent if she reports ‘Yes’ on observing ‘No’ can be determined to be:

$$\frac{a(p(1 - p)/2 + q(1 - q)/2) + b((1 - p)^2/2 + (1 - q)^2/2)}{(1 - p)/2 + (1 - q)/2}. \quad (86)$$

Thus reporting ‘No’ on observing ‘No’ yields a higher expected payment if

$$b((1 - p)^2 + (1 - q)^2) \leq (c - a)(p(1 - p) + q(1 - q)). \quad (87)$$

If we set  $p$  and  $q$  such that  $p^2 + q^2 = p(1 - p) + q(1 - q)$  (e.g.,  $p = q = 0.5$ ), then from Equation 84 we obtain  $b \geq c - a$  on the other hand, if we set  $p$  and  $q$  such that  $(1 - p)^2 + (1 - q)^2 = p(1 - p) + q(1 - q)$  (e.g.,  $p = q = 0.5$ ), then from Equation 91 we obtain  $b \leq c - a$ . Thus, we have  $b = c - a$ .

Next, if  $b = c - a > 0$ , then Equations 84 and 91, reduce to:

$$p^2 + q^2 \leq p(1 - p) + q(1 - q), \quad (88)$$

$$(1 - p)^2 + (1 - q)^2 \leq p(1 - p) + q(1 - q). \quad (89)$$

In this case, setting  $p = q = 0.25$  violates the second inequality. If  $b = c - a < 0$ , then Equations 84 and 91, reduce to:

$$p^2 + q^2 \geq p(1 - p) + q(1 - q), \quad (90)$$

$$(1 - p)^2 + (1 - q)^2 \geq p(1 - p) + q(1 - q). \quad (91)$$

In this case, setting  $p = q = 0.25$  violates the first inequality. Hence, we have that  $b = c - a = 0$ , i.e., the mechanism’s payments are independent of the reports of the agent.  $\square$

#### D.4. Infeasibility of a generic adaptation of the KS framework to multi-task, homogeneous responses settings and the special role of the square-root agreement measure (SRAM)

The design of SRA suggests that perhaps a generic adaptation of the KS mechanism to homogeneous responses setting that incentivizes single evaluations is possible under any mutual information measure. We argue that this is not true via the example of Shannon mutual information (Cover and Thomas 2012). For two random variables  $Y_1$  and  $Y_2$  taking values in finite sets  $\mathcal{Y}_1$  and  $\mathcal{Y}_2$  respectively, the Shannon mutual information is defined to be,

$$I(Y_1; Y_2) = \sum_{y \in \mathcal{Y}_1, y' \in \mathcal{Y}_2} P(Y_1 = y, Y_2 = y') \log \frac{P(Y_1 = y, Y_2 = y')}{P(Y_1 = y)P(Y_2 = y')}. \quad (92)$$

Suppose that the distribution of two agents’ responses to a common evaluation task is available to the platform (estimated from a large number of evaluation tasks). Then, along the lines of SRA, the mutual information measure above suggests the following mechanism.

1. Each agent  $j$  is paired with another randomly chosen agent  $j'$ , and their responses are compared.
2. If the response of agent  $j$  is  $y$  and that of agent  $j'$  is  $y'$ , then  $j$  gets a reward  $K \log \left( \frac{P(Y_j=y, Y_{j'}=y')}{P(Y_j=y)P(Y_{j'}=y')} \right)$ , where  $K$  is some positive constant.

Under this mechanism, if  $j$ 's true response is  $y$  and  $j'$  is truthful, her expected reward for a truthful report is,

$$K \sum_{y' \in \mathcal{Y}} P(Y_{j'} = y' | Y_j = y) \log \frac{P(Y_j = y, Y_{j'} = y')}{P(Y_j = y)P(Y_{j'} = y')} = K \sum_{y' \in \mathcal{Y}} \frac{P(Y_{j'} = y', Y_j = y)}{P(Y_j = y)} \log \frac{P(Y_j = y, Y_{j'} = y')}{P(Y_j = y)P(Y_{j'} = y')}. \quad (93)$$

