

Incentives for Privacy Tradeoff in Community Sensing

Adish Singla and Andreas Krause

ETH Zurich

Universitätsstrasse 6, 8092 Zürich, Switzerland

Abstract

Community sensing, fusing information from populations of privately-held sensors, presents a great opportunity to create efficient and cost-effective sensing applications. Yet, reasonable privacy concerns often limit the access to such data streams. How should systems value and negotiate access to private information, for example in return for monetary incentives? How should they optimally choose the participants from a large population of strategic users with privacy concerns, and compensate them for information shared?

In this paper, we address these questions and present a novel mechanism, SEQTGREEDY, for budgeted recruitment of participants in community sensing. We first show that privacy tradeoffs in community sensing can be cast as an adaptive submodular optimization problem. We then design a budget feasible, incentive compatible (truthful) mechanism for adaptive submodular maximization, which achieves near-optimal utility for a large class of sensing applications. This mechanism is general, and of independent interest. We demonstrate the effectiveness of our approach in a case study of air quality monitoring, using data collected from the Mechanical Turk platform. Compared to the state of the art, our approach achieves up to 30% reduction in cost in order to achieve a desired level of utility.

Introduction

Community sensing is a new paradigm for creating efficient and cost-effective sensing applications by harnessing the data of large populations of sensors. For example, the accelerometer data from smartphone users could be used for earthquake detection and fine grained analysis of seismic events. Velocity data from GPS devices (in smartphones or automobiles) could be used to provide real-time traffic maps or detect accidents. However, accessing this stream of private sensor data raises reasonable concerns about privacy of the individual users. For example, mobility patterns and the house or office locations of a user could possibly be inferred from their GPS tracks (Krumm 2007). Beyond concerns about sharing sensitive information, there are general anxieties among users about sharing data from their private devices. These concerns limit the practical applicability

of deploying such applications. In this paper, we propose a principled approach to negotiate access to certain private information in an incentive-compatible manner.

Applications of community sensing are numerous. Several case studies have demonstrated the principal feasibility and usefulness of community sensing. A number of research and commercial prototypes are built, often relying on special campaigns to recruit volunteers (Zheng, Xie, and Ma 2010) or on contracts with service providers to obtain anonymized data (Wunnava et al. 2007). The SenseWeb system (Kansal et al. 2007) has been developed as an infrastructure for sharing sensing data to enable various applications. Methods have been developed to estimate traffic (Yoon, Noble, and Liu 2007; Mobile-Millennium 2008; Krause et al. 2008), perform forecasts about future traffic situations (Horvitz et al. 2005) or predict a driver's trajectory (Krumm and Horvitz 2006). Cell tower signals obtained from the service providers are leveraged for travel time estimation on roadways (Wunnava et al. 2007). Additionally, captured images and video clips from smartphones have been used to link places with various categories (Chon et al. 2012). Clayton et al. (2012) describes the design of a *Community Seismic Network* to detect and monitor earthquakes using a dense network of low cost sensors hosted by volunteers from the community. Aberer et al. (2010) envisions a community driven sensing infrastructure for monitoring air quality.

Privacy concerns in community sensing are expected and reasonable (Lieb 2007; Wunnava et al. 2007; Olson, Grudin, and Horvitz 2005). Irrespective of the models of privacy we consider (Sweeney 2002; Dwork 2006; Machanavajjhala et al. 2006), the key concern is about identifiability as users become members of increasingly smaller groups of people sharing the same characteristics inferred from data. Beyond general anxieties about the sharing of location and mobility data, studies have demonstrated that, even with significant attempts at obfuscation, home and work locations of drivers can be inferred from GPS tracks (Krumm 2007).

Incentives to participants for privacy tradeoff. Olson, Grudin, and Horvitz (2005) show that people's willingness to share information depends greatly on the type of information being shared, with whom the information is shared, and how it is going to be used. They are willing to share certain private information if compensated in terms of their utility gain (Krause and Horvitz 2008). In this paper, we are ex-

ploring the design of intelligent systems that empower users to consciously share certain private information in return of, *e.g.*, monetary or other form of incentives. We model the users as strategic agents who are willing to negotiate access to certain private information, aiming to maximize the monetary incentives they receive in return. Empowering users to opt into such negotiations is the key idea that we explore in this paper.

Overview of our approach

Our goal is to design policies for selecting (and compensating) the participants, which provide near-optimal utility for the sensing application under strict budget constraints. As basis for selection, the community sensing system receives obfuscated estimates of the private attributes. For concreteness, we focus on *sensor location* as private information, but our approach generalizes to other attributes. The users also declare a bid or cost as the desired monetary incentive for participation and hence privacy tradeoff. After receiving the bids, the mechanism sequentially selects a participant, commits to make her the payment, receives the actual private information, selects the next participant and so on. At the end, all selected participants are provided the agreed payment. Figure 1 illustrates this protocol.

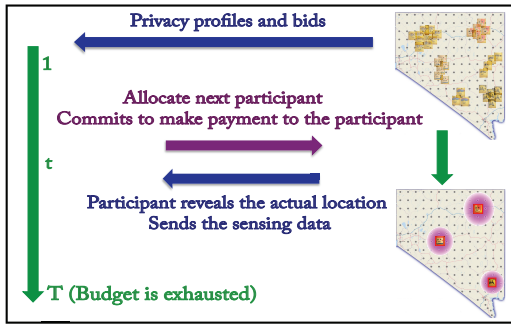


Figure 1: Illustration of the protocol by which the proposed system interacts with the users.

We model the participants as strategic agents who aim to maximize their profit, by possibly misreporting their private costs. As a consequence, we require the mechanism to be truthful. In order to capture a large class of sensing applications, we only require the utility function to satisfy submodularity, a natural diminishing returns condition (Nemhauser, Wolsey, and Fisher 1978; Krause and Guestrin 2007). To design our mechanism, we first reduce the sequential negotiation of the privacy tradeoff to the problem of adaptive submodular maximization (Asadpour, Nazerzadeh, and Saberi 2008; Golovin and Krause 2011). Then, we extend recent results on truthful budget feasible mechanisms for submodular functions (Singer 2010; Chen, Gravin, and Lu 2011; Singer 2012) to the adaptive setting.

Our main contributions are:

- An integrated approach to community sensing by incentivizing users to share certain private information.
- A novel mechanism, SEQTGREEDY, for budgeted recruitment of strategic participants, which achieves near-optimal utility for the community sensing application. The mechanism is general and of independent interest,

suitable also for other applications, *e.g.*, viral marketing.

- Evaluation of our approach on a realistic case study of air quality monitoring based on data obtained through Amazon Mechanical Turk ¹.

