



# Incentivizing Truthful Data Quality for Quality-Aware Mobile Data Crowdsourcing

Xiaowen Gong

Auburn University

Department of Electrical and Computer Engineering

Auburn, Alabama 36849, USA

xgong@auburn.edu

Ness Shroff

The Ohio State University

Department of Electrical and Computer Engineering &

Department of Computer Science and Engineering

Columbus, Ohio 43210, USA

shroff.11@osu.edu

## ABSTRACT

Mobile data crowdsourcing has found a broad range of applications (e.g., spectrum sensing, environmental monitoring) by leveraging the “wisdom” of a potentially large crowd of “workers” (i.e., mobile users). A key metric of crowdsourcing is data accuracy, which relies on the *quality* of the participating workers’ data (e.g., the probability that the data is equal to the ground truth). However, the data quality of a worker can be its own private information (which the worker learns, e.g., based on its location) that it may have incentive to misreport, which can in turn *mislead* the crowdsourcing requester about the accuracy of the data. This issue is further complicated by the fact that the worker can also manipulate its effort made in the crowdsourcing task and the data reported to the requester, which can also mislead the requester. In this paper, we devise truthful crowdsourcing mechanisms for *Quality, Effort, and Data Elicitation (QEDE)*, which incentivize strategic workers to truthfully report their private worker quality and data to the requester, and make truthful effort as desired by the requester. The truthful design of the QEDE mechanisms overcomes the lack of ground truth and the coupling in the joint elicitation of worker quality, effort, and data. Under the QEDE mechanisms, we characterize the socially optimal and the requester’s optimal task assignments, and analyze their performance. We show that the requester’s optimal assignment is determined by the largest “virtual valuation” rather than the highest quality among workers, which depends on the worker’s quality and the quality’s distribution. We evaluate the QEDE mechanisms using simulations which demonstrate the truthfulness of the mechanisms and the performance of the optimal task assignments.

## CCS CONCEPTS

• **Human-centered computing** → Ubiquitous and mobile computing; • **Networks** → Network economics;

## KEYWORDS

Mobile data crowdsourcing, data quality, incentive mechanism

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Mobihoc '18, June 26–29, 2018, Los Angeles, CA, USA

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-5770-8/18/06...\$15.00

<https://doi.org/10.1145/3209582.3209599>

## 1 INTRODUCTION

Mobile data crowdsourcing (referred to as “crowdsourcing” for brevity) has found a wide range of applications. Typical applications involves physical sensing tasks (also known as “crowdsensing”) such as spectrum sensing, traffic monitoring, and environmental monitoring. In principle, crowdsourcing leverages the “wisdom” of a potentially large crowd of workers (i.e., mobile users) for a crowdsourcing task. A key advantage of crowdsourcing lies in that it can exploit the diversity of inherently inaccurate data from many workers by aggregating the data obtained by the crowd, such that the data accuracy (also referred to as “data quality”) after aggregation can be substantially enhanced. With enormous opportunities brought by big data, crowdsourcing serves as an important first step for data mining tools to harness the power of big data in many application domains.

To fully exploit the potential of crowdsourcing, it is important to assign crowdsourcing tasks to workers based on their quality. A worker’s quality<sup>1</sup> captures the *intrinsic* accuracy of the worker’s data relative to the ground truth of the interested variable, and it generally varies for different workers depending on a worker’s characteristics (e.g., location, capabilities of sensors). For example, if the crowdsourcing task is to detect whether a licensed frequency band is idle or occupied by a licensed user (for opportunistic spectrum access by unlicensed users), then the quality of a worker’s data is the probability of correct detection, which depends on the worker’s location relative to the licensed user. Workers generally have *diverse* quality. A worker can learn its quality based on the knowledge of its characteristics, such as its location<sup>2</sup>. However, the quality of a worker’s can be its private information, which is unknown to and cannot be verified by the crowdsourcing requester. For example, a worker’s location is often its private information that is unknown to the requester. As a result, a strategic worker may have incentive to manipulate its quality revealed to the requester so as to gain an advantage. For example, a worker of low quality may pretend to have high quality in the hope of receiving a high reward for contributing high quality data to the task.

In addition to the worker quality, the data quality of a worker is also affected by its effort exerted in a crowdsourcing task. The data quality of a worker when it makes effort in the task is higher than when it makes no effort. For example, to detect whether a licensed

<sup>1</sup>We use “worker quality” and “quality” exchangeably in this paper. “Worker quality” should be distinguished from “data quality”.

<sup>2</sup>Alternatively, a worker can report its characteristics (e.g., location) that determines its quality to the requester, so that the requester can learn the worker’s quality. In this case, reporting the worker’s quality is equivalent to reporting its characteristics.

frequency band is idle, a worker should measure the signal in that band to make an estimate, rather than making a guess without any measuring. However, a worker’s effort can also be its hidden action that cannot be observed by the requester. Due to the inaccurate nature of the data, a strategic worker may report some arbitrary data to the requester without making effort in the task, while the requester is not able to verify whether effort was actually made. Furthermore, the data itself obtained by a worker from the task could also be its private information that it can manipulate in favor of itself.

In the presence of strategic workers with private worker quality, hidden effort, and private data, our goal is to incentivize workers to truthfully reveal their worker quality and data, and make truthful effort as desired by the crowdsourcing requester. Such a truthful mechanism is desirable as it eliminates the possibility of manipulation, which would encourage workers to participate in crowdsourcing. More importantly, the joint truthful elicitation of quality, effort, and data ensures that the requester can *correctly know the data accuracy of the collected data*, which is a key metric of crowdsourcing. This is in contrast to the situation of crowdsourcing with private participating cost, where manipulating the cost does not mislead the requester about the data accuracy.

The joint elicitation of quality, effort, and data calls for new truthful design that is different from existing mechanisms. First, a worker’s payoff as a function of its quality, effort, and data has a different structure from that of its private participating cost. As a result, existing designs for cost elicitation cannot work for the problem here. Second, due to the statistical dependency of a worker’s private data on its private quality and hidden effort, the joint elicitation of quality, effort, and data needs to overcome the coupling therein.

Given a truthful mechanism that can elicit quality, effort, and data from workers, an important question for the requester is to determine which worker(s) the task should be assigned to based on their quality, in order to maximize the social welfare or the requester’s payoff. This involves the tradeoff between assigning the task to more workers to improve the data accuracy, and assigning it to fewer workers to reduce the total cost incurred or total reward paid to the workers.

The main contributions of this paper can be summarized as follows.

- Under a quality-aware crowdsourcing framework, we devise truthful crowdsourcing mechanisms for Quality, Effort and Data Elicitation (QEDE). With general task assignment functions, the QEDE mechanisms incentivize strategic workers to truthfully reveal their private quality and data, and make truthful effort as desired by the crowdsourcing requester. The truthful design of the QEDE mechanisms overcomes the lack of ground truth and the coupling in the joint elicitation of worker quality, effort, and data, by exploiting the statistical dependency of a worker’s private data on its private worker quality and hidden effort.
- Under the QEDE mechanisms, we characterize the socially optimal (SO) and the requester’s optimal (RO) task assignments, and analyze their performance. We show that the RO assignment is determined by the largest virtual valuation

rather than the highest quality among workers, which depends on the worker’s quality and the quality’s distribution. We also show that, as the number of workers becomes large, the gap between the social welfare attained by the RO assignment and the SO assignment decreases and converges to 0.

- We evaluate the QEDE mechanisms using simulation results which demonstrate the truthfulness of the mechanisms and the performance of the RO and SO assignments.

The rest of this paper is organized as follows. Section 2 reviews related work. In Section 3, we describe the system model of quality-aware crowdsourcing with private data quality and formulate the problem of truthful mechanism design. In Section 4, we devise truthful mechanisms for Quality, Effort, and Data Elicitation (QEDE), and explain the ideas of the design and the rationale behind. In Section 5, we characterize the optimal assignments under the QEDE mechanisms and analyze their performance. Section 6 discusses modification of the QEDE mechanisms when there is no reference data from the requester. Simulation results are presented in Section 7. Section VII concludes this paper and discusses future work.

## 2 RELATED WORK

**Quality based data crowdsourcing.** The quality of data is important for allocating crowdsourcing tasks to workers, and has been studied in a few works [1–8]. One interesting line of work [6–8] in this direction has studied truthful mechanisms for information quality based task allocation where workers have private participating cost. Some other works have focused on learning the data quality of workers, e.g., by exploiting the correlation of their data for the same tasks [1, 3], or allocating tasks on the fly [5]. Different from these works, this paper focuses on the situation where quality is a worker’s private information that is unknown to the requester.

