

# Privacy Preservation for Context Sensing on Smartphone

Wei Wang, *Member, IEEE*, and Qian Zhang, *Fellow, IEEE*

**Abstract**—The proliferation of sensor-equipped smartphones has enabled an increasing number of context-aware applications that provide personalized services based on users’ contexts. However, most of these applications aggressively collect users’ sensing data without providing clear statements on the usage and disclosure strategies of such sensitive information, which raises severe privacy concerns and leads to some initial investigation on privacy preservation mechanisms design. While most prior studies have assumed static adversary models, we investigate the context dynamics and call attention to the existence of intelligent adversaries. In this paper, we identify the context privacy problem with consideration of the context dynamics and malicious adversaries with capabilities of adjusting their attacking strategies. Then, we formulate the interactive competition between users and adversaries as a competitive Markov decision process (MDP), in which the users attempt to preserve the context-based service quality and their context privacy in the long-term defense against the strategic adversaries with the opposite interests. In addition, we propose an efficient minimax learning algorithm to obtain the optimal policy of the users and prove that the algorithm quickly converges to the unique Nash equilibrium point. Our evaluations on real smartphone context traces of 94 users demonstrate that the proposed algorithm largely improves the convergence speed by three orders of magnitude compared with traditional algorithm and the optimal policy obtained by our minimax learning algorithm outperforms the baseline algorithms.

**Index Terms**—Context sensing, Privacy preservation, Markov decision process (MDP)

## I. INTRODUCTION

The increasing popularity of sensor-equipped smartphones provides new opportunities for proliferation of context-aware applications that offer personalized services based on the operating conditions of smartphone users and their surrounding environments. Such applications effectively use sensors such as GPS, accelerometer, proximity sensor and microphone to infer smartphone user’s current context including location, mobility mode (e.g., walking or driving), and social activities. Examples of context-aware applications include *GeoNote* [1] that reminds a user of something when he is at a particular location, *Running* [2] that keeps track of user’s jogging trajectory, and *AutoSilent* [3] that automatically mutes the phone when the user is in a meeting.

Although context-aware applications improve user experiences on smartphones, severe privacy issues arise with these

Wei Wang is with the School of Electronic Information and Communications, Huazhong University of Science and Technology, and the Fok Ying Tung Research Institute, Hong Kong University of Science and Technology. E-mail: gswwang@cse.ust.hk.

Qian Zhang is with the Department of Computer Science and Engineering, Hong Kong University of Science and Technology. E-mail: qianzh@cse.ust.hk.

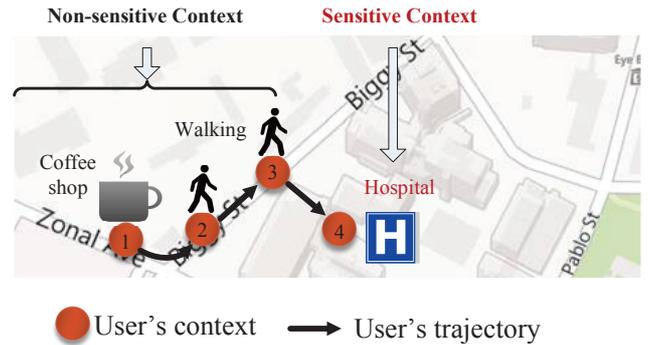


Fig. 1. An illustration on context privacy.

applications. Nowadays, the growing privacy threats of sharing location-related context information via context-aware applications on smartphones have been concerned by both consumers [4] and governments [5]. Such privacy threats come from the fact that many smartphone applications aggressively collect sensing data without clear statements about how to use the sensing data and whom the sensing data will be shared with. Untrusted applications may sell such personal information to advertisers without user’s permission. Enck et al. [6] studied 30 popular Android applications that have access to user’s location, camera, microphone data, and found that 15 of them sent users’ information to remote advertisement or analytics servers. Moreover, malicious adversaries with criminal intent could hack the applications with such information to pose a threat to individual security and privacy. Being aware of such risks, the smartphone users may not allow the applications to access their sensing data, which, however, disables the functionalities provided by the context-aware applications, and thus, causes inconvenience to the users.

To enable smartphone users to enjoy services provided by smartphone applications with privacy protection, many existing privacy preserving approaches have been proposed to explore better tradeoffs between service quality and individual privacy. Most of these approaches focus on location privacy [7]–[10], which, however, fall short when applied to context privacy analysis due to the dynamics of user behaviors and temporal correlations between contexts. Specifically, smartphone users usually transit between different contexts (e.g., a user goes to a particular hospital after eating at a coffee shop), whose sensitivities are different to the users. Moreover, the contexts are usually correlated, which has already been studied for different goals [11]–[13]. Thus, the adversaries

can learn the connections between contexts by exploiting the temporal correlations, and then use such correlations to infer user’s sensitive contexts based on their observations on non-sensitive contexts. For example, in Fig. 1, a context-aware application may learn that a user regularly follows a trajectory  $1 \rightarrow 2 \rightarrow 3 \rightarrow 4$ . Then, releasing the context information that the user is at the coffee shop at location 1 may reveal that the user is very likely to go to the hospital, which is sensitive to the user. However, the frameworks on location privacy [7]–[10] do not consider such inference attacks from adversaries knowing temporal correlations, and thus, are not directly applicable to context privacy analysis.

To the best of our knowledge, the only existing work on context privacy protection is *MaskIt* [12], which assumes that adversaries take fixed attacking strategies that do not change over time. However, some adversaries launch continuous online attacks [10] or long-term offline monitoring, in which case the adversaries may adapt their attacking strategies over time to gain more benefits. For example, a context-aware application may sell user’s sensing data to remote advertisement adversaries, who continuously push context-related ads or spam to users based on the user’s instant context information. Note that in the real world, context-based ads or spam need to be delivered in real time (e.g., NAVTEQ or AdLocal by Cirius Technologies) as users may lose interest if the ads do not match current context. In such case, it is highly possible that the adversaries will adapt their attacking strategies based on their observations of previous attacking results and context dynamics.

To satisfy the aforementioned requirements, in this paper, we model the strategic and dynamic competition between a smartphone user and a malicious adversary as a competitive Markov decision process (MDP), where the user preserves context-based service quality and context privacy against strategic adversaries. Both offline attack and online attack are considered. The user’s action is to control the released data granularity of each sensor used by context-aware applications in a long-term defense against the adversary, while the adversary’s action is to select which sensing data as the source for attacks. The interactive competition between the user and the adversary are considered to last for a number of stages with the contexts dynamics. Both the user and the adversary observe previous contexts and their transitions, based on which both players adjust their future strategies. An efficient minimax learning algorithm with proved convergence is proposed to obtain the user’s optimal defense strategy. Compared to traditional learning algorithm, the proposed algorithm significantly reduces the computational cost by reducing the dimensions of state values that need to be updated. We give both analytical results and evaluations on real smartphone traces to analyze the factors that affect the user’s optimal defense strategy.

The main contributions of this paper are summarized as below.

- We identify the context privacy problem in context-aware applications with adaptable adversaries. We consider dynamics of user’s context and powerful adversaries that know the temporal correlations between contexts and are capable of adjusting their attacking strategies over time.

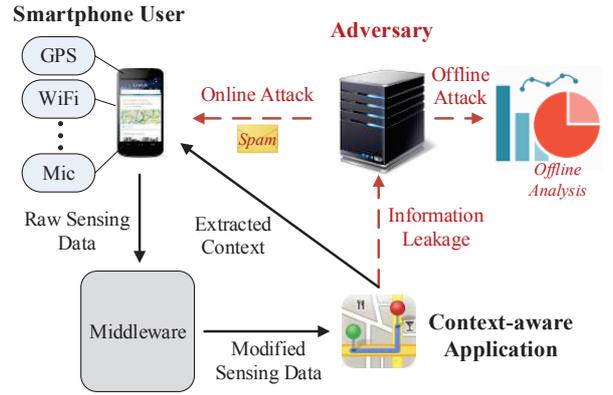


Fig. 2. A mobile phone context sensing system.

These distinct features make existing works inapplicable to this problem.

- We analyze the context privacy problem via a competitive MDP formulation. The interactive competition between smartphone users and adversaries are modeled as a competitive MDP, which captures the distinct features of the context privacy problem. In particular, we discuss the cases of both online and offline attacks.
- We devise an efficient minimax learning algorithm with provable convergence to obtain optimal policies. We improve the efficiency of the learning algorithm by solving an equivalent problem with reduced dimensions. The convergence speed is reduced by three orders of magnitude compared with the traditional learning algorithm. We also prove that the algorithm converges to the Nash equilibrium (NE) point.
- We use real smartphone context traces of 94 users to demonstrate the efficacy and efficiency of the proposed algorithm. The results show that the optimal policy obtained by our minimax learning algorithm outperforms the baseline algorithms. Promisingly, the results give guidance to the design of context privacy preserving mechanisms.

The rest of the paper is organized as follows. Section II introduces the system model. Section III presents the competitive MDP formulation for the context privacy problem. Section IV proposes a minimax learning algorithm to obtain the user’s optimal defense strategy under online attacks, and Section IV-G discusses the optimal policy under offline attacks. Section V describes the performance evaluations, and Section VI reviews the related works. Finally, Section VII concludes the paper.

## II. SYSTEM MODEL

In this section, we describe the model of the mobile phone context sensing system and the privacy issue when the context-aware application is untrusted.

### A. Context Sensing System

**User context.** A smartphone user encounters a set of contexts  $\mathcal{C} = \{c_1, \dots, c_n\}$ , including locations, mobility modes,

and social activities. Due to user's behaviors and activities, the user's contexts keep changing. We assume that all user activities can be classified into a finite number of elementary activities such as user's motion states or locations, and each user can only perform one activity in one time slot [14], [15]. The transitions between contexts can be captured by the Markov model: previous studies [14], [15] have shown that human behaviors and activities extracted from smartphone sensors can be modeled well with a two-state Markov chain. Specifically, at time  $t$ , the user's context is denoted as  $C_t \in \mathcal{C}$ , which is generated from a Markov model  $M$ . According to the independence property of Markov chains, we have  $\Pr[C^t = c_i | C^1, \dots, C^{t-1}] = \Pr[C^t = c_i | C^{t-1}]$ .