Related Work

Himmel et al. (2005) propose to provide users with rewards such as free minutes to motivate them to accept mobile advertisements. Hui et al. (2011) develop MobiAd, a system for targeted mobile advertisements, by utilizing the rich set of information available on the phone and suggesting the service providers to give discounts to the users, in order to incentivize use of the system. Liu, Krishnamachari, and Annaram (2008) propose a game theoretic model of privacy for social networking-based mobile applications and presents a tit-for-tat mechanism by which users take decisions about their exposed location obfuscation for increasing personal or social utility. Chorpapath and Alpcan (2012) study a privacy game in mobile commerce, where users choose the degree of granularity at which to report their location and the service providers offer them monetary incentives under budget constraints. The best users' response and the optimal strategy for the company are derived by analyzing the Nash equilibrium of the underlying privacy game. This is very different from our setting as we focus on algorithmic aspects of the mechanism in choosing the best set of users for participation in community sensing. Li and Faltings (2012) and Faltings, Jurca, and Li (2012) study the problem of incentivizing users in community sensing to report accurate measurements and place sensors in the most useful locations. While developing incentive-compatible mechanisms, they do not consider the privacy aspect. Singla and Krause (2013b) develops online incentive-compatible and budget feasible mechanisms for procurement. However, they consider a simple modular utility function where each participant provides a unit value. This is not applicable to our community sensing setting which deals with more complex utility functions. Carrascal et al. (2013) study how users value their personally identifiable information (PII) while browsing. The experiments demonstrate that users have different valuations, depending on the type and information content of private data. Higher valuations are chosen for offline PII, such as age and address, compared to browsing history. This work is complementary and supports the assertion that users indeed associate monetary valuations to certain private data.

Problem Statement

We now formalize the problem addressed in this paper.

Sensing phenomena. We focus on community sensing applications with the goal to monitor some spatial phenomenon, such as air quality or traffic. We discretize the environment as a finite set of locations \mathcal{V} , where each $v \in \mathcal{V}$ could, *e.g.*, denote a zip code or more fine grained street addresses, depending on the application. We quantify the utility $f(\mathcal{A})$ of obtaining measurements from a set of locations \mathcal{A} using a set function $f : 2^{\mathcal{V}} \rightarrow \mathbb{R}$. Formally, we

¹<https://www.mturk.com/mturk/>

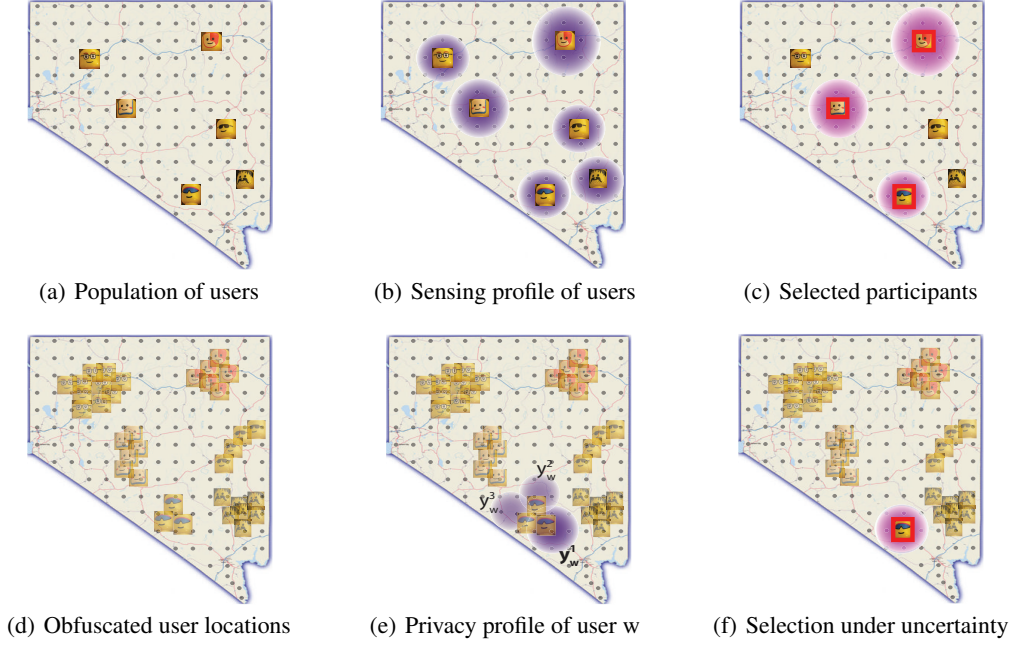


Figure 2: The sensing region is uniformly discretized into a set of locations \mathcal{V} indicated by the dots. (a) illustrates a population of users, along with their sensing profiles in (b). The set of users selected by the system in absence of privacy are shown in (c). However, to protect privacy, users only share an obfuscated location with the system in (d) and a collection of sensing profiles ($\{y_w^1, y_w^2$ and $y_w^3\}$ for user w) in (e). The privacy profile of user w , given by Y_w , is the uniform distribution over these sensing profiles, given by $P(Y_w = y_w^i) = \frac{1}{3}$. (f) shows the selection of the participants in presence of uncertainty introduced by privacy profiles. The actual sensing profile is only revealed to the system after a user has been selected.

only require that f is *nonnegative*, *monotone* (i.e., whenever $\mathcal{A} \subseteq \mathcal{A}' \subseteq \mathcal{V}$ it holds that $f(\mathcal{A}) \leq f(\mathcal{A}')$) and *submodular*. Submodularity is an intuitive notion of diminishing returns, stating that, for any sets $\mathcal{A} \subseteq \mathcal{A}' \subseteq \mathcal{V}$, and any fixed location $a \notin \mathcal{A}'$ it holds that $f(\mathcal{A} \cup \{a\}) - f(\mathcal{A}) \geq f(\mathcal{A}' \cup \{a\}) - f(\mathcal{A}')$. As a simple, concrete example, we may derive some nonnegative value d_a for observing each location $a \in \mathcal{A}$, and may define $f(\mathcal{A}) = \sum_{a \in \mathcal{A}} d_a$. More generally, sensing at location $a \in \mathcal{V}$ may actually cover a subset \mathcal{S}_a of nearby locations, and $f(\mathcal{A}) = \sum \{d_j : j \in \cup_{a \in \mathcal{A}} \mathcal{S}_a\}$. These conditions are rather general, satisfied by many sensing utility functions and f can capture much more complex notions, such as reduction of predictive uncertainty in a probabilistic model (Krause and Guestrin 2007).

Sensing profile of users. We consider a community \mathcal{W} of $|\mathcal{W}| = N$ users, owning some sensing device such as a smartphone. Each user can make observations at a set of locations depending on her geolocation or mobility as well as the type of device used. We model this through a collection of *sensing profiles* $\mathcal{O} \subseteq 2^{\mathcal{V}}$ whereby we associate each user $w \in \mathcal{W}$ with a profile $y_w \in \mathcal{O}$, specifying the set of locations covered by her. This set y_w could be a singleton $y_w = \{a\}$ for some $a \in \mathcal{V}$, modeling the location of the user at a particular point in time, or could model an entire trajectory, visiting multiple locations in \mathcal{V} . We denote a given set of users $\mathcal{S} \subseteq \mathcal{W}$ jointly with their sensing profiles as $\mathbf{y}_{\mathcal{S}} \subseteq \mathcal{W} \times \mathcal{O}$. The goal is to select set of users \mathcal{S} (also called *participants*) so as to maximize the utility of the sensing application given by $g(\mathbf{y}_{\mathcal{S}}) = f(\mathcal{A})$ where $\mathcal{A} = \bigcup_{s \in \mathcal{S}} y_s$. We

assume that each user's maximal contribution to the utility is bounded by a constant f_{\max} .