A recent work [9] has proposed a quality-aware crowdsourcing framework and devised truthful mechanisms for quality and effort elicitation. Compared to this paper, a key difference of [9] is that the data considered in [9] take continuous values and the quality is measured by the variance of the error, which can capture physical sensing tasks involving fine-grained measurements such as measuring temperature or air pollution. On the other hand, this paper focuses on data taking *discrete* values and the quality measured by the *correct probability*, which can capture physical sensing tasks as well as human intelligent tasks involving coarse-grained detection or classification, such as spectrum occupancy, image labeling. Moreover, the truthful mechanisms devised in this paper achieve joint elicitation of quality, effort, and data, which is stronger than in [9] which achieves quality and effort elicitation. As a result of these differences, the truthful mechanisms of in this paper and its analysis are non-trivially different from those in [9].

**Truthful crowdsourcing with private cost.** There have been a lot of recent research on incentive mechanisms for crowdsourcing [2, 4, 10–17]. Most of these mechanisms incentivize workers to truthfully reveal their participating cost. The cost is considered to be a strategic worker’s private information that it may not reveal truthfully without appropriate incentive. Different from these works, we study the situation where the quality of a worker’s data

obtained from a crowdsourcing task is the worker’s private information that it can manipulate. A worker’s payoff as a function of its private quality has a different structure than that of its private cost. As a result, existing designs for cost elicitation (such as the classical VCG auction and the characterization of truthful mechanisms [18, Theorem 9.36]) cannot work for quality elicitation, so that new design is needed. Furthermore, this paper aims at joint elicitation of quality, effort, and data. The statistical dependency of a worker’s private data on its private quality and hidden effort leads to the coupling in the elicitation of quality, effort, and data, which needs to be addressed.

### Mechanism design for hidden actions and data elicitation.

There have been many studies on mechanism design for hidden actions in the economics literature [19], which is concerned with strategic agents that can take hidden actions not desired by a principal who recruits the agents to work on a task. A few recent works have studied this problem in the context of crowdsourcing [14, 20–23]. Cai *et al.* [21] have designed truthful mechanisms to incentivize workers to make effort as desired in statistical estimation. Luo *et al.* [14] have designed mechanisms that not only elicit desired effort from workers but also truthful revelation of their private cost and data. This paper is different from these works as we aim to jointly elicit workers’ private quality, private data, and hidden effort, which cannot be achieved by existing truthful design.

Mechanism design for truthful elicitation of strategic agents’ data (e.g., opinions) has been extensively studied in various applications (e.g., [24]), more recently for crowdsourcing [14, 20, 22, 23, 25]. The data of an agent can be its private information that it can manipulate in favor of its benefit. Different from the existing works, in this paper we aim to design truthful mechanisms that jointly elicit workers’ private data, private quality, and hidden effort, which calls for new truthful design.

## 3 QUALITY-AWARE DATA CROWDSOURCING FRAMEWORK

We consider a crowdsourcing requester (also referred to as worker  $0^3$ ) recruiting a set of workers  $\mathcal{N} \triangleq \{1, \dots, N\}$  to work on a task. For convenience, let  $\mathcal{N}^+ \triangleq \mathcal{N} \cup \{0\}$ . The structure and procedure of the crowdsourcing system is illustrated in Fig. 1 and described in detail as follows.

### 3.1 Data crowdsourcing with private data quality

**Data observation.** The crowdsourcing task is to observe and estimate an unknown and random variable of interest  $X$ . The interested variable  $X$  takes discrete values (e.g., the answer of a multi-choice question). For ease of exposition, we assume that  $X$  takes one of two possible values<sup>4</sup> 0 and 1. We also assume that  $X$  takes value 0 or 1 equally likely in the prior distribution. After working on the task, each worker  $i \in \mathcal{N}^+$  (i.e., including the requester) obtains random data  $D_i$ . The accuracy of the data  $D_i$  is quantified by the *correct probability*  $p_i$ , which is the probability that  $D_i$  is equal to

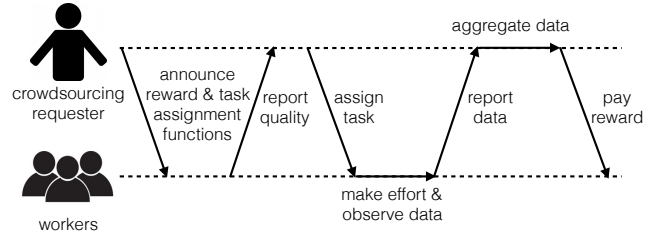


Figure 1: Structure and procedure of the quality-aware crowdsourcing framework.

the interested variable  $X$ , given by

$$p_i \triangleq \Pr(D_i = X) = q_i e_i + 0.5(1 - e_i). \quad (1)$$

Here the correct probability  $p_i$  depends on the *worker quality*  $q_i$  of worker  $i$  and the *effort*  $e_i$  exerted by worker  $i$  in the task, which is explained as follows.

**Worker quality.** Given that worker  $i$  makes effort in the task, the quality  $q_i \in [0, 1]$  determines the correct probability  $p_i$  which quantifies how accurate  $D_i$  is. The quality  $q_i$  is an *intrinsic* coefficient that captures worker  $i$ ’s capability for the task. Note that a larger  $q_i$  means higher quality. The quality generally varies for different workers. We assume that each worker  $i \in \mathcal{N}^+$  knows its quality  $q_i$  (e.g., by learning the correct probability based on its location). However, the quality of each worker  $i \in \mathcal{N}$  is unknown to the requester. For ease of exposition, we assume that each worker’s quality  $q_i$  is within the range of  $[\underline{q}, \bar{q}]$  which is known to the requester.

**Work effort.** The effort  $e_i \in \{0, 1\}$  represents whether worker  $i$  makes effort in the task, where  $e_i = 1$  and  $e_i = 0$  indicate making and not making effort, respectively. If worker  $i$  makes effort, then the correct probability  $p_i$  of worker  $i$  is equal to the worker quality  $q_i$ ; otherwise,  $p_i$  is equal to 0.5, which means that worker  $i$  simply makes a guess of  $X$  randomly according to the prior distribution. To ensure that making effort is meaningful, we assume that  $q_i > 0.5$ . Therefore, given the quality  $q_i$ , making effort  $e_i = 1$  means a larger correct probability  $p_i$  and thus higher accuracy of  $D_i$  than not making effort. The binary effort model (i.e., either making effort or not) is reasonable (also used in, e.g., [20, 22, 23]), as workers’ behavior tend to be simple in practice. We assume that each worker  $i$  can control its effort  $e_i$ , but it cannot be observed by the requester. We assume that the requester itself always makes effort in the task (i.e.,  $e_0 = 1$ ).

**Task assignment.** The requester assigns the crowdsourcing task to the workers by assigning effort  $e_i'$  to each worker  $i$ , which indicates whether it desires worker  $i$  to make effort in the task, based on the workers’ quality. To this end, each worker  $i$  reports its quality  $q_i'$  to the requester<sup>5</sup>. Since the true quality  $q_i$  is worker  $i$ ’s private information, it may manipulate the reported quality  $q_i'$  to its own advantage such that  $q_i' \neq q_i$ . Based on the quality reported by all workers, the requester determines the effort  $e_i'$  assigned to

<sup>3</sup>In Section 6, we will address the situation where the requester cannot work on the task as a worker.

<sup>4</sup>The results of this paper can be fairly easily extended to the case of multiple possible values of the interested variable  $X$ .

<sup>5</sup>Workers should report their worker quality  $\{q_i'\}$  rather than data quality  $\{p_i'\}$ , as it allows the requester to assign the task to workers based on their quality  $\{q_i'\}$ . This is desirable for achieving some particular task assignments, such as the socially optimal assignment.

each worker  $i$  according to some assignment function

$$e'_i(\mathbf{q}') \quad (2)$$

and notifies worker  $i$  of  $e'_i$ . The assignment function  $e'_i(\mathbf{q}')$  is pre-defined by the requester and announced to all the workers before they report their quality. A worker's assignment generally varies for different workers due to the diversity of their quality. Intuitively, a worker of high quality is preferred to be assigned to the task. Note that in general the assignment  $e'_i$  is not only dependent on the quality  $q'_i$  reported by worker  $i$  but also on the quality  $q'_{-i}$  reported by the other workers. After being assigned effort  $e'_i$  to, each worker  $i$  works on the task by making actual effort  $e_i$ . Since  $e_i$  is a hidden action of worker  $i$ , it may manipulate it against the assignment  $e'_i$  to its own advantage such that  $e_i \neq e'_i$ . After obtaining data  $d_i$  from the task (which is a sample realization of the random data  $D_i$ ), each worker  $i$  reports data  $d'_i$  to the requester. Since  $d_i$  is also private information of worker  $i$ , it may manipulate the reported data  $d'_i$  against the actual obtained data  $d_i$  to its own advantage such that  $d_i \neq d'_i$ .

