

A Distributed k -Anonymity Protocol for Location Privacy

Ge Zhong and Urs Hengartner

Cheriton School of Computer Science, University of Waterloo

{gzhong, uhengart}@cs.uwaterloo.ca

Abstract—To benefit from a location-based service, a person must reveal her location to the service. However, knowing the person's location might allow the service to re-identify the person. Location privacy based on k -anonymity addresses this threat by cloaking the person's location such that there are at least $k - 1$ other people within the cloaked area and by revealing only the cloaked area to a location-based service. Previous research has explored two ways of cloaking: First, have a central server that knows everybody's location determine the cloaked area. However, this server needs to be trusted by all users and is a single point of failure. Second, have users jointly determine the cloaked area. However, this approach requires that all users trust each other, which will likely not hold in practice. We propose a distributed approach that does not have these drawbacks. Our approach assumes that there are multiple servers, each deployed by a different organization. A user's location is known to only one of the servers (e.g., to her cellphone provider), so there is no single entity that knows everybody's location. With the help of cryptography, the servers and a user jointly determine whether the k -anonymity property holds for the user's area, without the servers learning any additional information, not even whether the property holds. A user learns whether the k -anonymity property is satisfied and no other information. The evaluation of our sample implementation shows that our distributed k -anonymity protocol is sufficiently fast to be practical. Moreover, our protocol integrates well with existing infrastructures for location-based services, as opposed to the previous research.

I. INTRODUCTION

With the advance of location technologies, people can now determine their location in various ways, for instance, with GPS or based on nearby cellphone towers. These technologies have led to the introduction of location-based services, which allow people to get information relevant to their current location. Location privacy is of utmost concern for such location-based services, since knowing a person's location can reveal information about her activities or her interests.

In this paper, we focus on location-based services that need to know only a person's location, but not her identity. For example, these can be services that provide road maps, nearby places (e.g., restaurants or gas stations), or current traffic conditions. As it turns out, even if a service learns only a person's location, it might still be able to re-identify the person [1]. For example, the location could be associated with the person (e.g., her home), or the location corresponds to a place that is under physical surveillance by the location-based service. Once a service has re-identified a person, the service can literally connect the dots and build a detailed location profile for this person (assuming the person uses the service

in a continuous way).

Location cloaking is a defence against re-identification. It is based on the idea of sending coarse-grained location information that covers multiple people to a location-based service. In a naïve approach, a user simply determines an area of a static size (e.g., four city blocks) that contains her current location and sends the area's coordinates to the service. The service returns information about the entire area, and the user discards any irrelevant information. Unfortunately, this approach could still allow re-identification. For example, in a rural, less populated region, this kind of cloaking might well result in an area that includes only one user.

Location cloaking based on k -anonymity does not have this disadvantage. Here, a user's current location is cloaked such that there are at least $k - 1$ other users within the cloaked area. A location-based service learns only the cloaked area, which allows the user to remain anonymous within the set of k users. Applying k -anonymity to location cloaking has been studied extensively [1]–[13]. Traditionally, this approach has been implemented with the help of a central trusted server [1]–[4], [6], [9]–[13]. Here, users register their current location with the trusted server. Whenever a user wants to access a location-based service, she has the trusted server compute a cloaked area that has the k -anonymity property. Then, the trusted server contacts the location-based service on the user's behalf. The drawback of this approach is that the trusted server knows everybody's location. Users must trust it not to leak their location information to unauthorized parties, maybe inadvertently. In short, the trusted server is a single point of failure. More recent research has proposed to get rid of the trusted server and to have (nearby) users jointly compute a cloaked area that has the k -anonymity property [5], [7], [8]. Then, the user (or, for increased privacy, another user on her behalf) contacts the location-based service. The drawback of this approach is that all previous solutions trust users to implement the proposed solution faithfully and not to leak location information learned during the computation. Whereas this requirement might hold in a closed environment, where users know each other, it will be difficult to satisfy in more open environments.