Privacy profile of users. In order to protect privacy, we consider the setting where the exact sensing profiles y_w of the users (containing, e.g., tracks of locations visited) are not known to the sensing system. Instead, y_w is only shared after obfuscation with a random perturbation intended to reduce the risk of identifiability (Sweeney 2002; Dwork 2006). The system's highly uncertain belief about the sensing profile of user w can therefore be represented as a (set-valued) random variable (also called *privacy profile*) Y_w with y_w being its realization. For example, suppose $y_w = \{a\}$ for some location a (i.e., the user's private location is $a \in \mathcal{V}$). In this case, the user may share with the system a collection of locations a_1, \dots, a_m containing a (but not revealing which one it is), w.l.o.g. $a = a_1$. In this case the distribution shared $P(Y_w = \{a_i\}) = \frac{1}{m}$ is simply the uniform distribution over the candidate locations. Figure 2 illustrates the notions of sensing and privacy profiles for a user.

We use $\mathbf{Y}_{\mathcal{W}} = [Y_1, \dots, Y_N]$ to refer to the collection of all (independent) variables associated with population \mathcal{W} and assume that $\mathbf{Y}_{\mathcal{W}}$ is distributed according to a factorial joint distribution $P(\mathbf{Y}_{\mathcal{W}}) = \prod_w P(Y_w)$. The sensing profile y_w (and the actual sensor data obtained from sensing at locations y_w) is revealed to the application only after it commits to provide the desired incentives to the user w . Then, the goal is to select a set of users \mathcal{S} to maximize $\mathbb{E}_{\mathbf{Y}_{\mathcal{W}}} [g(\mathbf{y}_{\mathcal{S}})]$, i.e., the expected utility, where the expectation is taken over the realizations of $\mathbf{Y}_{\mathcal{W}}$ w.r.t. $P(\mathbf{Y}_{\mathcal{W}})$.

Incentive structure for privacy tradeoff. We assume that users are willing to share certain non-sensitive private information in return for monetary incentives. Each user w has a private cost $c_w \in \mathbb{R}_{\geq 0}$ that she experiences for her privacy tradeoff. Instead of revealing c_w , she only reveals a *bid* $b_w \in \mathbb{R}_{\geq 0}$. We are interested in *truthful* mechanisms, where it is a dominant strategy for a user to report $b_w = c_w$, *i.e.*, users cannot increase their profit (in *expectation*) by lying about their true cost. We assume that costs have known bounded support, *i.e.*, $c_w \in [c_{\min}, c_{\max}]$.

Optimization problem. Given a strict budget constraint \mathcal{B} , the goal of the sensing application is to design a mechanism \mathcal{M} , which implements an allocation policy to select participants \mathcal{S} and a payment scheme to make *truthful* payments θ_s to each of the participants, with the goal of maximizing the expected utility. Instead of committing to a fixed set of participants \mathcal{S} in advance (*non-adaptive* policy), we are interested in mechanisms that implement an *adaptive* policy taking into account the observations made so far (revealed sensing profiles of participants already selected) when choosing the next user. Formally, the goal of the mechanism is to adaptively select participants \mathcal{S}^* along with the payments $\theta_{\mathcal{S}^*}$, such that

$$\mathcal{S}^* = \arg \max_{\mathcal{S} \subseteq \mathcal{W}} \mathbb{E}_{\mathbf{Y}_{\mathcal{W}}} [g(\mathbf{y}_{\mathcal{S}})] \text{ subject to } \sum_{s \in \mathcal{S}} \theta_s \leq \mathcal{B}. \quad (1)$$

Here, the set of participants \mathcal{S} selected and the payments $\theta_{\mathcal{S}}$ may depend on the realization of $\mathbf{Y}_{\mathcal{W}}$ as well. We formally introduce adaptive policies in subsequent sections.

Existing Mechanisms

We first review existing mechanisms that fall short of either privacy-preservation, adaptivity or truthfulness. In next section, we then build on these and present our main contribution: a privacy-respecting, truthful and adaptive mechanism.

Non-private mechanisms

Consider first an unrealistic setting, where the system has full information about the users' exact sensing profiles and their true costs. In such a setting, Problem 1 reduces to that of budgeted maximization of a monotone non-negative submodular function with non-uniform costs, studied by Sviridenko (2004). A simple algorithm combining partial enumeration with greedy selection guarantees a utility of at least $(1 - 1/e)$ ($= 0.63$) times that obtained by optimal selection OPT. This result is tight under reasonable complexity assumptions (Feige 1998). We denote this setting and mechanism as GREEDY. Note that each participant is paid their true cost in this untruthful setting. Now, consider the non-private setting with *unknown* true costs. The problem then requires designing a truthful budget feasible mechanism for monotone submodular set functions, as done by (Singer 2010; Chen, Gravin, and Lu 2011; Singer 2012). In this setting, a constant factor $1/7.91$ ($= 0.13$) approximation compared to OPT can be achieved, using a mechanism that we will refer to as TGREEDY. TGREEDY executes a greedy allocation on a reduced budget with carefully chosen stopping criteria (for ensuring budget

	Untruthful	Truthful
Priv. off	GREEDY	TGREEDY
Priv. on (Non-Ad.)	CONSTGREEDY	CONSTTGREEDY
Priv. on (Adaptive)	SEQGREEDY	SEQTGREEDY

Table 1: Different information settings and mechanisms.

feasibility), in order to select a set of participants and then computes the truthful payments to be made to them.

Non-adaptive mechanisms with privacy

In our case, where privacy is preserved through random obfuscation, one must deal with the stochasticity caused by the uncertainty about users' sensing profiles. Here, the objective

$$G(\mathcal{S}) \equiv \mathbb{E}_{\mathbf{Y}_{\mathcal{W}}} [g(\mathbf{y}_{\mathcal{S}})] = \sum_{\mathbf{y}_{\mathcal{W}}} P(\mathbf{Y}_{\mathcal{W}} = \mathbf{y}_{\mathcal{W}}) f \left(\bigcup_{s \in \mathcal{S}} y_s \right)$$

in (1) can be seen as an expectation over multiple submodular set functions, one for each realisation of the privacy profile variables $\mathbf{Y}_{\mathcal{W}}$. However, as submodularity is preserved under expectations, the set function $G(\mathcal{S})$ is submodular as well. One can therefore still apply the mechanisms GREEDY and TGREEDY in order to obtain near-optimal *non-adaptive* solutions (*i.e.*, the set of participants is fixed in advance) to Problem (1). We denote these non-adaptive (constant) mechanisms applied to our privacy-preserving setting as CONSTGREEDY and CONSTTGREEDY.

Untruthful, adaptive mechanisms with privacy

Instead of non-adaptively committing to the set \mathcal{S} of participants a priori, one may wish to obtain increased utility through adaptive (active/sequential) selection, *i.e.*, by taking into account the observations from the users selected so far when choosing the next user. Without assumptions, computing such an optimal policy for Problem (1) is intractable. Fortunately, as long as the sensing quality function f is monotone and submodular, Problem (1) satisfies a natural condition called *adaptive submodularity* (Golovin and Krause 2011). This condition generalizes the classical notion of submodularity to sequential decision / active selection problems as faced here.