**Data aggregation.** After collecting all the data  $\mathbf{d}$  reported by the workers, the requester aggregates the data  $\mathbf{d}$  by making the optimal estimate  $x_0$  of the interested variable  $X$  based on  $\mathbf{d}$ . The optimal estimate  $x_0$  maximizes the posterior probability that  $x_0$  is equal to the ground truth  $x$ , i.e.,

$$x_0(\mathbf{q}', \mathbf{e}', \mathbf{d}') \triangleq \arg \max_{d \in \{0,1\}} E_{X|d'(\mathbf{q}', \mathbf{e}')} [\mathbf{1}_{X=d}]. \quad (3)$$

Note that the distribution of  $X$  conditioned on  $\mathbf{d}'$  depends on workers' reported quality  $\mathbf{q}'$  and assigned effort  $\mathbf{e}'$ . Then the utility of crowdsourcing is represented by the correct probability  $p_c$  of the optimal estimate  $x_0$ , given by

$$p_c(\mathbf{q}', \mathbf{e}', \mathbf{d}') \triangleq E_{X|d(\mathbf{q}, \mathbf{e})} [\mathbf{1}_{X=x_0(\mathbf{q}', \mathbf{e}', \mathbf{d}')}]. \quad (4)$$

Note that the expectation is over the posteriori distribution  $X|d(\mathbf{q}, \mathbf{e})$  conditioned on the true data  $\mathbf{d}$  depending on the true quality  $\mathbf{q}$  and actual effort  $\mathbf{e}$ . If the task is not assigned to any worker (i.e.,  $e'_i = 0$ ,  $\forall i$ ), then the correct probability  $p_c$  is defined to be 0.

**Reward payment.** On the other hand, the requester pays a reward  $r_i$  to each worker  $i$  for working on the task, according to a certain reward function:

$$r_i(\mathbf{q}', e'_i, d'_i, d_0). \quad (5)$$

Note that the reward  $r_i$  depends on the *reference* data  $d_0$  obtained by the requester itself. The reward function is also pre-defined by the requester, and announced to all the workers before they report their quality (together with the assignment function  $e'_i(\mathbf{q}')$ ). Note that the reward function can only depend on the information that the requester knows, i.e.,  $\mathbf{q}'$ ,  $\mathbf{e}'$ ,  $\mathbf{d}'$ , and  $d_0$ .

### 3.2 Mechanism Design Objective

Based on the crowdsourcing system described above, each worker  $i$ 's payoff  $u_i$  is the reward  $r_i$  paid by the requester minus its cost in the task, given by,

$$u_i(\mathbf{q}', e_i, d'_i, d_0) \triangleq r_i(\mathbf{q}', e'_i, d'_i, d_0) - c_i e_i. \quad (6)$$

Here the cost  $c_i$  represents how much resource is consumed by worker  $i$  (e.g., how much time is spent by worker  $i$ ) if it makes effort  $e_i = 1$  in the task. If worker  $i$  make no effort  $e_i = 0$ , it incurs no

cost. Note that the relative weight of the cost  $c_i$  with respect to the reward  $r_i$  in (6) can be captured by  $c_i$ . We assume that workers have the same cost<sup>6</sup>  $c$  (i.e.,  $c_i = c, \forall i$ ) which is known to the requester. This assumption is reasonable when the cost  $c$  is determined by a uniform market price for working on a task.

The requester's payoff  $u_0$  is the crowdsourcing utility (i.e., the correct probability  $p_c$ ) minus the total reward paid to the workers, i.e.,

$$u_0(\mathbf{q}', \mathbf{e}', \mathbf{d}') \triangleq p_c(\mathbf{q}', \mathbf{e}', \mathbf{d}') - \sum_{i \in \mathcal{N}} r_i(\mathbf{q}', e'_i, d'_i, d_0). \quad (7)$$

As the workers have private quality and data and make hidden effort, if any worker manipulates its reported quality, reported data, or actual effort, then the estimate  $x_0$  found by the requester would be different from the correct estimator, i.e.,

$$x_0(\mathbf{q}', \mathbf{e}', \mathbf{d}') \neq x_0(\mathbf{q}, \mathbf{e}, \mathbf{d}).$$

More importantly, the correct probability  $p_c$  found by the requester would be different from the correct one, i.e.,

$$p_c(\mathbf{q}', \mathbf{e}', \mathbf{d}') \neq p_c(\mathbf{q}, \mathbf{e}, \mathbf{d}).$$

This means that manipulation would lead to the requester's *incorrect knowledge of the correct probability!* This is highly undesirable since the data accuracy is often a key performance metric that the requester needs to know correctly (e.g., to meet some threshold requirement). Note that this issue does not arise in the setting where workers have private cost only, since manipulating the cost can affect only the crowdsourcing utility and the reward payment but cannot affect the requester's knowledge of the data accuracy. Furthermore, the possibility of manipulation could result in concerns that discourage workers to participate in crowdsourcing. Thus motivated, we aim to design a mechanism, which is a pair of an assignment function  $e'_i(\mathbf{q}')$  and a reward function  $r_i(\mathbf{q}', e'_i, d'_i, d_0)$ , that can achieve the property of *incentive compatibility* as stated below.

**DEFINITION 1.** A mechanism is *dominant incentive-compatible (DIC)* if, given any quality reported by the other workers, the optimal strategy of each worker  $i$  for maximizing its expected payoff is to truthfully report its quality and data, and make the effort desired by the requester, i.e.,

$$E_{D_0|d_i(q_i, e_i)} [u_i(q_i, \mathbf{q}'_{-i}, e_i, d_i, D_0)] \geq E_{D_0|d_i(q_i, e_i)} [u_i(q'_i, \mathbf{q}'_{-i}, e'_i, d'_i, D_0)], \forall (q'_i, e_i, d'_i), \forall \mathbf{q}'_{-i}.$$

Another natural and desirable property we aim to achieve is that each worker's expected reward should at least compensate its cost (i.e., its expected payoff is nonnegative), since otherwise the worker would not participate in crowdsourcing for a payoff of 0. This property of *individual rationality* is stated as follows.

**DEFINITION 2.** A mechanism is *individually rational (IR)* if for each worker  $i$ , given that it truthfully reports its quality and makes the effort desired by the requester, its expected payoff is nonnegative, i.e.,

$$E_{D_0|d_i(q_i, e_i)} [u_i(q_i, \mathbf{q}'_{-i}, e'_i, d_i, D_0)] \geq 0, \forall \mathbf{q}'_{-i}.$$

<sup>6</sup>The truthful mechanisms still hold when workers have diverse costs  $c_i$  (i.e.,  $c_i \neq c_j, \forall i \neq j$ ) which are known to the requester.

## 4 TRUTHFUL QUALITY, EFFORT, AND DATA ELICITATION FOR CROWDSOURCING

In this section, we design truthful crowdsourcing mechanisms that achieve the DIC and IR properties.

We first present the QEDE mechanisms as follows.

**DEFINITION 3.** *A Quality, Effort, and Data Elicitation (QEDE) mechanism consists of any assignment function  $e'_i(\mathbf{q}')$  that satisfies the condition in (8) and the reward function  $r_i(d_0, d_i, \mathbf{q}', e'_i)$  given by (9) based on that  $e'_i(\mathbf{q}')$ :*

$$e'_i(q'_i, \mathbf{q}'_{-i}) \geq e'_i(q'_i, \mathbf{q}'_{-i}), \forall q'_i \leq q'_i, \forall \mathbf{q}'_{-i} \quad (8)$$

$$r_i(d_0, d'_i, \mathbf{q}', e'_i) = kq'_i e'_i(\mathbf{q}') \left[ \frac{1_{d_0=d_i} + q_0 - 1}{2q_0 - 1} \right] + ce'_i(\mathbf{q}') + \int_q^{q'_i} kqe'_i(q, \mathbf{q}'_{-i})dq - kq'_i e'_i(\mathbf{q}') [0.5 + (q'_i - 0.5)e'_i(\mathbf{q}')] \quad (9)$$

where  $k$  is any constant that satisfies the condition

$$k \geq \frac{c}{q(q-0.5)} \quad (10)$$

and  $1_A$  is the indicator function that is equal to 1 if condition  $A$  is true and 0 otherwise.