**Context sensing.** Fig. 2 illustrates a smartphone context sensing system, where a sensor-equipped smartphone runs untrusted context-aware applications. The smartphone user senses its environment with multiple sensors (e.g., GPS, Wi-Fi, microphone) and releases the sensing data to the application periodically for energy efficiency reasons [11], [12], where a period is referred to as a time slot in the context sensing system. It is worthwhile noting that there can be multiple sensor samples in one period, and the sampling rates can be different for different sensors. For example, we need multiple accelerometer samples to determine a user's motion state in one period while we only need a single GPS sample to determine the user's location. The context-aware applications provide services based on the user's contexts, which are extracted by the applications using certain context recognition approaches [11], [13].

### B. Privacy Issue in Context Sensing

A subset of contexts  $\mathcal{C}$  are considered to be private contexts whose disclosure is undesired by the smartphone user. The user claims a subset of  $\mathcal{C}$  to be sensitive via special applications (e.g., Locaccino [16]). The user's context privacy is breached if the adversary successfully infers that the user is in its sensitive context. To protect context privacy, the user can control the released data granularity of each sensor via the privacy-preserving middleware (e.g., *MaskIt* [12]). The privacy-preserving middleware employs a certain existing privacy-preserving technique (e.g., [12], [17], [18]) to modify the raw sensing data before releasing them to the applications. Context-aware applications can only access user's data via the privacy-preserving middleware while they do not have the permission to access raw sensing data. These applications provide services based on the user's contexts, which are extracted by the applications using certain context recognition approaches [11], [13]. Normally, the released sensing data with coarser granularity leaks less information about the user, while the accuracies of context recognition performed by the context-aware applications are also compromised.

The context-aware application is untrusted or corrupted, and is considered to be the adversary. The application collects users' data in an authorized manner, but tries to extract users' private information. As such, the adversary is able to obtain the released sensing data at the time when the untrusted application accesses the data. The adversary is assumed to know

the Markov chain of a user [12]. The sensing data retrieved by the adversary in a time slot is limited due to computational constraints (caused by curse of dimensionality when using private data [8]) or limited bandwidth used for retrieving data. As the contexts and user's released data granularity vary over time, the adversary can adaptively choose different subsets of sensors to maximize its long-term utility. To protect smartphone users against all kinds of adversaries, we make the worst case assumption: the adversary is a *malicious* attacker that aims at minimizing user's utility through a series of strategic attacks [10]. We consider the following two types of attacks:

- **Offline attack.** The adversary passively collects the user's sensing data, and infer the user's personal information, such as context behaviors and trajectories, via offline analysis. The adversary can sell the information, which is considered to be of great value [19]. The user is unaware of the attack results, i.e., to what extent its privacy is breached.
- **Online attack.** The adversary collects the user's sensing data and actively infers the user's instant context information based on collected data. Based on the instant context, the adversary may push context-based spams/scams to the user, or even make the user a victim of blackmail or physical violence.

It is worth noting that the common part of both attacks is that adversaries aggressively collect the user's sensing data, while the only departure is whether the adversaries actively interact with the user in real-time. The type of attacks can be easily discerned at the user end by checking whether there exist observable activities from the adversary.

### C. Problem Statement

Our goal is to find the optimal defense strategies for users to preserve privacy against the malicious adversary over a serial of correlated contexts. In this paper, we discuss the optimal strategies under both offline and online attacks. Since the context is considered to keep changing over time and both the user and the adversary make different actions at different times, the interactions between the user and the adversary is in a stochastic setting and should be formulated as a competitive MDP.

## III. PROBLEM FORMULATION

A competitive MDP, or stochastic game, is a dynamic game with probabilistic transitions played in a sequence of stages. A two player competitive MDP  $\Gamma$  consists of a six-tuple  $\langle \mathcal{S}, \mathcal{A}^1, \mathcal{A}^2, r^1, r^2, P \rangle$ .  $\mathcal{S}$  is the discrete state space.  $\mathcal{A}^k$  is the action space of player  $k$  for  $k = 1, 2$ .  $r^k : \mathcal{S} \times \mathcal{A}^1 \times \mathcal{A}^2 \mapsto \mathbb{R}$  is the stage payoff function for player  $k$ , where  $\mapsto$  denotes the input-output mapping of a function. Note that action spaces of different players in a stochastic game can be different [20].  $P : \mathcal{S} \times \mathcal{A}^1 \times \mathcal{A}^2 \mapsto \Delta(\mathcal{S})$  is the transition probability map, where  $\Delta(\mathcal{S})$  is the set of probability distributions over  $\mathcal{S}$ . Note that both the stage payoff function  $r^k$  and the transition probability map  $P$  take the set of state  $\mathcal{S}$  and actions  $\mathcal{A}^1, \mathcal{A}^2$  as the domain of definition, that is, the stage payoff and state transition rely

on the system state and players' actions. The game  $\Gamma$  is played in a sequence of stages, where each player  $k$  receives a stage payoff  $r^k(s, a^1, a^2)$  based on players actions  $a^k \in \mathcal{A}^k$  and current stage  $s \in \mathcal{S}$ . Each player  $k$  attempts to maximize its expected sum of discounted payoffs. In our formation, the players' actions are not observable to each other.

In this section, we formulate the privacy problem in mobile phone context sensing as a competitive MDP.

#### A. States and Actions

1) *System States*: In each time slot, the smartphone user is in a certain context and releases data of multiple sensors to the context-aware application. The user's context is included in the system state as the user's action depends on its observation of the current context. Note that the current context is only observable to the user, while the adversary can only infer the context based on the modified sensing data and the user's Markov model.

**System states in online attacks.** In the case of online attacks, previous attack results should also be included in the system state. As the adversary's strategy is not known by the user, the user can only conjecture the adversary's strategy from previous attack results, which are assumed to be observable to the user as online attacks have instant impact on the user. The reason why next action depends on the previous attack result is that i) the previous result determines whether the adversary has correctly inferred the last context, and ii) the current context is correlated to the last context according to the Markov chain. For instance, if a user receives an advertisement based on its current private context, then the user knows that the adversary successfully inferred this private context; if the user receives an advertisement based on a context that it has never been to, then the user knows that the adversary has failed to infer its true context. Thus, the user should maintain a record of which contexts the adversary has launched attacks on, and which contexts have been successfully attacked. The attack result observed at time  $t$ , namely the attack result in the last time slot, is denoted as  $Ar^t$ , whose value can be  $Ar^t = 1$  meaning the adversary successfully infers the context  $C_{t-1}$ , or  $Ar^t = 0$  meaning the adversary fails to infer the context  $C_{t-1}$ . In summary, the context and attack result are observable to the user and affect the user's decisions. Thus, the state of the competitive MDP at time  $t$  is defined by  $S_{\text{on}}^t = \{C^t, Ar^t\}$ .

**System states in offline attacks.** The adversary who launches offline attacks passively collects sensing data for offline analysis, and thus the user cannot observe attack results. Analogous to online attacks, the adversary's strategy of offline attacks is not known by the user. Hence, the user can only observe the current context and the system state in the case of offline attacks is  $S_{\text{off}}^t = C^t$ .

2) *User's Actions*: After observing the state  $S^t$  at each stage (note that the adversary can only infer  $C^t$  based on  $M$ ), the user decides its action for the current stage. As discussed in Section II, the user controls the granularity of the released sensing data to protect its context privacy while preserving the quality of context-based services. For simplicity, we use the accuracy of context recognition to measure the granularity

of the sensing data, which is assumed to be the weighted summation of the data granularity of each sensor. The rationale behind this assumption comes from the generalization technique, which is widely adopted in location-based services and data anonymization [7], [21], [22]. The intuition behind the generalization technique is that when the original data is generalized to be coarser, it incurs more information loss, while providing stronger privacy preservation. Formally, the action of the user at time  $t$  is defined as  $\mathbf{a}_u^t = \{a_{u,1}^t, \dots, a_{u,K}^t\}$ , with each sensor's data granularity  $a_{u,k}^t \in [0, 1], \forall k = 1, \dots, K$ , where  $K$  is the total number of sensors used for recognition. The accuracy of context recognition  $g$  ( $0 \leq g \leq 1$ ) based on  $\mathbf{a}_u^t$  is given by

$$g = \sum_{k=1}^K \kappa_k a_{u,k}^t, \quad (1)$$

where  $\{\kappa_k : \forall k\}$  are the weights measuring the sensitivity of the sensor's data granularity to the context recognition accuracy.

3) *Adversary's Actions*: On the other hand, due to the limited attacking capability, the adversary needs to select a proper subset of sensing data for retrieval. Mathematically, the adversary's actions at time  $t$  are defined as  $\mathbf{a}_a^t = \{a_{a,1}^t, \dots, a_{a,K}^t\}$ , where  $a_{a,k}^t$  is the probability of retrieving the data of the  $k$ th sensor. The power limitation constraints for the adversary's actions are as follows.

$$\begin{aligned} \sum_k a_{a,i}^t &\leq L, \\ 0 &\leq a_{a,i}^t \leq 1, \forall k, \end{aligned} \quad (2)$$

where  $L$  is the power limitation of the adversary. When  $L \geq K$ , the adversary is able to collect all sensing data. This type of adversary is referred to as *the adversary with unlimited power*, which is discussed in Section IV-E.

4) *State Transitions*: **State transitions in online attack.** The state  $S^t$  is uncertain (due to the uncertainty of  $C^t$ ) and depend on the actions of the user and the adversary ( $Ar^t$  depends on the player's actions). We assume that user behavior is independent of player's actions. Then, the state transition probability can be computed by

$$\begin{aligned} \Pr[S_{\text{on}}^{t+1} | S^t, \mathbf{a}_u^t, \mathbf{a}_a^t] &= \Pr[Ar^{t+1} | Ar^t, \mathbf{a}_u^t, \mathbf{a}_a^t] \Pr[C^{t+1} | C^t] \\ &= \Pr[Ar^{t+1} | \mathbf{a}_u^t, \mathbf{a}_a^t] \Pr[C^{t+1} | C^t]. \end{aligned} \quad (3)$$

The second equality holds because  $Ar^{t+1}$  is the attack results observed at time  $t + 1$ , which only depends on the actions players made at the last stage.