Adaptive submodularity requires, in our setting, that the expected benefit of any fixed user $w \in \mathcal{W}$ given a set of observations (*i.e.*, set of users and observed sensing profiles) can never increase as we make more observations. Formally, consider the *conditional expected marginal gain* of adding a user $w \in \mathcal{W} \setminus \mathcal{S}$ to an existing set of observations $\mathbf{y}_{\mathcal{S}} \subseteq \mathcal{W} \times \mathcal{O}$:

$$\begin{aligned} \Delta_g(w|\mathbf{y}_{\mathcal{S}}) &= \mathbb{E}_{Y_w} [g(\mathbf{y}_{\mathcal{S}} \cup \{(w, y_w)\}) - g(\mathbf{y}_{\mathcal{S}}) | \mathbf{y}_{\mathcal{S}}] \\ &= \sum_{y \in \mathcal{O}} P(Y_w = y | \mathbf{y}_{\mathcal{S}}) \cdot [g(\mathbf{y}_{\mathcal{S}} \cup \{(w, y)\}) - g(\mathbf{y}_{\mathcal{S}})]. \end{aligned}$$

Function g with distribution $P(\mathbf{Y}_{\mathcal{W}})$ is *adaptive submodular*, if $\Delta_g(w|\mathbf{y}_{\mathcal{S}}) \geq \Delta_g(w|\mathbf{y}_{\mathcal{S}'})$ whenever $\mathbf{y}_{\mathcal{S}} \subseteq \mathbf{y}_{\mathcal{S}'}$. Thus, the gain of a user w , in expectation over its unknown privacy profile, can never increase as we select and obtain data from more participants.

Proposition 1. Suppose f is monotone and submodular. Then the objective g and distribution P used in Problem 1 are adaptive submodular.

Above Proposition follows from Theorem 6.1 of Golovin and Krause (2011), assuming distribution P is factorial (i.e., the random obfuscation is independent between users). Given this problem structure, for the simpler, untruthful setting (i.e., *known* true costs), we can thus use the sequential greedy policy for stochastic submodular maximization studied by Golovin and Krause (2011). This approach is denoted by SEQGREEDY and obtains a utility of at least $(1 - 1/e)$ ($= 0.63$) times that of optimal sequential policy SEQOPT.

Table 1 summarizes the settings and mechanisms considered so far. They all fall short of at least one of the desired characteristics of privacy-preservation, truthfulness or adaptivity. In the next section, we present our main contribution – SEQTGREEDY, an adaptive mechanism for the realistic setting of privacy-sensitive and strategic agents.

Our main mechanism: SEQTGREEDY

We now describe our mechanism $\mathcal{M} = (\pi_{\mathcal{M}}, \theta_{\mathcal{M}})$, with allocation policy $\pi_{\mathcal{M}}$ and payment scheme $\theta_{\mathcal{M}}$. \mathcal{M} first obtains the bids $B_{\mathcal{W}}$ and privacy profiles $P(\mathbf{Y}_{\mathcal{W}})$ from all users, runs the allocation policy $\pi_{\mathcal{M}}$ to adaptively select participants \mathcal{S} and makes observations $\mathbf{y}_{\mathcal{S}}$ during selection. At the end, it computes payments $\theta_{\mathcal{S}}$ using scheme $\theta_{\mathcal{M}}$. The allocation policy $\pi_{\mathcal{M}}$ can be thought of as a decision tree. Formally, a policy $\pi : 2^{\mathcal{W} \times \mathcal{O}} \rightarrow \mathcal{W}$ is a partial mapping from observations $\mathbf{y}_{\mathcal{S}}$ made so far to the next user $w \in \mathcal{W} \setminus \mathcal{S}$ to be recruited, denoted by $\pi(\mathbf{y}_{\mathcal{S}}) = w$. We seek policies that are provably competitive with the optimal (intractable) sequential policy SEQOPT. $\theta_{\mathcal{M}}$ computes payments which are truthful in expectation (a user cannot increase her total expected profit by lying about her true cost, for a fixed set of bids of other users) and individually rational ($\theta_s \geq b_s$). For budget feasibility, the allocation policy needs to ensure that the budget \mathcal{B} is sufficient to make the payments $\theta_{\mathcal{S}}$ to all selected participants. Next, we describe in detail the allocation policy and payment scheme of SEQTGREEDY with these desirable properties.

Allocation policy of SEQTGREEDY

Policy 1 presents the allocation policy of SEQTGREEDY. The main ingredient of the policy is to greedily pick the next user that maximizes the expected marginal gain $\Delta_g(w|\mathbf{y}_{\mathcal{S}})$ per unit cost. The policy uses additional stopping criteria to enforce budget feasibility, similar to TGREEDY (Chen, Gravin, and Lu 2011). Firstly, it runs on a reduced budget \mathcal{B}/α . Secondly, it uses a proportional share rule ensuring that the expected marginal gain per unit cost for the next potential participant is at least equal to or greater than the expected utility of the new set of participants divided by the budget. We shall prove below that $\alpha = 2$ achieves the desired properties.

Payment characterization of SEQTGREEDY

The payment scheme is based on the characterization of threshold payments used by TGREEDY (Singer 2010). How-

Policy 1: Allocation policy of SEQTGREEDY

```

1 Input: budget  $\mathcal{B}$ ; users  $\mathcal{W}$ ; privacy profiles  $\mathbf{Y}_{\mathcal{W}}$ ;
   bids  $B_{\mathcal{W}}$ ; reduced budget factor  $\alpha$ ;
2 Initialize:
   • Outputs: participants  $\mathcal{S} \leftarrow \emptyset$ ; observations  $\mathbf{y}_{\mathcal{S}} \leftarrow \emptyset$ ;
     marginals  $\Delta_{\mathcal{S}} \leftarrow \emptyset$ ;
   • Variables: remaining users  $\mathcal{W}' \leftarrow \mathcal{W}$ ;
3 begin
4   while  $\mathcal{W}' \neq \emptyset$  do
5      $w^* \leftarrow \arg \max_{w \in \mathcal{W}'} \frac{\Delta_g(w|\mathbf{y}_{\mathcal{S}})}{b_w}$ ;
6      $\Delta_{w^*} \leftarrow \Delta_g(w^*|\mathbf{y}_{\mathcal{S}})$ ;
7     if  $B_{\mathcal{S}} + b_{w^*} \leq \mathcal{B}$  then
8       if  $b_{w^*} \leq \frac{\mathcal{B}}{\alpha} \cdot \frac{\Delta_{w^*}}{(\sum_{s \in \mathcal{S}} \Delta_s) + \Delta_{w^*}}$  then
9          $\mathcal{S} \leftarrow \mathcal{S} \cup \{w^*\}$ ;  $\Delta_{\mathcal{S}} \leftarrow \Delta_{\mathcal{S}} \cup \{\Delta_{w^*}\}$ ;
10        Observe  $y_{w^*}$ ;  $\mathbf{y}_{\mathcal{S}} \leftarrow \mathbf{y}_{\mathcal{S}} \cup \{(w^*, y_{w^*})\}$ ;
11         $\mathcal{W}' \leftarrow \mathcal{W}' \setminus \{w^*\}$ ;
12      else
13         $\mathcal{W}' \leftarrow \emptyset$ ;
14      else
15         $\mathcal{W}' \leftarrow \mathcal{W}' \setminus \{w^*\}$ ;
16 Output:  $\mathcal{S}$ ;  $\mathbf{y}_{\mathcal{S}}$ ;  $\Delta_{\mathcal{S}}$ 