The condition (8) is a general *monotonicity* property for the task assignment functions: given any quality of the other workers, if the task is assigned to a worker, then the task is still assigned to that worker when its quality improves. Intuitively, these assignment functions are natural and desirable for system efficiency. Next we will explain the main ideas of the design of the QEDE mechanisms (8) and (9). In particular, we will show successively that, for each worker, 1) it is optimal to report its true data (Lemma 1); 2) it is optimal to make actual effort as desired (Lemma 2); 3) it is optimal to report the true quality (Lemma 3). As a result, the DIC property is achieved (Theorem 1). Then we will explain the rationale behind the design in Remark 1.

In the following, we show how the QEDE mechanisms achieve the DIC property (with the proofs in Appendix). We first show that any worker's optimal reported data is to report the true data, independent of its reported quality and actual effort. Given the lack of the ground truth  $x$ , this is achieved by the peer prediction mechanism (see, e.g., [20, 22, 23, 25]) which compares the reported data  $d_i$  with the reference data  $d_0$  from the requester.

**LEMMA 1.** *Under the QEDE mechanisms, given that any worker  $i$  reports any quality  $q'_i$  and makes any effort  $e'_i$ , its optimal reported data is its true data  $d'_i = d_i$ .*

Using Lemma 1, given that worker  $i$  reports the optimal data  $d'_i = d_i$ , we can express its expected payoff as

$$E_{D_0|d_i(q_i, e_i)}[u_i(\mathbf{q}', e'_i, d_i, D_0)] = kq'_i e'_i(\mathbf{q}') [0.5 + (q_i - 0.5)e_i] + \int_q^{q'_i} kqe'_i(q, \mathbf{q}'_{-i})dq - kq'_i e'_i(\mathbf{q}') [0.5 + (q'_i - 0.5)e'_i(\mathbf{q}')] + ce'_i(\mathbf{q}') - ce_i. \quad (11)$$

Then we have

$$E_{D_0|d_i(q_i, e_i)}[u_i(\mathbf{q}', e'_i, D_i, D_0)] = E_{D_0|d_i(q_i, e_i)}[u_i(\mathbf{q}', e'_i, d_i, D_0)], \forall d'_i$$

since the right hand side of (11) is independent of  $d_i$ . For convenience, we can define

$$\bar{u}_i(\mathbf{q}', q_i, e'_i, e_i) \triangleq E_{D_0|d_i(q_i, e_i)}[u_i(\mathbf{q}', e'_i, D_i, D_0)] \quad (12)$$

according to (11). Then we show that, as worker  $i$  can only affect its payoff in (12) via its reported quality  $q'_i$  and actual effort  $e_i$ , its optimal actual effort is the desired effort  $e'_i$ , independent of  $q_i$  and true quality  $q_i$ .

**LEMMA 2.** *Under the QEDE mechanisms, given that any worker  $i$  reports any quality  $q'_i$  and truthfully reports its data  $d_i$ , its optimal actual effort is the desired effort  $e_i = e'_i$ .*

Using Lemma 2, given that worker  $i$  reports the optimal data  $d'_i = d_i$  and makes the optimal effort  $e_i = e'_i$ , we can express its expected payoff using (11) and (12) as

$$\bar{u}_i(\mathbf{q}', q_i, e'_i, e_i) = kq'_i e'_i(\mathbf{q}') (q_i - q'_i) + \int_q^{q'_i} kqe'_i(q, \mathbf{q}'_{-i})dq. \quad (13)$$

For convenience, we can define

$$\hat{u}_i(\mathbf{q}', q_i, e'_i) \triangleq \bar{u}_i(\mathbf{q}', q_i, e'_i, e_i) \quad (14)$$

according to (13).

Next we show that, as worker  $i$  can only affect its payoff in (14) via its reported quality  $q'_i$ , its optimal reported quality is its true quality  $q_i$ , under the general condition (8) on the assignment function  $e'_i(\mathbf{q}')$ .

**LEMMA 3.** *Under the QEDE mechanisms, given that any worker  $i$  truthfully reports its data  $d_i$  and makes the desired effort  $e'_i$ , its optimal reported quality is its true quality  $q'_i = q_i$ .*

Using Lemmas 1, 2, and 3, we can show that the DIC property is achieved as in the next theorem. Using (13) and (14), given that worker  $i$  reports the optimal data  $d'_i = d_i$ , makes the optimal effort  $e_i = e'_i$ , and reports the optimal quality  $q'_i = q_i$ , its payoff is given by

$$\hat{u}_i(q_i, \mathbf{q}'_{-i}, q_i, e'_i) = k \int_q^{q_i} qe'_i(q, \mathbf{q}'_{-i})dq. \quad (15)$$

It follows that the IR property is also achieved since  $\hat{u}_i(q_i, \mathbf{q}'_{-i}, q_i, e'_i) \geq 0$  due to that  $e'_i(\mathbf{q}') \geq 0, \forall \mathbf{q}'$ .

**THEOREM 1.** *The QEDE mechanisms are DIC and IR.*

*Remark 1:* We explain the design rationale of the QEDE mechanisms as follows. We first observe that the optimal reported data  $d'_i$  that maximizes  $E_{D_0|d_i(q_i, e_i)}[1_{D_0=d'_i}]$  which is the probability that  $d'_i$  is equal to  $d_0$  is always the true data  $d_i$  (as in Lemma 1). Then we can design the expected reward as a function of  $0.5 + (q_i - 0.5)e_i$  which is the probability that  $d'_i$  is equal to the ground truth  $x$ , such that worker  $i$ 's expected payoff depends on the true quality  $q_i$  and the actual effort  $e_i$ , and is independent of the true data  $d_i$  (as in (11)). Now the expected payoff only depends on  $q'_i, e'_i, q_i$ , and  $e_i$  (as in (12)). Then we can design the reward function such that the optimal actual effort  $e_i$  that maximizes the payoff is always the

desired effort  $e'_i$  and independent of  $q'_i$  and  $q_i$  (as in Lemma 2). As a result, worker  $i$ 's payoff now only depends on  $q'_i$ ,  $e'_i$ , and  $q_i$  (as in (13) and (14)). Next we further design the reward function such that, under the monotonicity condition (8) on  $e'_i$ , the optimal reported quality  $q'_i$  is always the true quality  $q_i$  and independent of  $e'_i$  (as in Lemma 3).

It follows from (15) that when all workers behave truthfully (i.e.,  $q'_i = q_i$ ,  $e_i = e'_i$ , and  $d'_i = d_i$ ), the total expected reward paid by the requester is

$$\sum_{i \in \mathcal{N}} \left( k \int_{\underline{q}}^{q_i} q e'_i(q, \mathbf{q}'_{-i}) dq + c e'_i(\mathbf{q}') \right). \quad (16)$$

It can be seen from (16) that, to minimize the requester's payoff,  $k$  should be minimized such that condition (10) is satisfied with equality. We assume that this equality holds in the rest of this paper.

*Remark 2:* We can see from (16) that the requester's payment for each worker consists of two parts: while the second part  $c e'_i(\mathbf{q}')$  is to compensate the worker's cost, the first part (i.e., the integral multiplied by  $k$ ) is to elicit the worker's truthful behavior. This shows that the requester pays more than needed to cover the cost by the truth-eliciting payment (also known as "information rent" [26]), which is due to the requester's uncertainty of workers' quality. We can also observe from (16) that, as the multiplier  $k$  (determined by (10)) and the integral are both decreasing in the lower bound  $\underline{q}$  of workers' quality, the truth-eliciting payment is also decreasing in  $\underline{q}$  (as illustrated by Fig. 5 in Section 7). Intuitively, this is because the requester knows more information (i.e., less uncertainty) of workers' quality with a larger  $\underline{q}$ . We further observe that both the truth-eliciting payment and the payment for compensating the cost are increasing in the cost  $c$ . This shows that the requester's payment decreases faster than the cost when  $c$  increases (as illustrated by Fig. 6 in Section 7).