**State transitions in offline attack.** As the adversary only passively collects sensing data – which has no instant impact on the user – the state transition probability only depends on the context correlations:

$$\Pr[S_{\text{off}}^{t+1} | S^t, \mathbf{a}_u^t, \mathbf{a}_a^t] = \Pr[C^{t+1} | C^t]. \quad (4)$$

#### B. Stage Payoff

After defining the states and actions, we give a concrete expression of stage payoffs. The payoff function of the user

is defined to be the quality of the context-based service with weighted penalty on privacy loss, which is written as

$$r_u(S^t, \mathbf{a}_u^t, \mathbf{a}_a^t) = QoS(\mathbf{a}_u^t) - \omega \cdot Pri(S^t), \quad (5)$$

where  $QoS(\mathbf{a}_u^t)$  is the quality of context-based service the user enjoys,  $\omega$  the equivalent service quality improvement caused by unit privacy loss, and  $Pri(S^t)$  the privacy loss.  $QoS(\mathbf{a}_u^t)$  is a measure of the user's degree of satisfaction with the context-based service and can be modeled as a sigmoid function of the context recognition accuracy. Sigmoid function has been widely used to approximate the user's satisfaction with respect to service qualities [23]. The rationale behind the sigmoid function is that the user's satisfaction remains low when the service quality grows in a low range; the user's satisfaction grows quickly when the accuracy further grows across a satisfaction threshold; when the accuracy enters a relatively high range, further improvement becomes marginal and brings little benefit to the user. Concretely,  $QoS(\mathbf{a}_u^t)$  is measured as

$$QoS(\mathbf{a}_u^t) = \frac{1}{1 + e^{-\theta(g-\eta)}}, \quad (6)$$

where  $\theta$  decides the steepness of the quality of service satisfactory curve,  $g$  the accuracy of context recognition, and  $\eta$  the satisfaction threshold below which the user has very limited satisfaction (the function curve is convex) and above which the user's satisfaction rapidly approaches an asymptotic value (the function curve is concave).

Next, we measure the privacy loss of the user based on the definition of context privacy in [12]. Consider a user over a day with a context space  $\mathcal{C}$  and a set of sensitive context  $\mathcal{C}_s \subseteq \mathcal{C}$ . We say that the released data preserves privacy if the adversary learns little information about the user being in a private state from the released data, meaning that for all sensitive contexts and all times the difference between the posterior and prior beliefs on the user being in a sensitive context at that time is limited. Normally, the adversary values the information of user's recent contexts more highly than the information about user's contexts in the faraway future. Based on the above intuition, we define *context sensitivity* as follows.

**Definition 1** (Context Sensitivity). *The sensitivity of a context  $c$  is defined to be the sum of the discounted differences between the prior belief and the posterior belief after observing current context on the user being in each sensitive context in the future, that is,*

$$Sens(c) = \sum_{t=0}^{\infty} \sum_{c_s \in \mathcal{C}_s} \gamma^t |\Pr[C^t = c_s | C^0 = c] - \Pr[C^t = c_s]|, \quad (8)$$

where  $0 < \gamma < 1$  is the discount factor of the context privacy.

The sensitivity of a context  $c$  measures the maximum information that the adversary can learn about the user's sensitive contexts in the future by observing the user being in  $c$ .

Based on the context sensitivity, we define the user's privacy loss. In the case of online attacks, if an adversary successfully infers a user's current context, the user's privacy loss is the

sensitivity of the current context. Otherwise, the privacy loss is zero, as the user's true context is still unknown to the adversary. Thus, the privacy loss in the case of online attacks is expressed as

$$Pri_{\text{on}}(S^t, \mathbf{a}_u^t, \mathbf{a}_a^t) = Sens(C^t) Ar^{t+1}, \quad (9)$$

where  $Ar^{t+1}$  is the attack result known at time  $t+1$ , i.e., the attack result for context  $C^t$ . The probability of a successful attack at time  $t$  is  $\Pr[Ar^{t+1}] = \sum_i \kappa_i a_{u,i}^t a_{a,i}^t$ .

When the adversary launches offline attacks, the attack results are unobservable to the user, and thus we cannot directly measure the privacy loss in each stage. To overcome this predicament, we consider the worst case: the adversary always selects the best strategy that cause the most privacy loss. Therefore, the privacy loss under offline attacks is

$$Pri_{\text{off}}(S^t, \mathbf{a}_u^t) = Sens(C^t) \max_{\mathbf{a}_a^t} \left\{ \sum_i \kappa_i a_{u,i}^t a_{a,i}^t \right\}. \quad (10)$$

Then, we decide  $\omega$ , i.e., the equivalent service quality improvement caused by unit privacy loss. For each context, we measure service quality improvement and privacy loss when the adversary can access all user's raw sensing data, compared with the case that the adversary knows nothing. We assume that the adversary has prior belief of a user's context based on its background knowledge (e.g., the adversary knows the user's behavior pattern or the Markov chain of the user's contexts). Therefore, we express  $\omega$  as (7), where  $\Pr[C^t = c]$  is the adversary's prior belief on user's context. Substituting (6) (8) (9) (7) back into (5), we can obtain the stage payoff for the user, while the stage payoff for the adversary is the negative of (5).

Generally, context applications run continuously on a smartphone all day long [11], [12]. Thus, we assume that there is an infinite number of time slots, i.e., the context privacy game is played for an infinite number of stages. Normally, the smartphone users care more about the current context or near future contexts than the faraway future contexts. For example, a user's current context is more private since the adversaries can cause immediate damage to the user. Therefore, the user's utility is to the expected sum of discounted stage payoffs, where the delayed payoffs value less to the user. The discounted stage payoff is defined to be the stage payoff weighted by a discount factor. Hence, the user's utility can be expressed as

$$U_u = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_u(S^t, \mathbf{a}_u^t, \mathbf{a}_a^t) \right], \quad (11)$$

where  $\gamma$  is the discount factor of the context privacy. Then, the user's objective is to derive an *optimal defense strategy* that maximizes  $U_u$ , which is discussed in the following sections. The adversary aims at minimizing user's utility through a series of strategic attacks, and thus its utility is the reverse of the user's utility.

#### IV. LEARNING THE OPTIMAL DEFENSE STRATEGY

Based on the problem formulation in Section III, the context privacy problem under online attacks is a two-player zero-sum stochastic game. In this section, we will first discuss the

$$\omega = \frac{\sum_{c:c \in \mathcal{C}} \Pr[C^t = c](QoS(a_u = \mathbf{1}) - QoS(\mathbf{a}_u^t = \mathbf{0}))}{\sum_{c:c \in \mathcal{C}} \Pr[C^t = c](\Pr(\mathbf{a}_u^t = \mathbf{1}, \mathbf{a}_a^t = \mathbf{1}, C^t = c) - \Pr(\mathbf{a}_u^t = \mathbf{0}, \mathbf{a}_a^t = \mathbf{1}, C^t = c))}, \quad (7)$$

algorithm to derive the NE of the stochastic game under online attacks, so as to obtain the optimal policy of the user. Then, we extend the results to offline attacks.

#### A. Minimax Equilibrium in the Context Privacy Game

Formally, a policy in a stochastic game is defined to be a probability distribution over the action set at any state. A policy  $\pi$  is said to be *stationary* if  $\pi^t = \pi$  for all  $t$ , that is, the policy is fixed over time. In this paper, we are interested in stationary policies. In the context privacy stochastic game, the user's policy is denoted by  $\pi_u : \mathcal{S} \mapsto \Delta(\mathcal{A}_u)$  and the adversary's policy is denoted by  $\pi_a : \mathcal{S} \mapsto \Delta(\mathcal{A}_a)$ , where  $\mathcal{S}$  is the state space for  $S^t$ ,  $\Delta(\mathcal{A}_u)$  and  $\Delta(\mathcal{A}_a)$  the probability distributions over the user's action space  $\mathcal{A}_u$  and the adversary's action space  $\mathcal{A}_a$ , respectively.

In stochastic games, utilities are expressed in the form of *state value*. Here, the initial state is defined to be the state at time  $t = 0$ , denoted by  $S^0$ . Given policies  $\pi_u, \pi_a$  and a state  $s \in \mathcal{S}$ , the user's utility can be written as

$$V^\pi(s) = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[r_u(S^t, \mathbf{a}_u^t, \mathbf{a}_a^t) | \pi_u, \pi_a, S^0 = s]. \quad (12)$$

Denote the actions  $\mathbf{a}_u^t, \mathbf{a}_a^t$  determined by policies  $\pi_u, \pi_a$  to be  $\mathbf{a}_u^\pi, \mathbf{a}_a^\pi$ , respectively. Then, we can rewrite (12) as

$$V^\pi(s) = r_u(s, \mathbf{a}_u^\pi, \mathbf{a}_a^\pi) + \gamma \sum_{s'} \Pr[s' | s, \mathbf{a}_u^\pi, \mathbf{a}_a^\pi] V^\pi(s'). \quad (13)$$

Both user and adversary follow their optimal policies  $\{\pi_u^*, \pi_a^*\}$  that maximize their own utilities, where the optimal policies are called an *optimal policy pair*  $\pi^* = \{\pi_u^*, \pi_a^*\}$ . An optimal policy pair in a stochastic game are the policies at an NE point, which is defined as follows.

**Definition 2** (NE in Stochastic Game). *In a zero-sum stochastic game  $\Gamma$ , an NE point is an optimal policy pair  $\pi^* = \{\pi_u^*, \pi_a^*\}$ , such that for all state  $s \in \mathcal{S}$*

$$V^{\pi^*}(s) \geq V^{\pi^a}(s), \quad (14)$$

and

$$V^{\pi^*}(s) \leq V^{\pi^u}(s), \quad (15)$$

where  $\pi^a = \{\pi_u, \pi_a^*\}, \forall \pi_u$ , and  $\pi^u = \{\pi_u^*, \pi_a\}, \forall \pi_a$ .

In the context privacy stochastic game, the user aims to find the minimax equilibria, where the user tries to determine an optimal policy  $\pi_u^*$  that maximizes  $\{V^\pi(s) : \forall s\}$ , while the adversary tries to find an optimal policy  $\pi_a^*$  that minimizes  $\{V^\pi(s) : \forall s\}$ . Thus, based on (13), we have

$$V^{\pi^*}(s) = \max_{\pi_u} \min_{\pi_a} \left\{ r_u(s, \mathbf{a}_u^\pi, \mathbf{a}_a^\pi) + \gamma \sum_{s'} \Pr[s' | s, \mathbf{a}_u^\pi, \mathbf{a}_a^\pi] V^{\pi^*}(s') \right\}, \quad (16)$$

where  $V^{\pi^*}(s)$  is referred to as the value of state  $s$ .