```

ever, a major difficulty arises from the fact that the computation of payments for a participant depends also on the unallocated users, whose sensing profiles are not known to the mechanism. Let \mathcal{S} denote the set of participants allocated by $\pi_{\mathcal{M}}$ along with making observations $\mathbf{y}_{\mathcal{S}}$. Let us consider the set of all possible realizations of $\mathbf{Y}_{\mathcal{W}} = \mathbf{y}_{\mathcal{W}} \subseteq \mathcal{W} \times \mathcal{O}$ consistent with $\mathbf{y}_{\mathcal{S}}$, i.e., $\mathbf{y}_{\mathcal{S}} \subseteq \mathbf{y}_{\mathcal{W}}$. We denote this set by $\mathbf{Z}_{\mathcal{W}, \mathcal{S}} = [\mathbf{y}^1, \mathbf{y}^2 \dots \mathbf{y}^r \dots \mathbf{y}^Z]$, where $Z = |\mathbf{Z}_{\mathcal{W}, \mathcal{S}}|$. We first discuss how to compute the payment for each one of these possible realizations $\mathbf{y}^r \in \mathbf{Z}_{\mathcal{W}, \mathcal{S}}$, denoted by $\theta_s^d(\mathbf{y}^r)$ (where d indicates here an association with the deterministic setting of knowing the exact sensing profiles of all users $w \in \mathcal{W}$). These payments for specific realizations are then combined together to compute the final payment to each participant.

Payment θ_s^d for a given $\mathbf{y}_{\mathcal{W}}$. Consider the case where the variables $\mathbf{Y}_{\mathcal{W}}$ are in state $\mathbf{y}_{\mathcal{W}} \in \mathbf{Z}_{\mathcal{W}, \mathcal{S}}$ and let \mathcal{S} be the set of participants allocated by the policy. We use the well-known characterization of Myerson (1981) of truthful payments in single-parameter domains. It states that a mechanism is truthful if *i*) the allocation rule is monotone (i.e., an already allocated user cannot be unallocated by lowering her bid, for a fixed set of bids of others) and *ii*) allocated users are paid threshold payments (i.e., the highest bid they can declare before being removed from the allocated set). Monotonicity follows naturally from the greedy allocation policy, which sorts users based on expected marginal gain per unit cost. To compute threshold payments, we need to consider a maximum of all the possible bids that a user can declare and still get allocated. We next explain how this can be done.

Let us renumber the users $\mathcal{S} = \{1, \dots, i, \dots, k\}$ in the order of their allocation. and let us analyze the payment for participant $s = i$. Consider running the policy on an alternate set $\mathcal{W}' = \mathcal{W} \setminus \{i\}$ and let $\mathcal{S}' = \{1, \dots, j, \dots, k'\}$ be the allocated set (users renumbered again based on order of al-

location when running the policy on \mathcal{W}'). Δ_S and $\Delta'_{S'}$ are the marginal contributions of the participants in the above two runs of the policy. We define $\Delta_{i(j)}$ to be the marginal contribution of i (from S) if it has to replace the position of j (in set S'). Now, consider the bid that i can declare to replace j in S' by making a marginal contribution per cost higher than j , given by $b_{i(j)} = \frac{\Delta_{i(j)} \cdot b_j}{\Delta'_j}$. Additionally, the bid that i can declare must satisfy the proportional share rule, denoted by $\rho_{i(j)} = \frac{B}{\alpha} \cdot \Delta_{i(j)} / ((\sum_{s' \in [j-1]} \Delta'_{s'}) + \Delta_{i(j)})$. By taking the minimum of these two values, we get $\theta_{i(j)}^d = \min(b_{i(j)}, \rho_{i(j)})$ as the bid that i can declare to replace j in S' . The threshold payment for participant $s = i$ is given by $\theta_i^d = \max_{j \in [k'+1]} \theta_{i(j)}^d$.

Computing the final payment θ_s . For each $\mathbf{y}^r \in \mathbf{Z}_{\mathcal{W},S}$, compute $\theta_i^{d,r} = \theta_i^d(\mathbf{y}^r)$. The final payment made to participant s is given by $\theta_s = \sum_{\mathbf{y}^r \in \mathbf{Z}_{\mathcal{W},S}} P(\mathbf{Y}_{\mathcal{W}} = \mathbf{y}^r | \mathbf{y}_S) \cdot \theta_s^{d,r}$. Note that the set $\mathbf{Z}_{\mathcal{W},S}$ could be exponentially large, and hence computing the exact θ_s may be intractable. However, one can use sampling to get estimates of θ_s in polynomial time (using Hoeffding's inequality to bound sample complexity) and thus implement an approximately truthful payment scheme to any desired accuracy. Further, note that the approximation guarantees of \mathcal{M} do not require computation of the payments at all, and only require execution of the allocation policy, which runs in polynomial time.

Analysis of SEQTGREEDY

We now analyze the mechanism and prove its desirable properties. The proofs of all theorems are presented in the extended version of the paper (Singla and Krause 2013a). We only sketch them here.

Theorem 1. *SEQTGREEDY is truthful in expectation, i.e., no user can increase her profit in expectation by lying about her true cost, for a fixed set of bids of other users.*

Firstly, truthfulness of payments $\theta_s^{d,r}$ is proved for a considered realization \mathbf{y}^r . This is done by showing the monotonicity property of the greedy allocation policy and proving the threshold nature of the payment $\theta_s^{d,r}$. Truthfulness of the actual payment θ_s follows from the fact that it is a linear combination of individually truthful payments $\theta_s^{d,r}$.

Theorem 2. *Payments made by SEQTGREEDY are individually rational, i.e. $\theta_s \geq b_s$.*

This is proved by showing a lower bound of b_s on each of the payments $\theta_s^{d,r}$ used to compute the final payment θ_s .

Theorem 3. *For $\alpha = 2$, SEQTGREEDY is budget feasible, i.e., $\theta_S \leq B$. Moreover, an application specific tighter bound on α can be computed to better utilize the budget.*

We first show that when full budget B is used by mechanism, the maximum raise in bid b'_s that a participant s can make, keeping the bids of other users to be the same, to still get selected by mechanism is upper-bounded by $\alpha \cdot B \cdot \Delta_s / (\sum_{s' \in S} \Delta_{s'})$. By adapting the proof of Chen, Gravin, and Lu (2011), we prove that α is bounded by 2. Surprisingly, this payment bound on α holds irrespectively of the payment scheme used by the mechanism. Hence, when the

budget is reduced by $\alpha = 2$, this results in an upper bound on the payments made to any participant by $B \cdot \Delta_s / (\sum_{s' \in S} \Delta_{s'})$. Summing over these payments ensures budget feasibility. Moreover, by adapting a proof from Singer (2010), we show that a tighter bound on α can be computed based on the characterization of threshold payments used by SEQTGREEDY. Intuitively, the proof is based on the fact that a raise in bid that a participant can make depends on how much utility the application would lose if she refused to participate.