## 5 OPTIMAL TASK ASSIGNMENT FOR TRUTHFUL CROWDSOURCING

In Section 4, we have shown that the DIC and IR properties can be achieved by all the QEDE mechanisms which have general assignment functions that satisfy condition (8). In this section, we will find the optimal assignment under the QEDE mechanisms that maximizes the social welfare and the requester's payoff, respectively. Because of the DIC property, in this section we assume that  $\mathbf{q}' = \mathbf{q}$ ,  $\mathbf{e} = \mathbf{e}'$ , and  $\mathbf{d}' = \mathbf{d}$ . Therefore, for brevity, we use  $\mathbf{q}$ ,  $\mathbf{e}$ , and  $\mathbf{d}$  instead of  $\mathbf{q}'$ ,  $\mathbf{e}'$ , and  $\mathbf{d}'$  respectively.

### 5.1 Socially optimal assignment

An important metric for the assignment  $\mathbf{e}(\mathbf{q})$  is system efficiency, which is measured by the social welfare (the requester may be also interested in this objective). The social welfare  $v$  is the crowdsourcing utility (i.e., the correct probability  $p_c$ ) minus the total cost of all workers, i.e.,

$$v(\mathbf{q}, \mathbf{e}(\mathbf{q})) \triangleq E_{D(\mathbf{q}, \mathbf{e})}[p_c(\mathbf{q}, \mathbf{e}, D)] - \sum_{i \in \mathcal{N}} c e_i. \quad (17)$$

**DEFINITION 4.** *The socially optimal (SO) assignment  $\mathbf{e}^{so}(\mathbf{q})$  for the QEDE mechanisms is the assignment function  $\mathbf{e}(\mathbf{q})$  satisfying*

*condition (8) that maximizes the social welfare, i.e.,*

$$\{\mathbf{e}^{so}(\mathbf{q}), \forall \mathbf{q}\} \triangleq \arg \max_{\{\mathbf{e}(\mathbf{q}), \forall \mathbf{q}\} \text{ s.t. (8)}} E_{\mathcal{Q}}[v(\mathcal{Q}, \mathbf{e}(\mathbf{q}))]. \quad (18)$$

We first consider single-worker assignment, which consists of the assignment functions that assign the task to at most one worker. The advantage of single-worker assignment is that it simplifies the implementation of crowdsourcing: the requester needs to collect data from only one worker rather than potentially many workers. We should note that single-worker assignment still *exploits the diversity of potentially many available workers* in crowdsourcing, as the worker is selected based on the quality of all the workers. Under single-worker assignment, the requester's optimal estimate  $x_0$  of the interested variable  $X$  is just equal to the data  $d_i$  reported by the worker  $i$  who works on the task, and the correct probability  $p_c$  is equal to the quality  $q_i$  of that worker  $i$ .

We can find the socially optimal assignment for single-worker assignment as follows.

**PROPOSITION 1.** *For single-worker assignment, the socially optimal assignment is given by*

$$e_i^{so}(\mathbf{q}) = \begin{cases} 1, & i = \arg \max_j q_j \text{ and } q_i \geq c \\ 0, & \text{otherwise} \end{cases}. \quad (19)$$

Proposition 1 shows that the task is assigned to the "best" worker  $i$  that has the highest quality  $q_i$  if and only if the cost  $c$  is less than the quality  $q_i$ . This is clearly because the best worker maximizes the correct probability  $p_c$  and thus the social welfare  $v$ . It is also clear that the SO assignment (19) satisfies the monotonicity condition (8) of the QEDE mechanisms (thus the proof of Proposition 1 follows and is omitted). We should note that although single-worker assignment involves only one worker to work on the task, it still exploits the diversity of multiple available workers, as it selects the worker of the highest quality.

Next we consider general assignment functions that can assign the task to multiple workers. It can be shown (e.g., see [5]) that the optimal estimate given in (3) is equivalent to that  $x_0 = 1$  if and only if

$$\prod_{i: e_i=1, d_i=1} \log \frac{q_i}{1-q_i} \geq \prod_{j: e_j=1, d_j=1} \log \frac{q_j}{1-q_j}$$

and  $x_0 = 0$  otherwise. It has been shown in [5] that, without imposing the condition (8) of the QEDE mechanisms, the optimal assignment that maximizes the social welfare satisfies an intuitive property: there exists some  $k$  such that the task is assigned to only the top  $k$  workers that have the highest quality. As a result, it can be found by an efficient exhaustive search algorithm with linear complexity as described in Algorithm 1. In the following, we show that the solution found by Algorithm 1 is also the SO assignment for the QEDE mechanisms. Due to space limitation, the proofs of the results in the rest of this paper are given in our online technical report [27].

**PROPOSITION 2.** *For multi-worker assignment, the socially optimal assignment is found by Algorithm 1.*

The main idea of the proof of Proposition 2 is to show that the output of Algorithm 1 satisfies the monotonicity condition (8) of the QEDE mechanisms.

---

**Algorithm 1:** Find the socially optimal assignment for multi-worker assignment

---

```

1 Index workers in the descending order of their quality, i.e.,
   $q_1 \geq q_2 \geq \dots \geq q_N$ ;
2  $e_j \leftarrow 0, \forall j, t \leftarrow v(\mathbf{q}, \mathbf{e}), i = 1$ ;
3 while  $i \leq N$ ;
4 do
5    $e_i \leftarrow 1$ ;
6   if  $v(\mathbf{q}, \mathbf{e}) > t$  then
7      $e^* \leftarrow \mathbf{e}$ ;
8   end
9    $i \leftarrow i + 1$ ;
10 end
11 return  $e^{so}$ ;

```

---

## 5.2 Requester's optimal assignment

A desirable objective for the requester is to find the optimal assignment that maximizes its expected payoff.

**DEFINITION 5.** *The crowdsourcing requester's optimal (RO) assignment  $e^*(\mathbf{q})$  for the QEDE mechanism is the assignment function  $e(\mathbf{q})$  satisfying condition (8) that maximizes the requester's expected payoff (7), i.e.,*

$$\{e^*(\mathbf{q}), \forall \mathbf{q}\} \triangleq \arg \max_{\{e(\mathbf{q}), \forall \mathbf{q}\} \text{ s.t. (8)}} E_{D(Q, \mathbf{e})}[u_0(Q, \mathbf{e}, D)]. \quad (20)$$

For ease of analysis, in the rest of this subsection we focus on single-worker assignment. One reason is that the characterization of the RO assignment for single-worker assignment and the corresponding performance analysis provide useful insights. We further assume that each worker's quality follows an independent and identical distribution over an interval  $[\underline{q}, \bar{q}]$ , which is known to the requester.

**PROPOSITION 3.** *For single-worker assignment, when*

$$\alpha(q) \triangleq q + kq \frac{F(q) - 1}{f(q)}$$

*is an increasing function of  $q$ , the RO assignment is given by*

$$e_i^*(\mathbf{q}) = \begin{cases} 1, & i = \arg \max_j \alpha(q_j) \text{ and } \alpha(q_i) \geq c \\ 0, & \text{otherwise} \end{cases} \quad (21)$$

*where  $f(q)$  and  $F(q)$  denote the probability density function (PDF) and cumulative density function (CDF) of each worker's quality, respectively.*

We should note that the property that  $\alpha(q)$  is an increasing function of  $q$  holds under mild conditions, e.g., when  $F(q)$  and  $f(q)$  follow a uniform distribution. We assume that this property holds in the rest of this section.

**Remark 3:** Proposition 3 shows that the task is assigned to the "best"<sup>7</sup> worker  $i$  that has the largest "virtual valuation"  $\alpha(q_i)$ , if and only if the cost  $c$  is less than the virtual valuation  $\alpha(q_i)$ . Note that each worker  $i$ 's virtual valuation depends on not only its quality  $q_i$  but also the quality's distribution  $F(q_i)$  and  $f(q_i)$ . This implies that

<sup>7</sup>If there are multiple "best" workers, only one of them is selected by breaking the tie randomly.

the range of a worker's possible quality, represented by  $\Delta q \triangleq \bar{q} - \underline{q}$ , affects the task assignment. For ease of analysis, suppose  $F(q)$  and  $f(q)$  follow a uniform distribution such that  $(F(q) - 1)/f(q) = q - \bar{q}$ . Given workers' quality, when the upper bound  $\bar{q}$  of workers' quality decreases so that the quality range  $\Delta q$  decreases, each worker's virtual valuation  $\alpha(q_i)$  increases, and thus the condition  $\alpha(q_i) \geq c$  for assigning the task to the best worker  $i$  is more likely to hold. Intuitively, this is because a smaller quality range incurs a lower truth-eliciting payment in (16) by the requester in order to achieve truthful elicitation, which increases the requester's payoff. In the special case of  $\Delta q = 0$ , a worker's virtual valuation is equal to its quality. The concept of virtual valuation was introduced by Myerson [18] and is in the same spirit as the result here.