It has been shown [24] that the equilibrium in a zero-sum stochastic game is the unique minimax equilibrium, and thus the optimal policy pair in the context privacy game is unique.

#### B. Equivalent State Value

Based on (31), the optimal policy pair can be derived via existing *reinforcement learning* algorithms, e.g. minimax-Q learning [25]. However, since cardinality of  $\mathcal{S}$  could be very large, the complexity of deriving  $\pi^*$  according to (31) would be very high. For example, the minimax-Q learning needs to solve  $|\mathcal{S}|$  bimatrix games, where  $|\mathcal{S}|$  is the cardinality of  $\mathcal{S}$ . In order to reduce the computational complexity, we solve an equivalent problem instead.

The equivalent state value  $\tilde{V}_u^{\pi^*}(Ar)$  is defined to be the expected state value over the context variable, i.e.,  $\tilde{V}_u^{\pi^*}(Ar) = \mathbb{E}_c[V_u^{\pi^*}(s)]$  where  $s = \{Ar, c\}$ . Then, we have the following observation.

**Lemma 1.** *The equivalent state value  $\tilde{V}_u^{\pi^*}(Ar)$  can be derived from (31) and enjoys an expression where context  $c$  is eliminated, i.e.,*

$$\tilde{V}^{\pi^*}(Ar) = \mathbb{E}_c \left[ r_u(s, \mathbf{a}^{\pi^*}) + \gamma \sum_{Ar'} \left( \Pr[Ar' | \mathbf{a}^{\pi^*}] \tilde{V}^{\pi^*}(Ar') \right) \right], \quad (17)$$

where  $\mathbf{a}^{\pi^*} = \{\mathbf{a}_u^{\pi^*}, \mathbf{a}_a^{\pi^*}\}$  is the action pair following the optimal policy pair  $\pi^*$ .

*Proof.* See Appendix A.  $\square$

We can see that  $\tilde{V}_u^{\pi^*}(Ar)$  largely reduces the number of state values from  $|\mathcal{S}|$  or  $2|\mathcal{C}|$  to 2 (since  $Ar$  is a binary variable). The following theorem proves that we can derive the optimal action pair from  $\tilde{V}_u^{\pi^*}(Ar)$ .

**Theorem 1.** *The optimal policy pair  $\pi^* = \{\pi_u^*, \pi_a^*\}$  for the context privacy stochastic game (31) can be obtained by solving the following equivalent problem*

$$\pi^* = \arg \max_{\pi_u} \min_{\pi_a} \left\{ r_u(s, \mathbf{a}^{\pi^*}) + \gamma \sum_{Ar'} \left( \Pr[Ar' | \mathbf{a}^{\pi^*}] \tilde{V}^{\pi^*}(Ar') \right) \right\}, \quad (18)$$

*Proof.* See Appendix B.  $\square$

#### C. Efficient Minimax Learning Algorithm

According to Theorem 1, we can derive  $\pi^*$  by learning  $\tilde{V}_u^{\pi^*}(Ar)$ , which can be obtained by the following update rule, which is modified from Q-learning [20].

$$\tilde{V}^{t+1}(Ar) = (1 - \alpha^{t+1}) \tilde{V}^t(Ar) + \alpha^{t+1} \mathbb{E}_c \left[ r_u(s, \mathbf{a}_u^t, \mathbf{a}_a^t) + \gamma \tilde{V}^t(Ar') \right], \quad (19)$$

---

**Algorithm 1** Minimax Learning Algorithm
 

---

**Input:** The context privacy stochastic game  $\Gamma$ 
**Output:**  $\pi^*$ 
*// 1. initialization*

- 1:  $t \leftarrow 0, Ar^t = 0$ ;
- 2:  $\tilde{V}^t(Ar = 0) \leftarrow 1, \tilde{V}^t(Ar = 1) \leftarrow 1$ ;
- 3: Initialize policy pair  $\pi^t$ : two uniform distributions where  $a_{u,i}^t = \frac{1}{K}, a_{a,i}^t = \frac{L}{K}, \forall i$ ;

*// 2. iteration*

- 4: **repeat**
  - 5:   Select an action pair  $\{\mathbf{a}_u^t, \mathbf{a}_a^t\}$  based on  $\pi^t$ ;
  - 6:   Update  $Ar^{t+1}$  after both players take their actions  $\{\mathbf{a}_u^t, \mathbf{a}_a^t\}$ ;
  - 7:   Update equivalent state value  $\tilde{V}^{t+1}(Ar)$  according to (19);
  - 8:   Update optimal policy  $\pi^{t+1}$  according to (18) with updated state values;
  - 9:    $t \leftarrow t + 1$ ;
  - 10: **until** Converge
- 

where  $\alpha^t \in [0, 1)$  is the learning rate, which needs to decay over time in order for the learning algorithm to converge. In this paper, we set  $\alpha^t = \frac{1}{t}$ .  $\tilde{V}_u^{t+1}(Ar)$  is used as an approximate of  $\tilde{V}_u^{\pi^*}(Ar)$  and iteratively updates according to (19) until converges.

Then, the learning algorithm for equivalent state value  $\tilde{V}_u^{\pi^*}(Ar)$  is described in Algorithm 1. First we initialize equivalent state values to be 1, and the policy of each player to be uniform distribution. Then, we iteratively update equivalent state values and policy pair according to (19) and (18), respectively, until the policy pair approaches the optimal policy pair.

In the following, we validate Algorithm 1 by proving that the iteratively updated  $\pi^t$  converges to the optimal policy pair. First, we show the convergence of Algorithm 1 by the following lemma.

**Lemma 2.** *In Algorithm 1,  $\tilde{V}^{t+1}(Ar)$  converges to  $\mathbb{E}_c [r_u(s, \mathbf{a}_u^t, \mathbf{a}_a^t) + \gamma \tilde{V}^t(Ar^t)]$ .*

*Proof.* See Appendix C.  $\square$

Next, we prove that the convergence point of Algorithm 1 is the true NE point.

**Theorem 2.** *In Algorithm 1, the equivalent state value  $\tilde{V}^{t+1}(Ar)$  updated by Line 7 converges to the NE  $\tilde{V}^{\pi^*}(Ar)$  defined by (17), and the corresponding optimal policy pair  $\pi^*$  is the unique NE solution for the context privacy stochastic game.*

*Proof.* See Appendix D.  $\square$

#### D. Properties of Optimal Policies

Note that in order to obtain the equivalent state value  $\tilde{V}^t(Ar)$  at stage  $t$ , the user needs to solve the equilibrium of a stage game, where the value of the game can be expressed by (19), and the equilibrium can be derived by solving the

minimax problem (18). By substituting (2) (5) (6) (8) (9) (7) into (18), we write the minimax problem as follows.

$$\max_{\pi_u} \min_{\pi_a} \left\{ \frac{1}{1 + e^{-\theta(\sum_i \kappa_i a_{u,i}^t - \eta)}} - C(c) \sum_i \kappa_i a_{u,i}^t a_{a,i}^t \right\}, \quad (20a)$$

$$\text{s.t.} \quad \sum_i a_{a,i}^t \leq L \quad (20b)$$

$$0 \leq a_{a,i}^t \leq 1, \forall i \quad (20c)$$

$$0 \leq a_{u,i}^t \leq 1, \forall i \quad (20d)$$

where  $C(c)$  is a function of the context  $c$  that  $C(c) = \omega \text{Sens}(c) + \gamma (\tilde{V}^{\pi^*}(Ar' = 0) - \tilde{V}^{\pi^*}(Ar' = 1))$ . Note that  $\tilde{V}^{\pi^*}(Ar' = 0)$  and  $\tilde{V}^{\pi^*}(Ar' = 1)$  are constant for a certain user, and thus only depends on the sensitivity of the context  $c$ . According to the definition of state value, i.e., Eq. (12), it can be seen that  $\tilde{V}^{\pi^*}(Ar' = 1) < \tilde{V}^{\pi^*}(Ar' = 0)$ . Since  $\text{Sens}(c) \geq 0$ , we can see that  $C(c) > 0$ .

To solve such minimax problem, we first assume that  $\pi_u$  is fixed. Since  $C(c) > 0$ , the adversaries will choose to attack  $L$  sensors with largest  $\kappa_i a_{u,i}^t$  so as to minimize (20a). Then, the problem becomes

$$\max_{\pi_u, \chi, \mathcal{I}} \left\{ \frac{1}{1 + e^{-\theta(\sum_i \kappa_i a_{u,i}^t - \eta)}} - C(c) \sum_{i \in \mathcal{I}} \kappa_i a_{u,i}^t \right\}, \quad (21a)$$

$$\text{s.t.} \quad a_{u,i}^t \geq \chi, \forall i \in \mathcal{I} \quad (21b)$$

$$a_{u,j}^t \leq \chi, \forall j \in \{1, \dots, K\} \setminus \mathcal{I} \quad (21c)$$

$$0 \leq a_{u,i}^t \leq 1, \forall i \quad (21d)$$

where  $\mathcal{I}$  is a subset of  $K$  that contains  $L$  sensors with largest  $\kappa_i a_{u,i}^t$ . Given a certain  $\mathcal{I}$ , the closed-form expressions for the optimal  $\pi_u$  can be easily derived. Based on (21a), we have the following observation on the NE of the context privacy stochastic game.