Another drawback of both the centralized and the distributed approach is that neither of them integrates well with existing infrastructures for location-based services. Namely, many existing location-based services are targeted at cellphone users, since the operator of a cellphone network knows the current

location of *its* customers and can provide this information to a location-based service. However, there is no single entity that knows the location of *all* cellphone users across all cellphone networks, as required by the centralized approach. The distributed approach fails to take advantage of the already existing location information that an operator has about its customers.

We propose a solution that requires neither a single trusted server nor trust in all users of the system and that integrates well with existing infrastructures. Namely, we have multiple servers, each deployed by a different organization (e.g., an operator of a cellphone network) and each knowing the location of only a *subset* of users (e.g., the operator's customers), with the subsets being disjoint. When a user wants to access a location-based service, she cloaks her location and asks each server for the number of people in the cloaked area. In a naïve solution, the servers simply give her these numbers, she sums them up and, if the sum is at least k , she accesses the location-based service. However, this approach has the flaw that it might allow the user to track people. For example, if the user learns that there is only a single person in an area and nobody in the surrounding areas, the user can likely follow the path of the person when the person leaves the area and enters one of the surrounding areas. As soon as the person enters an area that is associated with her identity or that is under surveillance by the user, the user can re-identify the person. In general, sophisticated data-mining algorithms might allow the tracking or re-identification of a person even if there are multiple people in an area.

Our solution avoids this problem with the help of cryptography and ensures that a user cannot learn the number of people in an area reported by a server. The user can learn only whether the sum of these numbers is at least k . Our contributions are:

- First, we introduce a distributed k -anonymity protocol for location privacy in which a user collaborates with multiple servers and a third party to learn whether there are at least k people in her area. Nobody, not even the servers and the third party, can learn the total number of people in the area.
- Second, we present a protocol that prevents users from registering multiple times with different servers and hence from skewing the total number of users in an area.
- Third, we present a sample implementation of our protocol. In its evaluation, we demonstrate that our protocol can be implemented efficiently.

A preliminary overview of our architecture appeared in a workshop paper [14]. The workshop paper misses the key components of the architecture and omits the evaluation and security analysis of the architecture.

The rest of this paper is organized as follows: In Section II, we discuss related work in the area of k -anonymity and location privacy. In Section III, we present our system and threat model. We introduce our distributed k -anonymity protocol for location privacy in Section IV. In Section V, we present our defence against multiple registrations. We give a

security analysis in Section VI and evaluate our architecture in Section VII.

II. RELATED WORK

Samarati and Sweeney [15] propose k -anonymity to enable the release of person-specific information from a database while maintaining individuals' privacy. Previous research has applied k -anonymity to the release of location information that occurs when a user queries a location-based service. We first discuss related work that is based on a central trusted server, then we review distributed approaches.

Gruteser and Grunwald [1] introduce location privacy based on k -anonymity. A trusted "location anonymizer" cloaks a user's location by subdividing space into quadrants until it finds a quadrant that contains the query issuer and fewer than $k - 1$ other users. The parent quadrant becomes the cloaked area. Gedik and Liu [6] let users have personalized values of k , and the cloaked area corresponds to the minimum bounding rectangle of k users. Mokbel et al. [12] observe that this approach can leak information about a user's location (e.g., some users will be on the boundary of the rectangle). They use a balanced quadtree that is traversed bottom-up for better performance until a quadrant with at least k users is found. In our approach, we choose the bottom-up strategy and allow users to personalize k .

Beresford [3] finds that, if a location-based service is familiar with the cloaking algorithm and knows the locations of all users within the cloaked area, the service could infer the identity of the query issuer from the shape of the cloaked area. Namely, this happens when the cloaked area generated for the query issuer is different from the cloaked areas that would have been generated for the other users in the cloaked area. Kalnis et al. [9] and Bettini et al. [4] later re-discover this finding. Kalnis et al. and Mascetti and Bettini [11] present (centralized) cloaking algorithms that are not susceptible to this attack. In our approach, we leave it up to a user to decide what kind of cloaking algorithm to use. She can use either an algorithm similar to Mokbel et al.'s algorithm that does not necessarily guarantee her privacy, but is easy to compute, or an algorithm similar to Mascetti and Bettini's that is robust in terms of privacy, but more expensive.