**Remark 4:** Comparing (19) and (21), we can see that the SO assignment is similar to the RO assignment in that the task can be assigned only to the best worker  $i$  that has the highest quality  $q_i$ . This is because the virtual valuation  $\alpha(q)$  is increasing in  $q$  and thus  $\arg \max_j q_j = \arg \max_j \alpha(q_j)$ . The difference is that the RO assignment assigns the task to the best worker  $i$  based on the condition  $\alpha(q_i) \geq c$  rather than the condition  $q_i \geq c$  for the SO assignment. Since it can be easily seen that  $\alpha(q_i) \leq q_i$  always holds, there exist some values of  $q_i$  such that  $\alpha(q_i) < c$  while  $q_i \geq c$ . In this case, the RO assignment  $e_i^*(\mathbf{q})$  is different from the SO assignment  $e_i^{so}(\mathbf{q})$  and attains lower social welfare than  $e_i^{so}(\mathbf{q})$ . Intuitively, this is because, although assigning the task increases the crowdsourcing utility and also the social welfare, it incurs a too high truth-eliciting payment. As a result, the RO assignment is not socially optimal, and the gap is essentially due to the asymmetry of workers' quality information between the workers and the requester.

## 5.3 Performance Analysis

Next we analyze the impact of system parameters on the performance of the SO and RO assignments.

**PROPOSITION 4.** *The expected RO payoff  $E_Q[u_0(e^*(Q))]$  attained by the RO assignment, the expected SO social welfare  $E_Q[v(e_1^{so}(Q))]$ , and the expected social welfare  $E_Q[v(e_1^*(Q))]$  attained by the RO assignment, all increase as the number of workers  $N$  increases, or the cost  $c$  decreases.*

**Remark 5:** Proposition 4 shows that the RO payoff and social welfare benefit from a greater diversity gain in workers' quality. This is because when there are more workers, the quality of the best worker is likely to be higher, which improves the crowdsourcing utility. On the other hand, a larger  $c$  increases the cost incurred to workers as well as the truth-eliciting payment (i.e., the first term in (16)), and thus reduces the RO payoff and social welfare.

**PROPOSITION 5.** *The gap between the expected social welfare of the SO assignment and the RO assignment  $E_Q[v(e_1^{so}(Q))] - E_Q[v(e_1^*(Q))]$  converges to 0 as the number of workers  $N$  goes to infinity.*

**Remark 6:** Proposition 5 shows that the performance gap between the RO assignment and the SO assignment decreases to 0 asymptotically as the number of workers increases. This is because when there are more workers, the quality of the best worker improves, so that the gap between the RO and SO assignments decreases to 0 (i.e., they are more often the same), and thus the gap between their social welfare also decreases to 0.

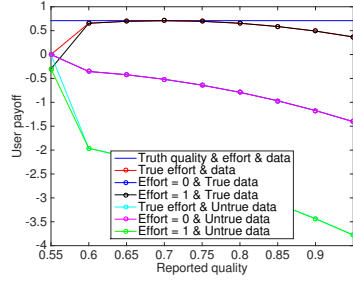


Figure 2: Impact of reported quality  $q'_1$

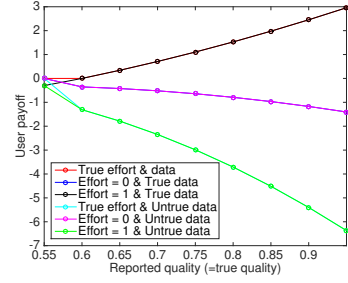


Figure 3: Impact of reported quality  $q'_1$  when it is truthful

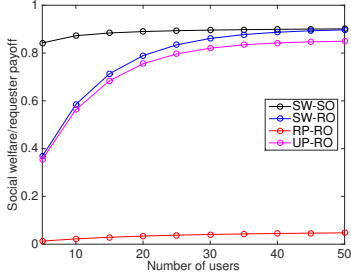


Figure 4: Impact of the number of workers  $n$

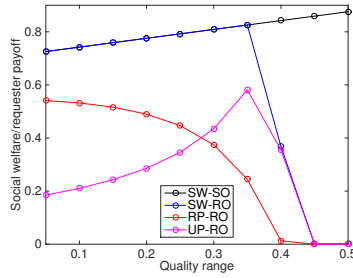


Figure 5: Impact of quality range  $\Delta q$

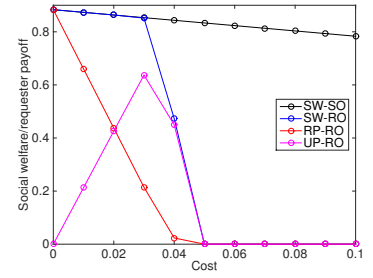


Figure 6: Impact of cost  $c$

## 6 DEALING WITH NO REFERENCE DATA FROM THE REQUESTER

In the previous sections, we have assumed that the requester itself can work on the task and obtains data  $d_0$  with quality  $q_0$  and effort  $e_0 = 1$ , which are (certainly) known by the requester. The reference data  $d_0$  and its quality  $q_0$  are necessary information needed to achieve the truthfulness of the QEDE mechanism. If the requester cannot work on the task (e.g., when it is too far away from the location of interest), we can modify the QEDE mechanism to deal with this situation, described as follows.

For each worker  $i$ , we pick any other worker  $j \neq i$  as a *reference worker*, and define the reward function  $r_i$  as

$$r_i(q', e'_i, d'_i, d'_j) \quad (22)$$

where  $d'_j$  is the data reported by worker  $j$ . We are interested in a mechanism under which truthful behavior of all workers is a *Nash equilibrium*, defined as follows.

**DEFINITION 6.** A mechanism achieves truthful strategies of all workers as a *Nash equilibrium (NE)* if, for each worker  $i$ , given that all other workers  $j \neq i, \forall j$  truthfully report their quality and data, and make the effort desired by the requester, the optimal strategy of worker  $i$  for maximizing its expected payoff is also to the truthful strategy, i.e.,

$$E_{D_j | d_i(q_i, e_i)} [u_i(q_i, \mathbf{q}_{-i}, e_i, d_i, D_j)] \geq E_{D_j | d_i(q_i, e_i)} [u_i(q'_i, \mathbf{q}_{-i}, e'_i, d'_i, D_j)], \forall (q'_i, e_i, d'_i), \forall \mathbf{q}_{-i}.$$

To deal with the lack of reference data  $d_0$  from the requester, we modify the reward function of the QEDE mechanism given in (9) by replacing  $d_0$  with the reported data  $d'_j$  of worker  $i$ 's reference

worker, worker  $j$ , and replacing  $q_0$  with worker  $j$ 's quality  $q'_j$ . To guarantee that each worker  $i$  working on the task (i.e.,  $e_i = 1$ ) has a reference worker  $j \neq i$  also working on the task (i.e.,  $e_j = 1$ ), we need to restrict the assignment function  $e'$  such that there are either at least two workers or no worker working on the task, i.e.,

$$\sum_{i \in N} e'_i(q') \neq 1, \forall q'. \quad (23)$$

The conditions (8) and (10) of the QEDE mechanism remain the same. We can show that the modified QEDE mechanism can achieve an NE where all workers behave truthfully, and also the IR property. The proof follows from the same argument as that of Theorem 1.

## 7 SIMULATION RESULTS

In this section, we evaluate the properties of the QEDE mechanisms and its performance with the RO assignment using simulations.

### 7.1 Worker's payoff

To illustrate the truthfulness of the QEDE mechanisms, we compare a worker's expected payoff when it truthfully reports its quality and data and makes its effort with when it untruthfully reports its quality and/or data and/or makes its effort. We use the SO assignment  $e_i^{so}(q)$  in (19) for the QEDE mechanisms. We set the default parameters as follows<sup>8</sup>:  $n = 2, c = 0.3, \mu_q \triangleq (\bar{q} + \underline{q})/2 = 0.75, \Delta q = 0.4, q_1 = 0.7, q_2 = 0.6$ .