**Proposition 1.** *The optimal policy of the adversary is independent of the context sensitivity and the state values, but depends on the sensor's weights  $\{\kappa_i\}$ , while the optimal policy of the user depends on the context sensitivity, the state values and the sensor's weights.*

#### E. Optimal Policies Against The Adversary With Unlimited Power

The above analyses are based on the assumption that the adversary's bandwidth or computational power limits the amount of sensing data it can access in each time slot. Now we discuss the case of adversaries with unlimited power, who can access the whole sensing data released by the user. In this case, the optimal policy can still be learned by Algorithm

1, while for each iteration (each stage game) in the learning process, the minimax problem (20) is changed to

$$\max_{\pi_u} \left\{ \frac{1}{1 + e^{-\theta(\sum_i \kappa_i a_{u,i}^t - \eta)}} - C(c) \max_{\pi_a} \sum_i \kappa_i a_{u,i}^t a_{a,i}^t \right\}, \quad (22a)$$

$$\text{s.t. } 0 \leq a_{a,i}^t \leq 1, \forall i \quad (22b)$$

$$0 \leq a_{u,i}^t \leq 1, \forall i \quad (22c)$$

where the  $L$  constraint for  $\mathbf{a}_a^t$  is removed. It can be easily  $\max_{\pi_a} \sum_i \kappa_i a_{u,i}^t a_{a,i}^t = \sum_i \kappa_i a_{u,i}^t$ . Hence, the objective function (22a) can be written as

$$\max_{\pi_u} \left\{ \frac{1}{1 + e^{-\theta(\sum_i \kappa_i a_{u,i}^t - \eta)}} - C(c) \sum_i \kappa_i a_{u,i}^t \right\}, \quad (23)$$

where the closed-form expression for the optimal policy  $\pi_u$  can be derived.

#### F. Integration with Smartphones

Our solution can be implemented as a middleware in smartphones to sanitize raw sensor data and releases the sanitized data to upper layer applications. We can leverage the sandbox mechanism in today's smartphone platforms, such as Android and iOS, to implement such a middleware. In particular, we can build a sandbox to confine all untrusted applications, and use the sandbox to sanitize raw sensor data according to certain privacy preserving mechanisms before providing the data to applications.

#### G. Optimal Strategy Under Offline Attacks

Through our investigations in Section IV, we obtain the optimal defense strategy under online attacks. In this section, we study the optimal defense strategy under offline attacks.

Different from the case of online attacks, the user cannot observe the attack results under offline attacks, and takes actions based on its own context transitions. The system state under offline attacks is  $S_{\text{off}}^t = C^t$ . Therefore, the context privacy problem becomes an optimization problem whose objective is to maximize the user's utility:

$$\begin{aligned} \max_{\mathbf{a}_u^t} U_u &= \max_{\mathbf{a}_u^t} \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_u(S^t, \mathbf{a}_u^t, \mathbf{a}_a^t) \right], \\ &= \max_{\mathbf{a}_u^t} \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t (QoS(\mathbf{a}_u^t) - \omega \cdot Pri_{\text{off}}(S^t, \mathbf{a}_u^t)) \right], \\ &= \max_{\mathbf{a}_u^t} \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t \left( QoS(\mathbf{a}_u^t) \right. \right. \\ &\quad \left. \left. - \omega \cdot Sens(C^t) \max_{\mathbf{a}_a^t} \left\{ \sum_i \kappa_i a_{u,i}^t a_{a,i}^t \right\} \right) \right]. \quad (24) \end{aligned}$$

To derive an optimal policy that maximizes the user's utility under offline attacks, we have the following observation.

**Theorem 3.** *The policy that maximizes each stage's payoff is the optimal policy that maximizes the user's utility. That is,*

$$\pi_u^* = \operatorname{argmax}_{\pi_u} \{r_u(s, \mathbf{a}_u^\pi, \mathbf{a}_a^\pi)\}. \quad (25)$$

*Proof.* See Appendix E.  $\square$

Then, the optimal policy can be obtained by solving a stage minimax problem similar to (20). Note that in the case of adversaries with unlimited power, the optimal policies still conform to Theorem 3.

## V. EVALUATION

In this section, we conduct trace-driven simulations to evaluate the smartphone user's payoffs under the privacy attack. Specifically, we construct the user's behaviors and contexts based on real dataset, while simulating the adversary's actions based on the model as defined in Section III. First, we show the proposed algorithm largely improves the convergence speed compared with the traditional learning algorithm. Then, we demonstrate the effectiveness of the proposed algorithm by comparing the sum of discounted payoffs when the user adopts different strategies. We also study how the user's utility and strategies are affected by some system parameters.

#### A. Setup

The user model, system parameters, and baselines used for evaluation are described as follows.

- **User Model.** We evaluate the performance of our proposed algorithm using the Reality Mining dataset<sup>1</sup>, which was collected by the MIT Media Laboratory from September 2004 to June 2005 [26]. The Reality Mining dataset is one of the most complete and large-scale smartphone mobility traces, and is actively be used as the real-world traces for human mobility model. It records the continuous activities of 94 students and staff at MIT equipped with Nokia 6600 smartphones, which are pre-installed with several pieces of software that collects data about call logs, Bluetooth devices in proximity of approximately five meters, location at granularity of cell tower, application usage, transportation model (e.g., driving, walking, stationary), etc. The total length of all subjects' traces combined is 266,200 hours, with average, minimum, and maximum length being 122 days, 30 days, and 269 days, respectively. As location is the most complete and fine-grained context in the dataset, we select location traces as the user's contexts in our evaluation. The average, minimum, and maximum numbers of locations per user is 19, 7, and 40, respectively. Based on the location traces, we train a Markov chain for each user. We use the first half of each user's trace as the training set, and the other half as the testing set. We construct a Markov model for each user based on the transition probabilities computed from the training set. Then, we simulate user's behaviors based on the trained Markov

<sup>1</sup><http://realitycommons.media.mit.edu/realitymining.html>

chain. For each user, a certain percentage  $p$  of contexts are selected as sensitive contexts.

- **System Parameters.** Unless explicitly otherwise stated, we use the following system parameters in our simulations. For each user, the percentage of sensitive contexts  $p$  is set to 0.5, satisfaction threshold  $\eta$  set to 0.7, QoS steepness  $\theta$  set to 10, the discount factor  $\gamma$  set to 0.8. According to [11], there are three sensors (i.e., GPS, WiFi, and Bluetooth) used to identify user's location contexts. Thus, we set the number of sensors needed to identify the context to 3, and the power limitation of the adversary  $L$  is set to 2. The weights of sensors  $\{\kappa_i : i = 1, \dots, K\}$  are set to the normalized values drawing from a uniform distribution.
- **Baselines.** We compare the convergence speed of the proposed algorithm and that of the traditional learning algorithm that learns state values directly according to (31). We also compare the performance of users adopting different strategies. We compare the optimal policies obtained by the proposed algorithm (denoted by *proposed*) with *fixed* strategy and *myopic* strategy. The fixed strategy draws an action that uniformly sets the granularity of each sensor to  $\frac{1}{K}$ . And the myopic strategy is the optimal policy obtained by myopic learning, where the effects of current actions on the future payoffs are ignored, i.e.,  $\gamma$  is considered to be 0 in the myopic learning.

## B. Results

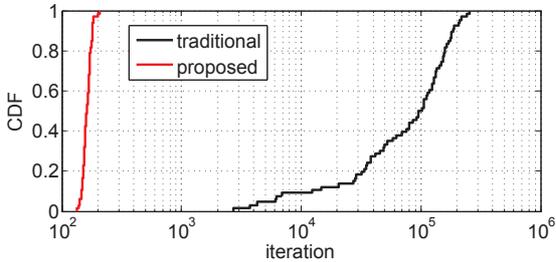


Fig. 3. CDF of convergence speed.

1) *Convergence Speed:* We first show the convergence speed improved by the proposed algorithm in Fig. 3. Fig. 3 depicts the cumulative distribution function (CDF) of iterations needed to learn the optimal policies for all users in the Reality Mining dataset. We can see that the convergence speed of the proposed algorithm for all users are less than 220 iterations, while for more than half of the users, the convergence speed of the traditional algorithm are more than  $10^5$  iterations, which demonstrates that the proposed algorithm largely improves the convergence speed compared with traditional learning algorithm. The improvement of the proposed algorithm comes from the smaller cardinality of the equivalent state value, which eliminate the context dimension in the learning process.

2) *Comparison of Different Strategies Under Online Attacks:* Fig. 4 compares the performance of the smartphone user when it adopts different strategies to evaluate the proposed algorithm under online attacks. It is assumed that the

adversaries use their optimal stationary policy learned by the minimax algorithm. As shown in Fig. 4, the proposed and the myopic strategies achieve higher sum of discounted payoffs than the fixed strategy against the adversaries with different power limitations, since the former two strategies maximize the worst-case performance, while the fixed strategy takes actions without considering the adversary's actions. Moreover, the proposed strategy achieves highest sum of discounted payoff. This is because the proposed strategy also takes the future payoff into consideration when optimizing the current strategy. Therefore, when smartphone users are under attack from adversaries that are capable of dynamically changing their strategies, the best choice is to adopt the strategy learned from the proposed algorithm that considers future payoff and the dynamics of the adversaries.

Moreover, comparing Fig. 4(a), Fig. 4(b), and Fig. 4(c), we can see that sum of discounted payoff achieved by the proposed strategy goes down as the power limitation of the adversaries  $L$  increases. This is because as  $L$  increases, the adversaries are able to access more sensing data, it is more likely for the adversaries to successfully attack the user. In such situation, the user may take more conservative actions (i.e., releasing data with less granularity), which results in lower service quality, or the user take the same action to preserve service quality, which, however, causes more privacy loss. As such, either case leads to lower payoff.

In the following, we show how the percentage of sensitive contexts and satisfaction threshold affect the sum of discounted payoff, and we also depict the optimal policies in different contexts. These evaluation results can provide some guidance in the design of the context privacy preserving schemes.

The average sums of discounted payoff of all users are reported in Fig. 5 and Fig. 6. From Fig. 5, we can see that the sums of discounted payoff achieved by the proposed and myopic strategies get lower as the percentage of sensitive contexts increases, since it would cost more privacy loss to release the same amount of data when the users have more sensitive contexts. The sum of discounted payoff achieved by the fixed strategy stays relatively the same over different percentage of sensitive contexts, because the service quality is invariant and dominates the payoff when adopting the fixed strategy. Moreover, the gap between the sums of discounted payoff obtained by adopting the proposed and myopic strategies approaches to zero when the percentage of sensitive contexts goes down. This is because the consideration of future payoff only affect the weights of privacy loss in the sum of discounted payoff according to (20a), and both strategies pay more attention to the service quality part when there are fewer sensitive contexts, which reduces the difference between with (the proposed strategy) and without (the myopic strategy) consideration of future payoff. This observation can provide some guidance for the context privacy preserving schemes that for the users with a small fraction of sensitive contexts, the impact of current actions on the future payoff can be neglected so as to design more efficient algorithm.