Theorem 4. *For $\alpha = 2$, SEQTGREEDY achieves a utility of at least $\left(\frac{e-1}{3e} - \gamma\right)$ times that obtained by the optimal policy SEQOPT with full knowledge of the true costs. Hereby, γ is the ratio of the participants' largest marginal contribution f_{\max} and the expected utility achieved by SEQOPT.*

We show that, because of the diminishing returns property of the utility function, the stopping criteria used by the mechanism based on proportional share and using only an α proportion of the budget still allows the allocation of sufficiently many participants to achieve a competitive amount of utility. As a concrete example, if each participant can contribute at most 1% to the optimal utility (i.e., $\gamma = 0.01$), Theorem 4 guarantees a constant approximation factor of 0.20.

Experimental Evaluation

In this section, we carry out extensive experiments to understand the practical performance of our mechanism on a realistic community sensing case study.

Benchmarks. We compare against the following benchmarks and state-of-the-art mechanisms.

- SEQGREEDY (unrealistically) assumes access to the true costs of the users, thus measuring the loss incurred by SEQTGREEDY for enforcing truthfulness and serving as upper bound benchmark on untruthful mechanisms.
- RANDOM allocates users randomly until the budget is exhausted and pays each participant its true cost. This represents a lower bound benchmark on untruthful mechanisms.
- CONSTTGREEDY is the non-adaptive variant of SEQTGREEDY and the state-of-the-art truthful mechanism.
- TGREEDY (unrealistically) assumes access to the exact sensing profiles of the users and hence provides insights in measuring the loss incurred due to privacy protection.

Metrics and experiments. The primary metric we measure is the utility acquired by the application. We also measure budget required to achieve a specified utility. To this end, we conduct experiments by varying the given budget and then varying the specified utility, for a fixed obfuscation level. To further understand the impact of random obfuscation, we then vary the level of obfuscation and measure i) % Gain from adaptivity (SEQTGREEDY vs. CONSTTGREEDY), ii) % Loss from truthfulness (SEQTGREEDY vs. SEQGREEDY), and iii) % Loss from privacy (SEQTGREEDY vs. TGREEDY). We present below the results obtained based on data gathered from Mechanical Turk (henceforth MTurk). The primary purpose of using Mechanical Turk (MTurk) data is to evaluate on realistic distributions rather than making assumptions

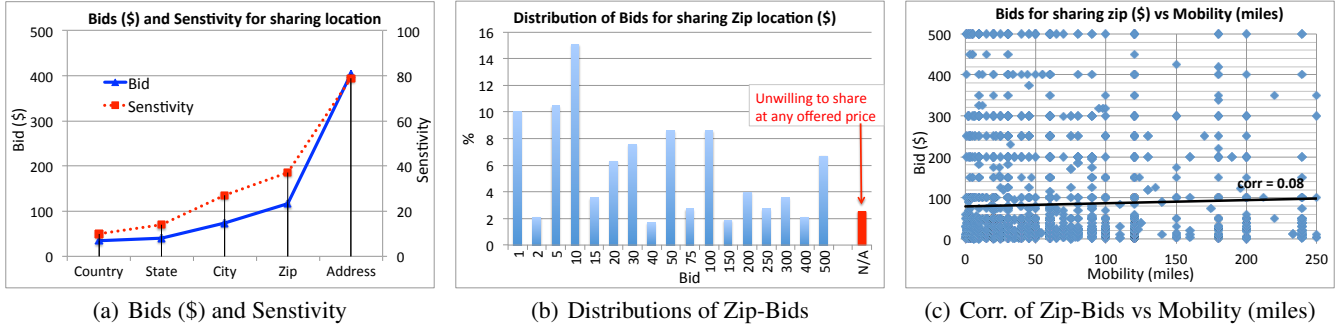


Figure 3: (a) Bids (\$) and sensitivity ([1-100]) for different levels of privacy tradeoff; (b) Distribution of bids (\$) for sharing location at a granularity level of zip codes; (c) Correlation of bids (\$) (for sharing zip) with mobility (daily distance in miles).

about bids and participants' mobility. We carried out experiments on simulated distributions as well with qualitatively similar results.

Experimental setup and data sets

We now describe our setup and data collection from MTurk.

Community sensing application. Suppose we wish to monitor air quality using mobile sensors (Aberer et al. 2010). We consider a granularity level of zip codes and locations \mathcal{V} correspond to the zip codes of state Nevada, USA. We obtained information related to latitude, longitude, city and county of these zips from publicly available data². This represents a total of 220 zip codes located in 98 cities and 17 counties. In order to encourage spatial coverage, we choose our objective f such that one unit utility is obtained for every zip code location observed by the selected participants. To simulate a realistic population of the N users, we also obtained the population statistics for these zip codes³.

MTurk data and user attributes. We posted a Human Intelligence Task (HIT) on MTurk in form of a survey, where workers were told about an option to participate in a community sensing application. Our HIT on MTurk clearly stated the purpose as purely academic, requesting workers to provide correct and honest information. The HIT presented the application scenario and asked workers about their willingness ("yes/no") to participate in such applications. 75% (487 out of 650) responded positively. Workers were asked to express their sensitivity (on scale of [1-100]), as well as the payment bids (in range of [1-500] \$) they desire to receive about exposing their location at the granularity of home address, zip, city, state or country respectively. Additionally, workers were asked about their daily mobility to gather data for defining the sensing radii of the users in our experiments.

A total of 650 workers participated in our HIT, restricted to workers from the USA with more than 90% approval rate and were paid a fixed amount each. We used the data of 487 workers for our experiments, who responded positively to participate in the envisioned application. Figure 4(a) shows the mean bids and expressed sensitivity for different levels of obfuscation. Figure 3(b) shows the distribution of bids for exposing zip level location information. A mean daily mo-

bility of 18 miles was reported. Figure 3(c) shows no correlation between their daily mobility (related to user's sensing radius and hence utility) and bids for exposing zip code information (related to user's bid).

Parameter choices and user profiles. We consider a population of size $N = 500$, distributed according to the population statistics for the zip codes. We used the distribution of bids reported for sharing location at a granularity level of zip codes. We set $c_{\min} = 0.01$ and $c_{\max} = 1$ by scaling the bids in this range. For a given location of a user, we used the distributions of daily mobility to define the sensing radius of the users. We set the maximum possible utility obtained from each user to $f_{\max} = 15$ by limiting the maximal number of observable zip code locations of each user to 15, which are randomly sampled from the locations covered by the user's sensing radius.

Given a user's zip location, the sensing profile of the user is uniquely specified. To create privacy profiles, we used obfuscated user locations, by considering obfuscation at city or state level in which the user is located. We also considered obfuscation within a fixed radius, centered around the user's location. For each of the obfuscated zip codes, multiple corresponding sensing profiles are generated, which collectively define the user's privacy profile.

Results

We now discuss the findings from our experiments.