Figs. 2 and 3 illustrate worker 1's expected payoff as it reports varying quality  $q'_1$  (Fig. 2) or varying true quality  $q_1$  (Fig. 3) while making desired effort  $e'_1 = e_1^*(q'_1, q_2)$  or undesired effort  $e'_1 \neq$

<sup>8</sup>It suffices to consider 2 workers only as the RO assignment only depends on the best worker's quality.



$e_1^*(q'_1, q_2)$ , and reporting true data  $d'_1 = d_1$  or untrue data  $d'_1 \neq d_1$ , compared to when it truthfully reports its quality and data and makes its effort. We can see that the worker's payoff when its behavior is untruthful is always less than when truthful. Furthermore, the worker's payoff gap due to untruthfulness often increases when it is more untruthful (i.e., the difference between the reported quality and true quality increases). This confirms that the DIC property is achieved by the QEDE mechanisms so that workers have incentive to behave truthfully. We also observe from Figs. 2-3 that the worker's payoff is always greater than 0 when it behaves truthfully. This confirms that the IR property is achieved by the QEDE mechanisms.

## 7.2 Requester's payoff

To illustrate the system efficiency of the RO assignment, we compare the expected requester's payoff (RP), workers' total payoff (UP), and social welfare (SW) attained by the RO assignment (RP-RO, UP-RO, SW-RO) with the expected social welfare (SW) attained by the SO assignment (SW-SO). Note that by definition, UP-RO is always equal to SW-RO minus RP-RO. We set the default parameters as follows:  $N = 5$ ,  $c = 0.04$ ,  $\mu_q \triangleq (\bar{q} + q)/2 = 0.75$ ,  $\Delta q = 0.4$ . We assume that each worker's quality follows an i.i.d. uniform distribution over  $[q, \bar{q}]$ .

Fig. 4 illustrates the impact of the number of workers  $N$  on the performance. We observe that all the curves are increasing in  $N$ , which is because they benefit from a greater diversity in workers' quality when there are more workers. We also observe that the gap between SW-RO and SW-SO converges to 0 as  $N$  increases, which confirms our result in Proposition 5.

Fig. 5 illustrates the impact of the quality range  $\Delta q$  on the performance. We observe that SW-SO is increasing in  $\Delta q$ . This is because the social welfare benefits from a greater diversity of workers' quality. We also observe that RP-RO is decreasing in  $\Delta q$ . Intuitively, this is due to that a larger quality range requires a higher truth-eliciting payment in (16). We further observe that as  $\Delta q$  increases, the gap between SW-RO and SW-SO is first 0 and then increases. This is because when  $\Delta q$  is small, the RO assignment is always the same as the SO assignment so that SW-RO is always equal to SW-SO; when  $\Delta q$  is greater than some value and increases, the RO assignment more often differs from the SO assignment due to the higher truth-eliciting payment, so that the gap between SW-RO and SW-SO increases.

Fig. 6 illustrates the impact of the cost  $c$  on the performance. We observe that all the curves except for UP-RO are decreasing in  $c$ , which is because a higher cost results in lower social welfare or the requester's payoff. We also observe that RP-RO decreases faster than SW-RO as  $c$  increases. This is because a larger  $c$  not only results in a higher payment for compensating workers' cost, but also a higher truth-eliciting payment in (16) (as discussed in Remark 2). We further observe that as  $c$  increases, the gap between SW-RO and SW-SO is first 0 and then increases. This is because when  $c$  is small, the RO assignment is always the same as the SO assignment; when  $c$  is greater than some value and increases, the RO assignment more often differs from the SO assignment due to the higher truth-eliciting payment.

## 8 CONCLUSION AND FUTURE WORK

In this paper, we have devised the QEDE mechanisms for quality-aware crowdsourcing, to incentivize strategic workers to truthfully report their private quality and data, and make truthful effort as desired by the crowdsourcing requester. The QEDE mechanisms have achieved the truthful design by exploiting the statistical dependency of a worker's private data on its private worker quality and hidden effort, while addressing the coupling in the joint elicitation of quality, effort, and data. Under the QEDE mechanisms, we have characterized the socially optimal and requester's optimal assignments and analyzed their performance, which provide useful insight.

For future work, one interesting direction is to consider workers that have no knowledge of their quality. In this case, the requester needs to learn the quality of strategic workers which may not truthfully provide data to the requester for the purpose of learning. In this paper, we have focused on the truthful elicitation of quality, effort, and data under the assumption that workers' cost is known to the requester. The truthful design when workers' cost is also their private information is still an open problem and will be studied in our future work.

## REFERENCES

- [1] D. R. Karger, S. Oh, and D. Shah, "Budget-optimal crowdsourcing using low-rank matrix approximations," in *IEEE Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, 2011.
- [2] I. Koutsopoulos, "Optimal incentive-driven design of participatory sensing systems," in *IEEE International Conference on Computer Communications (INFOCOM)*, 2013.
- [3] D. Lee, J. Kim, H. Lee, and K. Jung, "Reliable multiple-choice iterative algorithm for crowdsourcing systems," in *ACM International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS)*, 2015.
- [4] N. B. Shah and D. Zhou, "Double or nothing: Multiplicative incentive mechanisms for crowdsourcing," in *Conference on Neural Information Processing Systems (NIPS)*, 2015.
- [5] Y. Liu and M. Liu, "An online learning approach to improving the quality of crowd-sourcing," in *ACM International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS)*, 2015.
- [6] H. Jin, L. Su, D. Chen, K. Nahrstedt, and J. Xu, "Quality of information aware incentive mechanisms for mobile crowd sensing systems," in *ACM International Symposium on Mobile Ad Hoc Networking and Computing (MobiHoc)*, 2015.
- [7] H. Jin, L. Su, H. Xiao, and K. Nahrstedt, "INCEPTION: Incentivizing privacy-preserving data aggregation for mobile crowd sensing systems," in *ACM International Symposium on Mobile Ad Hoc Networking and Computing (MobiHoc)*, 2016.
- [8] H. Jin, L. Su, and K. Nahrstedt, "CENTURION: Incentivizing multi-requester mobile crowd sensing," in *IEEE International Conference on Computer Communications (INFOCOM)*, 2017.
- [9] X. Gong and N. Shroff, "Truthful mobile crowdsensing for strategic users with private qualities," in *International Symposium on Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks (WiOpt)*, 2017.
- [10] L. Duan, T. Kubo, K. Sugiyama, J. Huang, T. Hasegawa, and J. Walrand, "Incentive mechanisms for smartphone collaboration in data acquisition and distributed computing," in *IEEE International Conference on Computer Communications (INFOCOM)*, 2012.
- [11] D. Yang, G. Xue, X. Fang, and J. Tang, "Crowdsourcing to smartphones: Incentive mechanism design for mobile phone sensing," in *ACM Annual International Conference on Mobile Computing and Networking (MobiCom)*, 2012.
- [12] Z. Feng, Y. Zhu, Q. Zhang, L. M. Ni, and A. V. Vasilakos, "Trac: Truthful auction for location-aware collaborative sensing in mobile crowdsourcing," in *IEEE International Conference on Computer Communications (INFOCOM)*, 2014.
- [13] A. Tarable, A. Nordio, E. Leonardi, and M. A. Marsan, "The importance of being earnest in crowdsourcing systems," in *IEEE International Conference on Computer Communications (INFOCOM)*, 2015.
- [14] Y. Luo, N. B. Shah, J. Huang, and J. Walrand, "Parametric prediction from parametric agents," in *The 10th Workshop on the Economics of Networks, Systems and Computation (NetEcon)*, 2015.
- [15] W. Wang, L. Ying, and J. Zhang, "The value of privacy: Strategic data subjects, incentive mechanisms and fundamental limits," in *ACM International Conference*