Fig. 6 reports the sums of discounted payoff achieved by different strategies over the applications with different satisfaction thresholds. It can be seen that for the applications

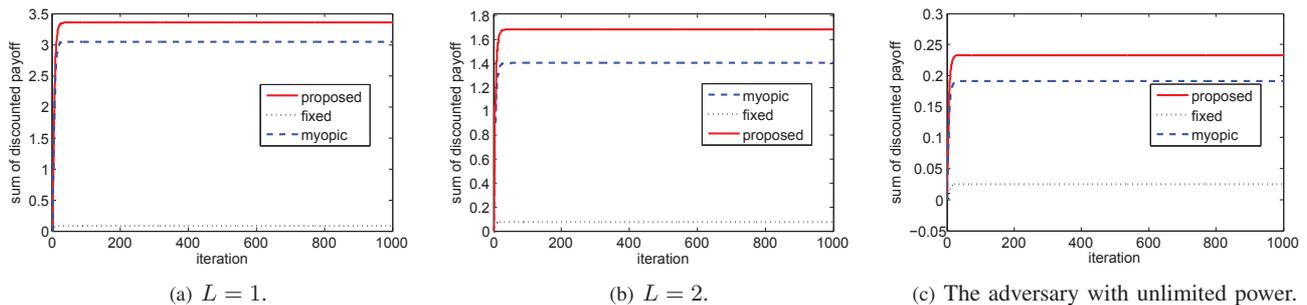


Fig. 4. Sum of discounted payoff of different strategies.

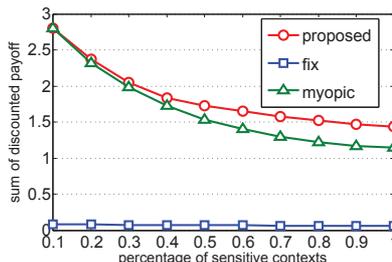


Fig. 5. Sum of discounted payoff vs. percentage of sensitive context.

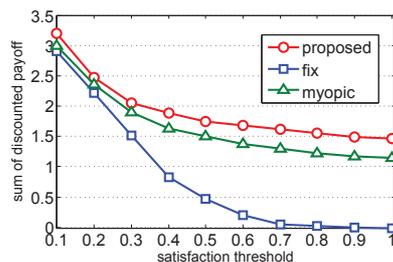


Fig. 6. Sum of discounted payoff vs. satisfaction threshold.

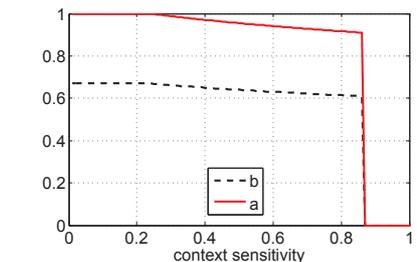


Fig. 7. Optimal policies in contexts of different sensitivities.

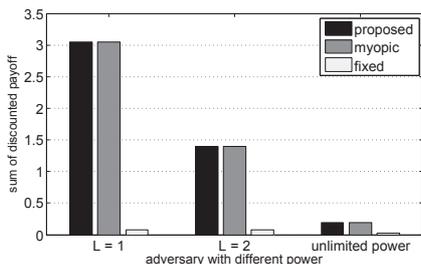


Fig. 8. Optimal policies in contexts of different adversaries.

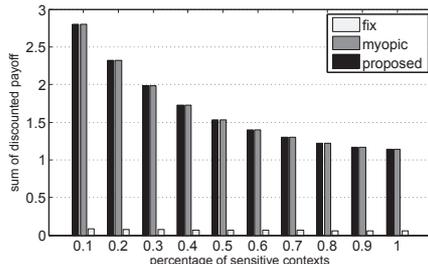


Fig. 9. Sum of discounted payoff vs. percentage of sensitive context.

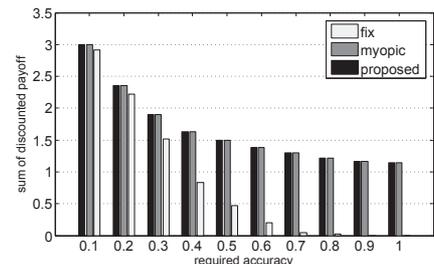


Fig. 10. Sum of discounted payoff vs. satisfaction threshold.

with higher satisfaction threshold, the sums of discounted payoff achieved by all strategies go down, since the service quality is lower as satisfaction threshold increases. We can also see that when the satisfaction threshold is very low, say 0.2, the differences in the sums of discounted payoff are achieved by different strategies are quite small. It can be seen according to (5) (6) (9) that with low satisfaction threshold, high service quality is easily achieved with only slight privacy loss by contributing a small amount of data, which are the cases of adopting the proposed and the myopic strategies. As such, the service quality dominates the payoff and stays relatively the same over different strategies. Thus, the privacy leaked by the applications that require high accuracy is hard to preserve, and the privacy preserving schemes need to be carefully designed to find a good tradeoff between privacy and utility since different strategies have significant impact on the user's total payoff.

Next, we study the optimal strategy in contexts with different sensitivities. To control the value of context sensitivity, we use the average values of the state values of all users

as the state values. We denote the total amount of released data  $a = \sum_i \kappa_i a_{u,i}^t$  and the amount of information leaked to the adversary  $b = \sum_i \kappa_i a_{u,i}^t a_{a,i}^t$ , which represent the optimal strategies. Fig. 7 depicts the variance of optimal  $a, b$  obtained by the proposed algorithm when the users are in different contexts. It can be seen that when the sensitivity of current context is smaller than 0.25 or larger than 0.87, the optimal  $a$  equals to 1 or 0, respectively. In such cases, either the variance of the service quality or the variance of the privacy loss dominates. While when the context sensitivity falls between 0.25 and 0.87,  $a$  and  $b$  slightly go down with the increment of the context sensitivity, since the users choose a more conservative strategy (releasing less data) as the privacy values more to the users. An interesting observation is that  $a$  stays larger than the satisfaction threshold (set to 0.7 by default) when the context sensitivity falls between 0.25 and 0.87. The reason is that below the satisfaction threshold the user enjoys only very limited service quality. Therefore, it is very important to identify the satisfaction threshold when designing the privacy preserving schemes.

3) *Comparison of Different Strategies Under Offline Attacks*: Fig. 8-10 compare the performance of different strategies under offline attacks. The results show that the proposed strategy achieves the same performance as the myopic strategy, which conforms to Theorem 3, that is the optimal policies are obtained when each stage payoff is maximized.

Fig. 8 shows that both the proposed and myopic strategies achieve much higher sum of discounted payoffs than the fixed strategy in the cases of adversaries with different power limitations, which demonstrates the benefits of considering the adversary's actions. In addition, the performance of all strategies against the adversary with unlimited power is worse compared with the adversaries with limited power. This is because the adversary with unlimited power gains more information, which forces the user takes more conservative strategies to minimize privacy leakage.

Fig. 9 reports the sum of discounted payoff under different percentage of sensitive context. The results show that the sum of discounted payoff diminishes when the percentage of sensitive context grows larger, as more sensitive users have to release coarser data to protect their privacy. Fig. 10 depicts the performance of all strategies in applications with different satisfaction thresholds. The results show that when the satisfaction threshold grows, the performance of the fixed strategy drops significantly while the proposed and myopic strategies drop much slower, which shows that the proposed and myopic strategies perform well in applications of different satisfaction thresholds.

## VI. RELATED WORK

**Privacy preservation techniques.** Numerous techniques have been proposed for preserving privacy in LBSs and participatory sensing on mobile phone. Spatial cloaking and anonymization are widely adopted [7], [8], [21], [22], where a value provided by a user is indistinguishable from those of  $k - 1$  other users to provide privacy guarantee, known as  $k$ -anonymity. Gedik et al. [7] devise a framework which provides  $k$ -anonymity with different context-sensitive personalized privacy requirements. Several clique-cloak algorithms are proposed in [7] to implement the framework by constructing a constraint graph. In [8], locality-sensitive hashing is utilized to partition user locations into groups that contain at least  $k$  users. A form of generalization based on the division of a geographic area is adopted by *Anonymsense* [21], where a map of wireless LAN access points is partitioned. *KIPDA* [22] enables  $k$ -anonymity for data aggregation with a maximum or minimum aggregation function in wireless sensor networks. However, these privacy techniques focus on the single shot scenario, which do not protect user's privacy against adversaries knowing temporal correlations.

Differential privacy has been considered as a major axis in data publishing. Publishing different types of data has been studied, such as histogram [27], [28], set-valued data [29], decision trees [30], as well as complex data format [31]. Among these studies, the data type related to our work is histogram. Blum et al. [27] divides the input counts into bins of roughly the same count to construct a one-dimensional

histogram. By observing that the accuracy of a differential privacy compliant histogram depends heavily on its structure, Xu et al. [28] propose two algorithms with different priorities for information loss and noise scales. Wang. et al. [31] propose a differential privacy based framework to outsource health data to hybrid cloud with personalized protection. However, these techniques focus on data modifications but do not environmental dynamics and adversaries' adjustable strategies.

Another category preserves privacy via cryptographic techniques. Girao et al. [32] aggregate data based on homomorphic encryption, which preserves privacy by performing certain computations on ciphertext. The limitation of homomorphic encryption is that a server must know all the users that have reported data to compute the final aggregated results. Secure information aggregation frameworks are proposed in [33]. However, the cryptographic techniques fail to cope with context privacy since the adversaries can decode the true sensing data by compromising context-aware applications.

**Context privacy.** Recent studies have investigated the context-related privacy leakage on smartphones. MaskIt [12] is a middleware that employs a privacy check to decide whether to release or suppress the current user context. As such, MaskIt limits the adversaries from knowing the user being in a sensitive context even when the adversaries have knowledge about the temporal correlation between user's contexts. Nevertheless, MaskIt does not consider the adversaries' capability of adjusting their attacking strategies. *CQue* focuses on modeling fine-grained correlations among contexts rather than providing specific privacy guarantees, which is orthogonal to our work. Context-related issues have also been discussed in specific applications. Context information is leveraged in [34] to generate fingerprints for secure pairing of users' co-located devices. Zhu et al. [35] develop a novel context-free attack to infer the keystrokes on smartphones using the Time Difference of Arrival (TDoA) of acoustic emanations.