Computing tighter bounds on payment. Based on Theorem 3, we compute tighter bounds on the payment and optimized the budget reduction factor α used by our mechanism in an application specific manner.

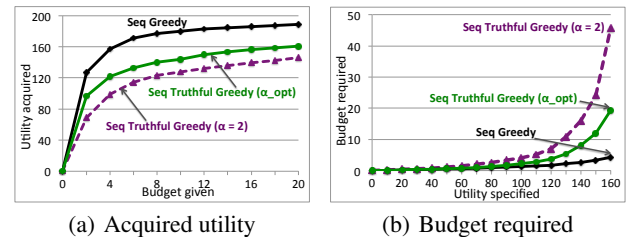


Figure 4: (a) and (d) compares SEQTGREEDY using $\alpha = 2$ w.r.t. to a variant using an optimized value of α .

In community sensing applications with a large number of users and bounded maximal contribution from each user,

²<http://www.populardata.com/downloads.html>

³<http://mcdc2.missouri.edu/>

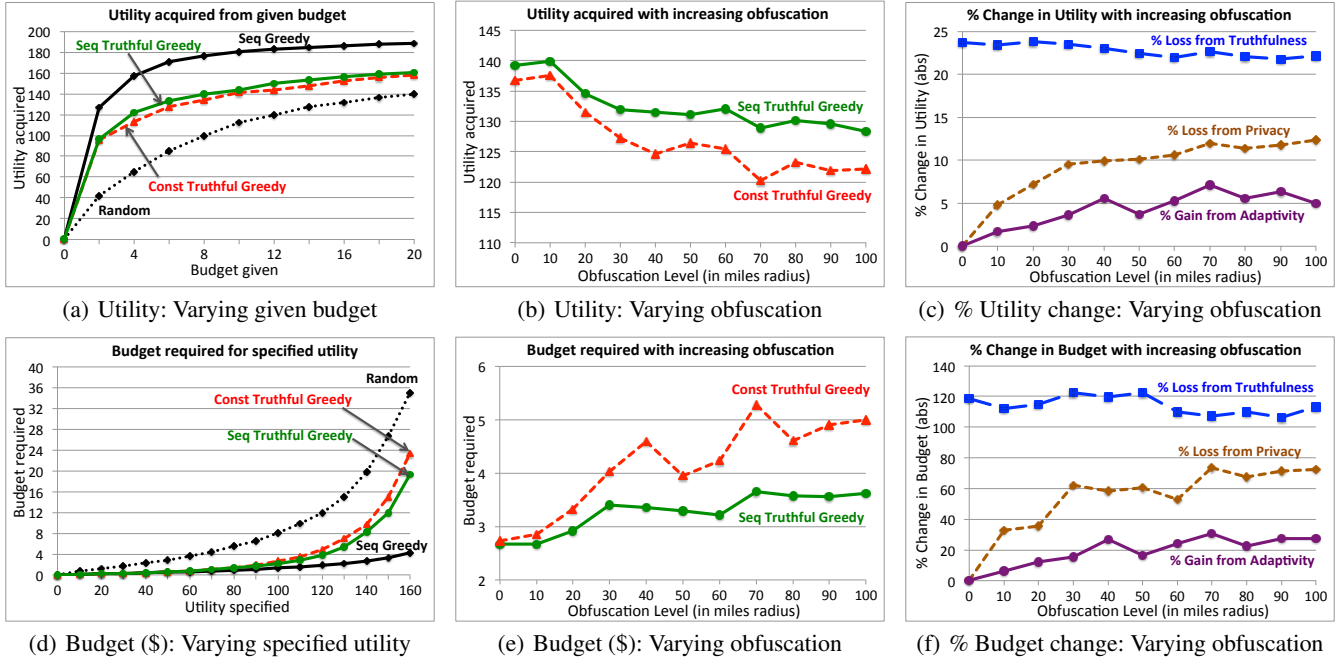


Figure 5: In (a) and (d), for a fixed obfuscation level of 100 miles radius, budget given and desired utility are varied. In (b), (c), (e) and (f) the obfuscation level is varied. (b) and (c) measure utility acquired for a given budget of \$5 and show about 5% adaptivity gain. (e) and (f) measure the budget required (in \$) to achieve a utility of 120 and show up to 30% adaptivity gain.

α is close to 1, resulting in a utilization of almost the entire budget. Figure 4 demonstrates the benefit of using tighter payment bounds (optimized α), compared to a mechanism simply using $\alpha = 2$. Henceforth, in the results, we use the optimized α for all the truthful mechanisms.

Varying the given budget and specified utility. For a fixed obfuscation level of 100 miles radius, Figures 5(a) and 5(d) show the effect of varying the given budget and desired utility respectively. Figure 5(a) illustrates the bounded approximation of our mechanism SEQTGREEDY w.r.t. SEQGREEDY and up to 5% improvement over CONSTTGREEDY in terms of acquired utility. Figure 5(d) shows that the budget required to achieve a specified utility by our mechanism is larger w.r.t. SEQGREEDY and we achieve up to 20% reduction in required budget by using the adaptive mechanism.

Utility acquired at different obfuscation levels. In Figures 5(b) and 5(c), the acquired utility is measured for a given budget of \$5 by varying the obfuscation level. We can see that adaptivity helps acquire about 5% higher utility and this adaptivity gain increases with higher obfuscation (more privacy). The loss from truthfulness is bounded (by 25%), agreeing with our approximation guarantees. The loss from the lack of private information grows, but so also does the gain from adaptivity, which helps to reduce the loss we incur due to privacy protection.

Budget required at different obfuscation levels. In Figures 5(e) and 5(f), the required budget is computed for a desired utility value of 120 by varying the obfuscation level. We can see an increasing adaptivity gain, up to a total of 30% reduction in required budget. As the privacy level increases, the adaptivity gain increases to help partially recover the incurred loss from privacy in terms of budget requirement.

Conclusions and Future Work

There is much potential in intelligent systems that incentivize and empower their users to consciously share certain private information. We presented a principled approach for negotiating access to such private information in community sensing. By using insights from mechanism design and adaptive submodular optimization, we designed the first adaptive, truthful and budget feasible mechanism guaranteed to recruit a near-optimal subset of participants. We demonstrated the feasibility and efficiency of our approach in a realistic case study. Privacy tradeoff is a personal choice and sensitive issue. In realistic deployments of the proposed approach, the choice of participation ultimately lies with the users. We believe that this integrated approach connecting privacy, utility and incentives provides an important step towards developing practical, yet theoretically well-founded techniques for community sensing.

There are some natural extensions for future work. Here, we considered a fairly simple utility function for the sensing phenomena. More complex objectives, e.g., reduction in predictive variance in a statistical model, can be readily incorporated. Further, we would like to design an application (e.g., smartphone app) for deploying our approach in a real world sensing application. It would also be interesting to apply our mechanisms to other application domains that involve uncertainty, sequential decision-making and strategic interactions, e.g., viral marketing.

Acknowledgments. We would like to thank Yuxin Chen and Gábor Bartók for helpful discussions. This research was supported in part by SNSF grant 200021_137971, ERC StG 307036 and a Microsoft Research Faculty Fellowship.