- on *Measurement and Modeling of Computer Systems (SIGMETRICS)*, 2016.
- [16] L. Pu, X. Chen, J. Xu, and X. Fu, "Crowdlet: Optimal worker recruitment for self-organized mobile crowdsourcing," in *IEEE Conference on Computer Communications (INFOCOM)*, 2016.
- [17] H. Zhang, B. Liu, H. Susanto, G. Xue, and T. Sun, "Incentive mechanism for proximity-based mobile crowd service systems," in *IEEE International Conference on Computer Communications (INFOCOM)*, 2016.
- [18] N. Nisan, T. Roughgarden, E. Tardos, and V. V. Vazirani, *Algorithmic game theory*. Cambridge University Press, 2007, vol. 1.
- [19] P. Bolton and M. Dewatripont, *Contract theory*. MIT press, 2005.
- [20] A. Dasgupta and A. Ghosh, "Crowdsourced judgement elicitation with endogenous proficiency," in *International World Wide Web Conference (WWW)*, 2013.
- [21] Y. Cai, C. Daskalakis, and C. H. Papadimitriou, "Optimum statistical estimation with strategic data sources," in *Conference on Learning Theory (COLT)*, 2015.
- [22] Y. Liu and Y. Chen, "Learning to incentivize: Eliciting effort via output agreement," *International Joint Conference on Artificial Intelligence (IJCAI)*, 2016.
- [23] —, "Sequential peer prediction: Learning to elicit effort using posted prices," in *AAAI Conference on Artificial Intelligence (AAAI)*, 2017.
- [24] D. Prelec, "A Bayesian truth serum for subjective data," *Science*, vol. 306, no. 5695, pp. 462–466, 2004.
- [25] H. Jin, L. Su, and K. Nahrstedt, "Theseus: Incentivizing truth discovery in mobile crowd sensing systems," in *ACM International Symposium on Mobile Ad Hoc Networking and Computing (MobiHoc)*, 2017.
- [26] V. Krishna, *Auction theory*. Academic press, 2009.
- [27] "Incentivizing truthful data quality for quality-aware mobile data crowdsourcing," Technical Report, 2017. [Online]. Available: <https://www.dropbox.com/s/y99bim7hbaxht4j/crowdsourcing-quality-mobihoc-TR.pdf?dl=0>

## APPENDIX

### Proof of Lemma 1

Let  $\bar{d}_i$  be the complementary value of  $d_i$ , i.e.,  $\bar{d}_i \neq d_i$ . For convenience, let  $P_{X|d_i(q_i, e_i)}(d)$  denote the probability that  $X$  is equal to  $d$  conditioned on data  $d_i$  given quality  $q_i$  and effort  $e_i$ . We observe that when  $e_i = 1$  we have

$$P_{X|d_i(q_i, 1)}(d_i) = q_i \geq 1 - q_i = P_{X|d_i(q_i, 1)}(\bar{d}_i),$$

and when  $e_i = 0$  we have

$$P_{X|d_i(q_i, 0)}(d_i) = 0.5 \geq 1 - 0.5 = P_{X|d_i(q_i, 0)}(\bar{d}_i).$$

Since

$$\begin{aligned} E_{D_0|d_i(q_i, e_i)} \left[ \mathbf{1}_{D_0=d_i} \right] &= P_{D_0|d_i(q_i, e_i)}(d_i) \\ &= q_0 P_{X|d_i(q_i, e_i)}(d_i) + (1 - q_0)(1 - P_{X|d_i(q_i, e_i)}(d_i)) \\ &= (2q_0 - 1)P_{X|d_i(q_i, e_i)}(d_i) + 1 - q_0, \end{aligned}$$

we have

$$E_{D_0|d_i(q_i, e_i)} \left[ \frac{\mathbf{1}_{D_0=d_i} + q_0 - 1}{2q_0 - 1} \right] = P_{X|d_i(q_i, e_i)}(d_i).$$

Similarly, we can show that

$$E_{D_0|d_i(q_i, e_i)} \left[ \frac{\mathbf{1}_{D_0=\bar{d}_i} + q_0 - 1}{2q_0 - 1} \right] = P_{X|d_i(q_i, e_i)}(\bar{d}_i).$$

Then it follows from (9) that, for any reported quality  $q'_i$  and any actual effort  $e_i$ , the optimal reported data is given by

$$\begin{aligned} d'_i &= \arg \max_{d \in \{0, 1\}} E_{D_0|d_i(q_i, e_i)} \left[ \frac{\mathbf{1}_{D_0=d} + q_0 - 1}{2q_0 - 1} \right] \\ &= \arg \max_{d \in \{0, 1\}} P_{X|d_i(q_i, e_i)}(d) \\ &= \arg \max_{d \in \{0, 1\}} \left[ P_{X|d_i(q_i, 0)}(d) + (P_{X|d_i(q_i, 1)}(d) - P_{X|d_i(q_i, 0)}(d))e_i \right] \\ &= d_i. \end{aligned}$$

### Proof of Lemma 2

Using (12), when  $e'_i = 1$  we have

$$\bar{u}_i(q'_i, q_i, 1, 1) - \bar{u}_i(q'_i, q_i, 1, 0) = kq'_i(q_i - 0.5) - c \geq 0$$

where the inequality follows from (10), and when  $e'_i = 0$  we have

$$\bar{u}_i(q'_i, q_i, 0, 0) - \bar{u}_i(q'_i, q_i, 0, 1) = c \geq 0.$$

Hence the optimal effort to make is  $e_i = e'_i$ .

### Proof of Lemma 3

For convenience, we write  $\hat{u}_i(q', q_i, e'_i)$  as  $\hat{u}_i(q'_i, q'_{-i}, q_i, e'_i)$ . It suffices to show that  $\hat{u}_i(q_i, q'_{-i}, q_i, e'_i) \geq \hat{u}_i(q', q'_{-i}, q_i, e'_i)$ ,  $\forall q' \neq q_i$ . Let  $q'_i > q_i$ . Using (14), we have

$$\begin{aligned} &\hat{u}_i(q_i, q'_{-i}, q_i, e'_i) - \hat{u}_i(q'_i, q'_{-i}, q_i, e'_i) \\ &= kq_i e'_i(q_i, q'_{-i})(q_i - q_i) + \int_{\underline{q}}^{q_i} kq e'_i(q, q'_{-i}) dq \\ &\quad - \left( kq'_i e'_i(q_i, q'_{-i})(q_i - q'_i) + \int_{\underline{q}}^{q'_i} kq e'_i(q, q'_{-i}) dq \right) \\ &= kq'_i e'_i(q_i, q'_{-i})(q'_i - q_i) - \int_{q_i}^{q'_i} kq e'_i(q, q'_{-i}) dq \\ &\geq kq'_i e'_i(q_i, q'_{-i})(q'_i - q_i) - kq'_i e'_i(q_i, q'_{-i})(q'_i - q_i) = 0 \end{aligned}$$

where the inequality follows from (8). Now let  $q'_i < q_i$ . Using (14), we have

$$\begin{aligned} &\hat{u}_i(q_i, q'_{-i}, q_i, e'_i) - \hat{u}_i(q'_i, q'_{-i}, q_i, e'_i) \\ &= kq_i e'_i(q_i, q'_{-i})(q_i - q_i) + \int_{\underline{q}}^{q_i} kq e'_i(q, q'_{-i}) dq \\ &\quad - \left( kq'_i e'_i(q_i, q'_{-i})(q_i - q'_i) + \int_{\underline{q}}^{q'_i} kq e'_i(q, q'_{-i}) dq \right) \\ &= kq'_i e'_i(q_i, q'_{-i})(q'_i - q_i) + \int_{q'_i}^{q_i} kq e'_i(q, q'_{-i}) dq \\ &\geq kq'_i e'_i(q_i, q'_{-i})(q'_i - q_i) + kq'_i e'_i(q_i, q'_{-i})(q_i - q'_i) = 0 \end{aligned}$$

where the inequality follows from (8).

### Proof of Theorem 1

As the IR property has been proved using (15), we only show that the DIC property is achieved. Choose and fix any  $(q'_i, e_i, d'_i)$ . It follows from Lemma 1 that

$$E_{D_0} \left[ u_i(q'_i, q'_{-i}, e_i, d_i, D_0) \right] \geq E_{D_0} \left[ u_i(q'_i, q'_{-i}, e_i, d'_i, D_0) \right].$$

Using (11) and (12), it follows from Lemma 2 that

$$\bar{u}_i(q', q'_i, e'_i, e'_i) \geq \bar{u}_i(q', q_i, e'_i, e_i).$$

Using (13) and (14), it follows from Lemma 3 that

$$\hat{u}_i(q', q_i, e'_i) \geq \hat{u}_i(q', q'_i, e'_i).$$

Therefore, we have

$$E_{D_0} \left[ u_i(q_i, q'_{-i}, e'_i, d_i, D_0) \right] \geq E_{D_0} \left[ u_i(q'_i, q'_{-i}, e_i, d'_i, D_0) \right].$$