Similarly, her reward for any other report  $\bar{y}$  is,

$$K \sum_{y' \in \mathcal{Y}} \frac{P(Y_{j'} = y', Y_j = y)}{P(Y_j = y)} \log \frac{P(Y_j = \bar{y}, Y_{j'} = y')}{P(Y_j = \bar{y})P(Y_{j'} = y')}. \quad (94)$$

Thus being truthful yields a higher reward if for any  $\bar{y} \neq y$ , expression in Equation 93 is higher than the one in Equation 94, which simplifies to the condition,

$$K \sum_{y' \in \mathcal{Y}} P(Y_j = y, Y_{j'} = y') \log \frac{P(Y_j = y, Y_{j'} = y')}{P(Y_j = \bar{y}, Y_{j'} = y')} - K P(Y_j = y) \log \frac{P(Y_j = y)}{P(Y_j = \bar{y})} \geq 0. \quad (95)$$

This inequality is not satisfied in general for homogeneous responses. We tested this condition in our experimental setup of Section 6. Assuming that there are  $|\mathcal{Y}| = 5$  responses as defined in that section, and two types of moving companies with delays exponentially distributed and mean delays drawn uniformly in  $[0, 60]$  (in minutes), we found that 629 of 10000 instances we generated violated the inequality in Equation 95.

## E. Auxillary results

**Proposition E.1** *If responses are categorical then they are self-predicting.*

*Proof.* For any two responses  $y$  and  $y'$ , the categorical responses condition says that,

$$P(Y_{j'} = y' | Y_j = y) \leq P(Y_{j'} = y'). \quad (96)$$

However, this implies that  $P(Y_{j'} = y') \leq P(Y_{j'} = y' | Y_j = y)$ . This means that for any two responses  $y$  and  $y'$ ,

$$P(Y_{j'} = y' | Y_j = y) \leq P(Y_{j'} = y' | Y_j = y'). \quad (97)$$

But this is exactly the self-prediction condition.  $\square$

**Proposition E.2** *Consider two exchangeable random variables,  $Y_1$  and  $Y_2$ , taking values in a finite set  $\mathcal{Y}$ . If their distribution satisfies the strict Cauchy-Schwarz property:*

$$\sqrt{P(Y_1 = Y_2 = y)} \sqrt{P(Y_1 = Y_2 = y')} > P(Y_1 = y, Y_2 = y'), \quad (98)$$

*for each  $y, y' \in \mathcal{Y}$ , then  $Y_1$  and  $Y_2$  are stochastically relevant random variables.*

*Proof.* We will show that stochastic irrelevance for two values  $y$  and  $y'$  implies that the CS property is satisfied for these values with an equality. Stochastic irrelevance for  $y$  and  $y'$  implies that the conditional distributions of  $Y_2$  given  $Y_1 = y$  and  $Y_1 = y'$  are identical. This implies that there is some constant  $C > 0$  such that  $(P(Y_1 = y', Y_2 = a); a \in \mathcal{Y}) = C \times (P(Y_1 = y, Y_2 = a); a \in \mathcal{Y})$ . In particular we have that:

$$P(Y_1 = Y_2 = y') = C \times P(Y_1 = y, Y_2 = y') \text{ and} \quad (99)$$

$$P(Y_1 = y', Y_2 = y) = C \times P(Y_1 = Y_2 = y). \quad (100)$$

We thus have,

$$\sqrt{P(Y_1 = Y_2 = y)}\sqrt{P(Y_1 = Y_2 = y')} = \sqrt{P(Y_1 = Y_2 = y)}\sqrt{C \times P(Y_1 = y, Y_2 = y')} \quad (101)$$

$$= \sqrt{P(Y_1 = Y_2 = y)}\sqrt{\frac{P(Y_1 = y', Y_2 = y)}{P(Y_1 = Y_2 = y)} \times P(Y_1 = y, Y_2 = y')} \quad (102)$$

$$= \sqrt{P(Y_1 = y', Y_2 = y)P(Y_1 = y, Y_2 = y')} \quad (103)$$

$$\stackrel{(a)}{=} P(Y_1 = y, Y_2 = y'). \quad (104)$$

Here (a) follows from exchangeability of  $Y_1$  and  $Y_2$ . Thus the CS property is satisfied with an equality for  $y$  and  $y'$ .  $\square$