References

- [Aberer et al. 2010] Aberer, K.; Sathé, S.; Chakraborty, D.; Martinoli, A.; Barrenetxea, G.; Faltings, B.; and Thiele, L. 2010. Opensense: Open community driven sensing of environment. *IWGS*.
- [Asadpour, Nazerzadeh, and Saberi 2008] Asadpour, A.; Nazerzadeh, H.; and Saberi, A. 2008. Maximizing stochastic monotone submodular functions.
- [Carrascal et al. 2013] Carrascal, J. P.; Riederer, C.; Erramilli, V.; Cherubini, M.; and de Oliveira, R. 2013. Your browsing behavior for a big mac: economics of personal information online. *WWW '13*, 189–200.
- [Chen, Gravin, and Lu 2011] Chen, N.; Gravin, N.; and Lu, P. 2011. On the approximability of budget feasible mechanisms. In *SODA*.
- [Chon et al. 2012] Chon, Y.; Lane, N. D.; Li, F.; Cha, H.; and Zhao, F. 2012. Automatically characterizing places with opportunistic crowdsensing using smartphones. In *UbiComp*.
- [Chorppath and Alpcan 2012] Chorppath, A. K., and Alpcan, T. 2012. Trading privacy with incentives in mobile commerce: A game theoretic approach. *Pervasive and Mobile Computing*.
- [Clayton et al. 2012] Clayton, R.; Heaton, T.; Chandy, M.; Krause, A.; Kohler, M.; Bunn, J.; Olson, M.; Faulkner, M.; Cheng, M.; Strand, L.; Chandy, R.; Obenshain, D.; Liu, A.; Aivazis, M.; and Guy, R. 2012. Community seismic network. *Annals of Geophysics* 54(6):738–747.
- [Dwork 2006] Dwork, C. 2006. Differential privacy. In *ICALP*, volume 4052, 1–12.
- [Faltings, Jurca, and Li 2012] Faltings, B.; Jurca, R.; and Li, J. J. 2012. Eliciting truthful measurements from a community of sensors. *3rd Int. Conference on Internet of Things* 51–18.
- [Feige 1998] Feige, U. 1998. A threshold of $\ln n$ for approximating set cover. *Journal of the ACM* 45:314–318.
- [Golovin and Krause 2011] Golovin, D., and Krause, A. 2011. Adaptive submodularity: Theory and applications in active learning and stochastic optimization. *Journal of Artificial Intelligence Research (JAIR)* 42:427–486.
- [Himmel et al. 2005] Himmel, M.; Rodriguez, H.; Smith, N.; and Spinac, C. 2005. Method and system for schedule based advertising on a mobile phone.
- [Horvitz et al. 2005] Horvitz, E.; Apacible, J.; Sarin, R.; and Liao, L. 2005. Prediction, expectation, and surprise: Methods, designs, and study of a deployed traffic forecasting service. In *UAI*.
- [Hui et al. 2011] Hui, P.; Henderson, T.; Brown, I.; and Haddadi, H. 2011. Targeted advertising on the handset : privacy and security challenges. In *Pervasive Advertising, HCI'11*.
- [Kansal et al. 2007] Kansal, A.; Nath, S.; Liu, J.; and Zhao, F. 2007. Senseweb: An infrastructure for shared sensing. *IEEE Multimedia* 14(4).
- [Krause and Guestrin 2007] Krause, A., and Guestrin, C. 2007. Near-optimal observation selection using submodular functions. In *AAAI, Nectar track*.
- [Krause and Horvitz 2008] Krause, A., and Horvitz, E. 2008. A utility-theoretic approach to privacy and personalization. In *AAAI*.
- [Krause et al. 2008] Krause, A.; Horvitz, E.; Kansal, A.; and Zhao, F. 2008. Toward community sensing. In *IPSN*.
- [Krumm and Horvitz 2006] Krumm, J., and Horvitz, E. 2006. Predestination: Inferring destinations from partial trajectories. In *UbiComp*, 243–260.
- [Krumm 2007] Krumm, J. 2007. Inference attacks on location tracks. In *PERVASIVE*, 127–143.
- [Li and Faltings 2012] Li, J. J., and Faltings, B. 2012. Incentive schemes for community sensing. *The 3rd International Conference in Computational Sustainability*.
- [Lieb 2007] Lieb, D. A. 2007. MoDOT tracking cell phone signals to monitor traffic speed, congestion.
- [Liu, Krishnamachari, and Annavaram 2008] Liu, H.; Krishnamachari, B.; and Annavaram, M. 2008. Game theoretic approach to location sharing with privacy in a community based mobile safety application. In *MSWiM*, 229–238.
- [Machanavajjhala et al. 2006] Machanavajjhala, A.; Kifer, D.; Gehrke, J.; and Venkitasubramaniam, M. 2006. L-diversity: Privacy beyond k-anonymity. In *ICDE*.
- [Mobile-Millennium 2008] Mobile-Millennium. 2008. Mobile millennium traffic-monitoring system. <http://traffic.berkeley.edu/>.
- [Myerson 1981] Myerson, R. 1981. Optimal auction design. *Mathematics of Operations Research* 6(1).
- [Nemhauser, Wolsey, and Fisher 1978] Nemhauser, G.; Wolsey, L.; and Fisher, M. 1978. An analysis of the approximations for maximizing submodular set functions. *Math. Prog.* 14:265–294.
- [Olson, Grudin, and Horvitz 2005] Olson, J.; Grudin, J.; and Horvitz, E. 2005. A study of preferences for sharing and privacy. In *CHI*.
- [Singer 2010] Singer, Y. 2010. Budget feasible mechanisms. In *FOCS*, 765–774.
- [Singer 2012] Singer, Y. 2012. How to win friends and influence people, truthfully: Influence maximization mechanisms for social networks. In *WSDM*.
- [Singla and Krause 2013a] Singla, A., and Krause, A. 2013a. Incentives for privacy tradeoff in community sensing (extended version). <http://arxiv.org/abs/1308.4013>.
- [Singla and Krause 2013b] Singla, A., and Krause, A. 2013b. Truthful incentives in crowdsourcing tasks using regret minimization mechanisms. *WWW '13*, 1167–1178.
- [Sviridenko 2004] Sviridenko, M. 2004. A note on maximizing a submodular set function subject to knapsack constraint. *Operations Research Letters* v.(32):41–43.
- [Sweeney 2002] Sweeney, L. 2002. k-anonymity: a model for protecting privacy. *Int. Journal on Uncertainty, Fuzziness and Knowledge-based Systems* 10(5):557–570.
- [Wunnavva et al. 2007] Wunnavva, S.; Yen, K.; Babij, T.; Zavaleta, R.; Romero, R.; and Archilla, C. 2007. Travel time estimation using cell phones (TTECP) for highways and roadways. Technical report, Florida Department of Transportation.
- [Yoon, Noble, and Liu 2007] Yoon, J.; Noble, B.; and Liu, M. 2007. Surface street traffic estimation. In *MobiSys*, 220–232.
- [Zheng, Xie, and Ma 2010] Zheng, Y.; Xie, X.; and Ma, W.-Y. 2010. Geolife: A collaborative social networking service among user, location and trajectory. *IEEE Data Engineering Bulletin* 32–40.