Chow et al. [5] propose the first distributed approach for location k -anonymity. A user who wants to access a location-based service broadcasts a message with Bluetooth or WiFi. Nearby users respond to this message with their current location. If the number of responses is smaller than $k - 1$, the user repeats the process, but has the nearby users forward the message, maybe iteratively. The user then computes her cloaked location and, for increased privacy, asks a nearby user to send her query for the cloaked location to the location-based service. Ghinita et al. [7] show that this approach often fails to achieve location privacy, since the query issuer tends to be in the center of the cloaked area. The same authors [8] later show that their earlier approach can be slow and propose an approach based on a distributed hash table. Here, a user's 2-D location is mapped to a 1-D position, used as index in the hash

table, such that two users who are nearby in 2-D are likely also close in 1-D. A user’s 1-D position leaks information about her 2-D location. A user knows the positions of the two users that immediately follow and precede her in the 1-D sequence. Furthermore, for robustness reasons, a user also needs to know the positions of $\log_2(n)$ other users, where n is the number of users. In summary, the proposed distributed approaches for location k -anonymity have the drawback that nodes can learn location information about other nodes, so the nodes have to trust each other [8].

Kapadia et al. [10] propose “statistical k -anonymity”. They assume the global availability of statistical data about the number of people who are present in an area with high probability at a particular time of the day. When a user wants to access a location-based service, she independently decides based on this data whether her area is likely to be visited by at least k people. The drawbacks of this approach are that there remains a chance that fewer than k people are actually in the area and the requirement of extensive data collection (across different communication technologies and providers and during different times of the day, days of the week,...) to compute the provided statistical data, which raises privacy issues of its own. Our approach is always accurate and requires no such data collection.

Zhong et al. [16] study a scenario where database records are horizontally distributed among different sites. They present an algorithm that allows a data miner to learn the sensitive part of a record only if there are least $k - 1$ other records, maybe at different sites, whose non-sensitive part is identical to the non-sensitive part of the record in question. The data miner (i.e., the user in our problem setting) always learns the overall number of records that have a particular non-sensitive part, which makes the algorithm inapplicable to our problem.

k -anonymity is not the only approach that has been suggested for location privacy. Another option are pseudonyms, where a user assumes a pseudonym when contacting a location-based service. Previous work (e.g., by Beresford and Stajano [17] or by Jiang et al. [18]) has explored the challenges of pseudonym-based approaches, such as changing pseudonyms in an unlinkable way. k -anonymity-based and pseudonym-based approaches for location privacy complement each other; the former one is attractive for scenarios where a location is associated with a particular person, the other one for scenarios where locations are public and visited by many people.

III. SYSTEM AND THREAT MODEL

In this section, we present our system and threat model.

A. System Model

We present our system model in Figure 1. The figure omits actual location-based services, which a user would access once she learns that at least $k - 1$ other users are in her area, likely via a proxy or an anonymous communication channel to hide her identity from the service. A possible anonymous communication channel is Tor [19], which lets a user hide her

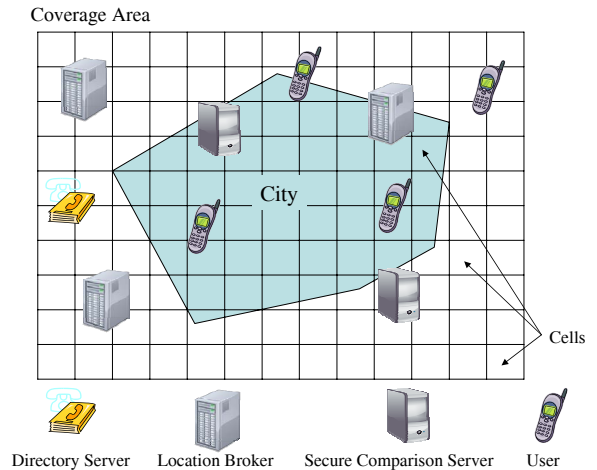


Fig. 1. System model. A user registers her location (cell) with a location broker, whose contact information is provided by the directory server. The user can learn whether there are at least k registered users in her cell by contacting all location brokers and one of the secure comparison servers.