**Location Privacy.** Our work is closely related to location privacy in LBSs. Homomorphic encryption is leveraged in [36] to allow the provider answer encrypted queries without knowing the location information. Similarly, Li et al. [37] utilize homomorphic encryption to enable Wi-Fi fingerprint-based localization without leaking users' locations. Tao et al. [38] defend against adversaries that are capable of inferring a user's location using localization techniques. These studies have focused on single-shot location privacy, while the temporal correlations in different contexts are not considered.

**Game theoretic analyses on privacy preservation.** Several game theoretic analyses on location privacy have also been discussed. Freudiger et al. [9] study the problem of selfishness in location privacy schemes based on pseudonym changes, and analyze the non-cooperative behavior of mobile nodes with an  $n$ -player complete information game. Shokri et al. [10] formulate the location privacy problem as Stackelberg Bayesian games with the consideration of user's service quality and adversary's cost. However, these location privacy problems are quite different from the context privacy discussed in this paper, where the stochastic dynamics and temporal correlation of user's behaviors and environments are considered.

## VII. CONCLUSION

This paper studied the privacy problem of context-aware applications on smartphones. Considering the distinct features of context privacy problem including the context dynamics and adversaries with knowledge of temporal correlations between contexts and capabilities of adjusting their attacking strategies, we formulate the interactive competition between users and adversaries as a competitive MDP, in which the users aim to maintain the context-based service quality and their context privacy by deciding the data granularity of each sensor that are accessed by the context-aware applications. On the other hand, the adversaries adjust their strategies on which sensing data are selected as the source to launch attacks. To obtain the optimal policy of the users efficiently under offline attacks, we propose a minimax learning algorithm to solve an equivalent problem with reduced dimensions. The proposed algorithm is proved to converge to the unique NE point of the stochastic game. In addition, we discuss the optimal policy under offline attacks.

We have conducted evaluation on real smartphone traces to demonstrate the effectiveness of the optimal policy obtained by the proposed algorithm. The results show the merits of considering the temporal correlations and future impacts. In addition, new observations about how user sensitivity and satisfaction threshold affect user's utility can provide some guidelines to the design of privacy preserving mechanisms for context privacy protection.

## ACKNOWLEDGEMENT

The research was supported in part by grants from 973 project 2013CB329006, China NSFC under Grant 61502114, China NSFC under Grant 61173156, RGC under the contracts CERG 622613, 16212714, HKUST6/CRF/12R, and M-HKUST609/13, as well as the grant from Huawei-HKUST joint lab.

## REFERENCES

- [1] [Online]. Available: <http://geonotehelp.blogspot.hk>
- [2] [Online]. Available: <http://nikeplus.nike.com/plus/products>
- [3] [Online]. Available: <http://www.novniv.com>
- [4] Microsoft, "Location based services usage and perceptions survey," Apr. 2011. [Online]. Available: <http://www.microsoft.com/en-hk/download/details.aspx?id=3250>
- [5] D. Weitzner, "Obama administration calls for a consumer privacy bill of rights for the digital age," Feb. 2012. [Online]. Available: <http://www.whitehouse.gov/blog/2012/02/23/we-can-t-wait-obama-administration-calls-consumer-privacy-bill-rights-digital-age>
- [6] W. Enck and et al., "Taintdroid: an information-flow tracking system for realtime privacy monitoring on smartphones," in *Proc. OSDI*, Oct. 2010, pp. 1–6.
- [7] B. Gedik and L. Liu, "Location privacy in mobile systems: A personalized anonymization model," in *Proc. IEEE ICDCS*, Jun. 2005.
- [8] K. Vu, R. Zheng, and J. Gao, "Efficient algorithms for k-anonymous location privacy in participatory sensing," in *Proc. IEEE INFOCOM*, Mar. 2012.
- [9] J. Freidiger, M. Manshaei, J. Hubaux, and D. Parkes, "On non-cooperative location privacy: a game-theoretic analysis," in *Proc. ACM CCS*, Nov. 2009, pp. 324–337.
- [10] R. Shokri, G. Theodorakopoulos, C. Troncoso, J. Hubaux, and J. Le Boudec, "Protecting location privacy: Optimal strategy against localization attacks," in *Proc. ACM CCS*, Oct. 2012, pp. 617–627.
- [11] S. Nath, "ACE: exploiting correlation for energy-efficient and continuous context sensing," in *Proc. ACM MobiSys*, 2012.
- [12] M. Götz, S. Nath, and J. Gehrke, "MaskIt: Privately releasing user context streams for personalized mobile applications," in *Proc. ACM SIGMOD*, May 2012, pp. 289–300.
- [13] A. Parate, M.-C. Chiu, D. Ganesan, and B. M. Marlin, "Leveraging graphical models to improve accuracy and reduce privacy risks of mobile sensing," in *Proc. ACM MobiSys*, 2013.
- [14] E. Kim, S. Helal, and D. Cook, "Human activity recognition and pattern discovery," *IEEE Perv. Comp.*, vol. 9, no. 1, pp. 48–53, 2010.
- [15] A. Mannini and A. Sabatini, "Accelerometry-based classification of human activities using markov modeling," *Computational Intelligence and Neuroscience*, 2011.
- [16] E. Toch and et al., "Empirical models of privacy in location sharing," in *Proc. ACM Ubicomp*, Sep. 2010, pp. 129–138.
- [17] M. Gruteser and D. Grunwald, "Anonymous usage of location-based services through spatial and temporal cloaking," in *Proc. ACM MobiSys*, May 2003, pp. 31–42.
- [18] Y. Wang and et al., "L2P2: Location-aware location privacy protection for location-based services," in *Proc. IEEE INFOCOM*, Mar. 2012, pp. 1996–2004.
- [19] A. Ghosh and A. Roth, "Selling privacy at auction," *Elsevier Games and Economic Behavior*, 2013.
- [20] J. Hu, M. Wellman et al., "Multiagent reinforcement learning: Theoretical framework and an algorithm," in *Proc. ICML*, 1998, pp. 242–250.
- [21] M. Shin, C. Cornelius, D. Peebles, A. Kapadia, D. Kotz, and N. Triandopoulos, "Anonymsense: A system for anonymous opportunistic sensing," *J. Perv. Mobile Comput.*, vol. 7, no. 1, pp. 16–30, 2011.
- [22] M. Groat, W. He, and S. Forrest, "KIPDA: k-indistinguishable privacy-preserving data aggregation in wireless sensor networks," in *Proc. IEEE INFOCOM*, Apr. 2011.
- [23] H. Lin, M. Chatterjee, S. Das, and K. Basu, "Arc: an integrated admission and rate control framework for competitive wireless cdma data networks using noncooperative games," *IEEE Trans. Mobile Comput.*, vol. 4, no. 3, pp. 243–258, 2005.
- [24] B. Wang, Y. Wu, K. Liu, and T. Clancy, "An anti-jamming stochastic game for cognitive radio networks," *IEEE J. Sel. Areas Commun.*, vol. 29, no. 4, pp. 877–889, 2011.
- [25] M. Littman, "Markov games as a framework for multi-agent reinforcement learning," in *Proc. ICML*, Jul. 1994, pp. 157–163.
- [26] N. Eagle, A. S. Pentland, and D. Lazer, "Inferring friendship network structure by using mobile phone data," *Proceedings of the National Academy of Sciences (PNAS)*, vol. 106, no. 36, pp. 15 274–15 278, 2009.
- [27] A. Blum, K. Ligett, and A. Roth, "A learning theory approach to non-interactive database privacy," in *Proc. ACM STOC*, May 2008.
- [28] J. Xu, Z. Zhang, X. Xiao, Y. Yang, and G. Yu, "Differentially private histogram publication," in *Proc. IEEE ICDE*, Apr. 2012.
- [29] R. Chen, N. Mohammed, B. Fung, B. Desai, and L. Xiong, "Publishing set-valued data via differential privacy," in *VLDB*, vol. 4, no. 11, 2011.
- [30] A. Friedman and A. Schuster, "Data mining with differential privacy," in *Proc. ACM SIGKDD*, Jul. 2010.
- [31] W. Wang, L. Chen, and Q. Zhang, "Outsourcing high-dimensional healthcare data to cloud with personalized privacy preservation," *Elsevier Comput. Netw.*, vol. 88, pp. 136–148, 2015.
- [32] J. Giroa, D. Westhoff, and M. Schneider, "Cda: Concealed data aggregation for reverse multicast traffic in wireless sensor networks," in *Proc. IEEE ICC*, May 2005.
- [33] B. Przydatek, D. Song, and A. Perrig, "Sia: Secure information aggregation in sensor networks," in *Proc. ACM SenSys*, Nov. 2003.
- [34] M. Miettinen, N. Asokan, T. D. Nguyen, A.-R. Sadeghi, and M. Sobhani, "Context-based zero-interaction pairing and key evolution for advanced personal devices," in *Proc. ACM CCS*, 2014, pp. 880–891.
- [35] T. Zhu, Q. Ma, S. Zhang, and Y. Liu, "Context-free attacks using keyboard acoustic emanations," in *Proc. ACM CCS*, 2014, pp. 453–464.
- [36] X.-Y. Li and T. Jung, "Search me if you can: privacy-preserving location query service," in *IEEE INFOCOM*, 2013, pp. 2760–2768.
- [37] H. Li, L. Sun, H. Zhu, X. Lu, and X. Cheng, "Achieving privacy preservation in WiFi fingerprint-based localization," in *Proc. IEEE INFOCOM*, 2014, pp. 2337–2345.
- [38] T. Shu, Y. Chen, J. Yang, and A. Williams, "Multi-lateral privacy-preserving localization in pervasive environments," in *Proc. IEEE INFOCOM*, 2014, pp. 2319–2327.
- [39] D. Bertsekas, *Dynamic programming and optimal control*. Athena Scientific, 1995.
- [40] X. Cao, *Stochastic learning and optimization: a sensitivity-based approach*. New Yorks Springer, 2007.
- [41] M. Littman and C. Szepesvári, "A generalized reinforcement-learning model: Convergence and applications," in *Proc. ICML*, Jul. 1996, pp. 310–318.