IP address from a server by routing the user’s traffic through a series of intermediate nodes of the user’s choice. These intermediate nodes are computers maintained by individuals, and each node knows only its immediate predecessor and successor node.

For scalability reasons, there are multiple *coverage areas*, where a coverage area corresponds to the area covered by a particular instantiation of our system (e.g., a city or a province). A coverage area is divided into a well-known grid of equally sized, square cells. The width of a cell is chosen such that, for most cells, there is a realistic chance that multiple users can be located in the cell. For example, a cell could have a width of 100 meters. Moreover, there are four kinds of parties: location brokers, users, secure comparison servers, and a directory server. (Whereas there are solutions that do not require secure comparison servers, they are much less efficient than the one proposed here [20].)

A **location broker** keeps track of the current location (i.e., the current cell) of a *subset* of the users in the coverage area. There are multiple location brokers, each keeping track of the location of a *different* subset of users, with the intersection of any two subsets being empty. Each broker is maintained by a different organization. For example, the operator of a cellphone network could maintain a location broker that keeps track of the location of the operator’s customers in the coverage area. A location broker does not necessarily provide coverage for all cells in the coverage area. For example, whereas a broker maintained by a cellphone network operator would likely cover most cells, a broker operated by the provider of a WiFi network might provide coverage only for a small subset of the cells.

Users carry a mobile device (e.g., a cellphone or a laptop) with them that can locate itself (e.g., using GPS or nearby WiFi base stations). A user registers her current location (i.e., her current cell) with exactly one of the location brokers, where

she can choose with which one. Likely, if the provider of the communication service exploited by the user’s device runs a location broker, the user will (maybe implicitly) register her location with this broker, since the provider already knows or at least has an estimate of the user’s location. We assume that users always register their location with a broker. This assumption is also made in the earlier work. More registrations will lead to smaller cloaked areas, which in turn will increase the quality of service obtained from a location-based service and will give users an incentive to register. If a location broker is run by the provider of the user’s communication service, registering her location continuously does not lead to additional loss of privacy for the user, since the provider already has this information.

A **secure comparison server** interacts with a user to let the user learn whether there are at least k users who have registered the user’s current cell as their location across *all* location brokers. (See Section IV for the detailed protocol.) Each secure comparison server is maintained by a different organization. An organization can maintain both a location broker and a secure comparison server. A secure comparison server provides coverage for the entire coverage area.

The **directory server** publishes contact information for the location brokers and for the secure comparison servers in the coverage area. Moreover, it publishes coverage information for location brokers, that is, which broker provides coverage for which cells in the coverage area. This way, users can choose a location broker to register with and a secure comparison server to interact with.

B. Threat Model

In our threat model, the location brokers and the secure comparison servers are honest-but-curious, that is, they honestly follow our protocol, but are curious about learning location information. We discuss malicious brokers and servers in Section VI-B.

A **location broker** can learn the location (i.e., cell) of users who register their location with this particular broker. Accordingly, a broker can learn the number of users in a cell that have registered this cell as their location with this particular broker. However, a location broker should not learn the *total* number of users in a cell that have registered this cell as their location across *all* location brokers. Knowing this number makes possible tracking attacks, similar to the one presented in Section I. Similarly, a location broker should not learn the location of users who register their location with any other location broker. We assume that the organizations that run the location brokers do not collude with each other. Legal means (e.g., privacy laws or a contract between a user and a location broker) can enforce this assumption. Technical enforcement means make less sense here, since today’s cellphone network operators know their customers’ location and could potentially share this information with each other. For the same reason, we assume that location brokers do not collude with users.