[42] C. Szepesvári and M. Littman, "A unified analysis of value-function-based reinforcement-learning algorithms," *Neural computation*, vol. 11, no. 8, pp. 2017–2060, 1999.

#### APPENDIX A PROOF OF LEMMA 1

By taking expectation over  $c$  on both sides of (31), we have

$$\begin{aligned}
& \tilde{V}^{\pi^*}(Ar) \\
&= \mathbb{E}_c \left[ r_u(s, \mathbf{a}^{\pi^*}) + \gamma \sum_{s'} \Pr[s'|s, \mathbf{a}^{\pi^*}] V^{\pi^*}(s') \right] \\
&= \gamma \sum_{Ar', c'} \mathbb{E}_c \left[ r_u(s, \mathbf{a}^{\pi^*}) + \Pr[Ar'|\mathbf{a}^{\pi^*}] \Pr[c'|c] V_u^{\pi^*}(Ar', c') \right] \\
&= \mathbb{E}_c \left[ r_u(s, \mathbf{a}^{\pi^*}) \right] \\
&\quad + \gamma \sum_{Ar', c', c} \left( \Pr[Ar'|\mathbf{a}^{\pi^*}] \Pr[c'|c] \Pr[c] V_u^{\pi^*}(Ar', c') \right) \\
&= \mathbb{E}_c \left[ r_u(s, \mathbf{a}^{\pi^*}) \right] + \gamma \sum_{Ar', c'} \left( \Pr[Ar'|\mathbf{a}^{\pi^*}] \Pr[c'] V_u^{\pi^*}(Ar', c') \right) \\
&= \mathbb{E}_c \left[ r_u(s, \mathbf{a}^{\pi^*}) + \gamma \sum_{Ar'} \left( \Pr[Ar'|\mathbf{a}^{\pi^*}] \tilde{V}^{\pi^*}(Ar') \right) \right]. \quad (26)
\end{aligned}$$

This completes the proof.

#### APPENDIX B PROOF OF THEOREM 1

By standard Markov decision process (MDP) techniques [39], [40], the problem (31) can be expressed as an equivalent MDP  $\min_{\pi_a} \max_{\pi_u} \mathbb{E}_c[V_u^{\pi}(s)]$  with the state space  $\mathcal{S}$ , the action space  $\{\{\mathbf{a}_u\}, \{\mathbf{a}_a\}\}$ , the transition kernel  $\Pr[Ar'|\mathbf{a}^{\pi^*}] = \mathbb{E}_c[\Pr[s'|s]]$ , and the stage payoff function  $\mathbb{E}_c[r_u(s, \mathbf{a}^{\pi^*})]$ . It is known that the optimal policy pair  $\pi^*$  can be obtained by solving

$$\begin{aligned}
\min_{\pi_a} \max_{\pi_u} \mathbb{E}_c[V^{\pi}(s)] &= \mathbb{E}_c \left[ \min_{\pi_a} \max_{\pi_u} \left\{ r_u(s, \mathbf{a}^{\pi^*}) \right. \right. \\
&\quad \left. \left. + \gamma \sum_{Ar'} \left( \Pr[Ar'|\mathbf{a}^{\pi^*}] \tilde{V}^{\pi^*}(Ar') \right) \right\} \right], \quad (27)
\end{aligned}$$

which completes the proof.

#### APPENDIX C PROOF OF LEMMA 2

First, we show that both states  $Ar = 0$  and  $Ar = 1$  will be visited infinite often. Since  $\pi^t$  is updated according to the minimax problem (18), it is obvious that  $0 < \Pr[Ar^t = 1] < 1$  will appear infinite times [41], meaning that both  $Ar^t = 0$  and  $Ar^t = 1$  will appear infinite times. Since  $\alpha^t = \frac{1}{t}$ , we can see that  $\sum_{t=0}^{\infty} \alpha^{t+1} = \infty$  and  $\sum_{t=0}^{\infty} (\alpha^{t+1})^2 < 1 + \sum_{t=1}^{\infty} \frac{1}{t(t-1)} < 1 + \sum_{t=1}^{\infty} \left( \frac{1}{t-1} - \frac{1}{t} \right) < \infty$ .

Then, according to *conditional average lemma* [42], the process of updating  $\tilde{V}^{t+1}(Ar^{t+1})$  by (19) converges to the component with factor  $\alpha^{t+1}$ , which proves the convergence of Algorithm 1.

#### APPENDIX D PROOF OF THEOREM 2

From Lemma 2,  $\tilde{V}^{t+1}(Ar)$  converges to

$$\begin{aligned}
& \mathbb{E}_{c, Ar'} \left[ r_u(s, \mathbf{a}_u^t, \mathbf{a}_a^t) + \gamma \tilde{V}^t(Ar') \right] \\
&= \sum_{c, Ar'} \Pr[c] \Pr[Ar'|\mathbf{a}^{\pi^*}] \left( r_u(s, \mathbf{a}_u^t, \mathbf{a}_a^t) + \gamma \tilde{V}^t(Ar') \right). \quad (28)
\end{aligned}$$

Denote  $\Omega^t \tilde{V}^t(Ar) = \mathbb{E}_{c, Ar'} \left[ r_u(s, \mathbf{a}_u^t, \mathbf{a}_a^t) + \gamma \tilde{V}^t(Ar') \right]$ , and  $\Delta^t \tilde{V}^t(Ar) = r_u(s, \mathbf{a}_u^t, \mathbf{a}_a^t) + \gamma \tilde{V}^t(Ar')$ . It has been proven that  $\Delta^t$  is a contract mapping of  $\tilde{V}^t(Ar)$  [20], i.e., we have

$$\left\| \Delta^t \tilde{V}^t(Ar) - \Delta^t \tilde{V}^{\pi^*}(Ar) \right\| \leq \gamma \left\| \tilde{V}^t(Ar) - \tilde{V}^{\pi^*}(Ar) \right\|. \quad (29)$$

Since  $\Omega^t \tilde{V}^t(Ar) = \sum_{c, Ar'} \Pr[c] \Pr[Ar'|\mathbf{a}_u^t, \mathbf{a}_a^t] \Delta^t \tilde{V}^t(Ar)$  and  $\Pr[c] \Pr[Ar'|\mathbf{a}_u^t, \mathbf{a}_a^t] \geq 0$ , we can see that  $\Omega^t \tilde{V}^t(Ar)$  is also a contract mapping of  $\tilde{V}^t(Ar)$ .

Next, we show that the fixed point of  $\Omega^t$  is  $\tilde{V}^{\pi^*}(Ar)$ . According to (28), we have

$$\begin{aligned}
& \Omega^t \tilde{V}^{\pi^*}(Ar) \\
&= \sum_{c, Ar'} \Pr[c] \Pr[Ar'|\mathbf{a}^{\pi^*}] \left( r_u(s, \mathbf{a}_u^{\pi^*}, \mathbf{a}_a^{\pi^*}) + \gamma \tilde{V}^{\pi^*}(Ar') \right) \\
&= \mathbb{E}_c \left[ r_u(s, \mathbf{a}_u^{\pi^*}, \mathbf{a}_a^{\pi^*}) + \sum_{Ar'} \Pr[c] \Pr[Ar'|\mathbf{a}^{\pi^*}] \gamma \tilde{V}^{\pi^*}(Ar') \right] \\
&= \tilde{V}^{\pi^*}(Ar), \quad (30)
\end{aligned}$$

which proves that the fixed point of  $\Omega^t$  is  $\tilde{V}^{\pi^*}(Ar)$ .

Therefore, the equivalent state value  $\tilde{V}^{t+1}(Ar)$  updated by Line 7 converges to the NE  $\tilde{V}^{\pi^*}(Ar)$  defined by (17), in which the optimal policy pair  $\pi^*$  is an NE solution for the context privacy stochastic game.

Since it has been proven [24] that the equilibrium in a zero-sum game is the unique minimax equilibrium, and thus the optimal policy pair  $\pi^*$  is the unique NE solution for the context privacy stochastic game.

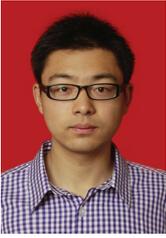
#### APPENDIX E PROOF OF THEOREM 3

According to Section IV-A, the utility can be expressed in the form of state values  $\{V^{\pi}(s) : \forall s\}$ , and thus the optimization problem becomes finding the optimal policy  $\pi_u^*$  that maximizes  $\{V^{\pi}(s) : \forall s\}$ .

$$\begin{aligned}
V^{\pi_u^*}(s) &= \max_{\pi_u} \left\{ r_u(s, \mathbf{a}_u^{\pi}, \mathbf{a}_a^{\pi}) + \gamma \sum_{s'} \Pr[s'|s, \mathbf{a}_u^{\pi}] V^{\pi^*}(s') \right\}, \\
&= \max_{\pi_u} \left\{ r_u(s, \mathbf{a}_u^{\pi}, \mathbf{a}_a^{\pi}) \right\} + \gamma \sum_{s'} \Pr[s'|s] V^{\pi_u^*}(s'). \quad (31)
\end{aligned}$$

As the term  $\gamma \sum_{s'} \Pr[s'|s] V^{\pi^*}(s')$  is independent of  $\pi_u^*$ , the optimal policy can be derived by

$$\pi_u^* = \operatorname{argmax}_{\pi_u} \left\{ r_u(s, \mathbf{a}_u^{\pi}, \mathbf{a}_a^{\pi}) \right\}. \quad (32)$$



**Wei Wang (S'10-M'15)** is currently a Research Assistant Professor in Fok Ying Tung Graduate School, Hong Kong University of Science and Technology (HKUST). He received his Ph.D. degree in Department of Computer Science and Engineering from HKUST. Before he joined HKUST, he received his bachelor degree in Electronics and Information Engineering from Huazhong University of Science and Technology, Hubei, China, in June 2010.



**Qian Zhang (M'00-SM'04-F'12)** joined Hong Kong University of Science and Technology in Sept. 2005 where she is a full Professor in the Department of Computer Science and Engineering. Before that, she was in Microsoft Research Asia, Beijing, from July 1999, where she was the research manager of the Wireless and Networking Group. She is a Fellow of IEEE for "contribution to the mobility and spectrum management of wireless networks and mobile communications". Dr. Zhang received the B.S., M.S., and Ph.D. degrees from Wuhan University, China, in 1994, 1996, and 1999, respectively, all in computer science.