A **secure comparison server** should learn neither a user’s location nor the total number of registered users in a cell. This

implies that the server should not learn the individual number for a location broker, either (except using back channels if a secure comparison server is run by the same organization as a location broker). A secure comparison server might collude with other secure comparison servers to learn additional information. Due to the same reason given above, we assume that secure comparison servers do not collude with location brokers (except in the implicit case where a broker and a server are run by the same organization, here it can learn at most the location and number of users registered with this broker).

A **user** should learn only her own location and whether the number of people in her cell (or superset of cells) is at least k , where k is a value of her choice. A user carries only one mobile device with her, and the device faithfully reports its location to a broker. A user cannot register multiple times with a single broker at the same time, since the broker authenticates the user’s device. All these assumptions are also made in the earlier work. (We discuss weaker ones in Section VI-C.) A user might still try to register multiple times with *different* brokers at the same time. This could let other users erroneously conclude that k -anonymity holds. Finally, a user might collude with a secure comparison server to learn the number of people in her cell (or superset of cells).

The **directory server** should not learn any location information about users. The server might misbehave, for example, it might list a location broker multiple times as providing coverage for a single cell, it might fail to vet location brokers or secure comparison servers (see Section IV-D), or it might try to track clients by providing them different information.

IV. DISTRIBUTED k -ANONYMITY PROTOCOL

In this section, we first give an overview of our distributed k -anonymity protocol and then present its key components.

A. Overview

The goal of a user is to learn whether there are at least k registered users (including herself) in the user’s *query area*, where k is a value chosen by the user and where the query area initially corresponds to the user’s current cell. If the user learns that there are fewer than k users in this cell, she can enlarge the query area to a superset of cells that contains the user’s current cell and re-execute the protocol for the enlarged area. This process can be repeated multiple times. As mentioned in Section II, a user can choose between different types of enlargement algorithms to determine the query area, which lets her trade off between privacy and cost. A user registers her current cell as her location with exactly one of the location brokers, but there is no need for the user to register additional cells when enlarging the query area.

To learn whether there are at least k users in her query area, a user first needs to identify the location brokers that provide coverage for (maybe parts of) the query area. The user must not ask the directory server for a list of brokers that provide this coverage, else the server could learn the user’s location. Instead, the user should download the entire directory (or recent changes to it) from the server on a regular basis,

such as once a day. The directory is signed, which allows retroactive detection of misbehaviour by the directory server. Another option is to have multiple directory servers, where users accept information only if it is signed by a threshold of the servers, similar to the directory servers in Tor.

The user then executes our distributed k -anonymity protocol with the relevant location brokers and one of the secure comparison servers, l . Our protocol uses the techniques of public-key cryptography, but we require the cryptosystem to have a special algebraic property: that it is *additive homomorphic*. Here, given only $\mathcal{E}_A(m_1)$ and $\mathcal{E}_A(m_2)$, where $\mathcal{E}_A(m)$ is an encryption of message m under public key A , one can efficiently compute $\mathcal{E}_A(m_1 + m_2)$. There are several cryptosystems with this property, such as the Paillier cryptosystem [21].

In our protocol, the user first asks each broker covering the query area for the number of users who have registered a cell in the query area as their current location with this particular broker. A broker gives this number to the user such that the user cannot learn it. Namely, if there are v_j users in the query area who have registered with broker j , broker j encrypts v_j with public key C_l of secure comparison server l , as published by the directory server, and sends $\mathcal{E}_{C_l}(v_j)$ to the user. Then, the user sums up the received numbers without being able to learn the sum. In particular, the user calculates $\mathcal{E}_{C_l}(r + \sum_i v_i)$ using the additive homomorphic property of the encryption scheme, where r is a random number generated by the user that will keep the total number of users hidden from secure comparison server l . Finally, with the help of secure comparison server l , the user determines whether this sum is at least k (see Section IV-B).

Our protocol gracefully deals with crashes of a location broker or of a secure comparison server. In the first case, the user contacts the remaining brokers, which might still report a sufficient number of registered users. To work around the second case, we can let the user and the broker choose a set of candidate secure comparison servers, instead of only a single one. Over time, the directory server will learn of the crash of a location broker or of a secure comparison server and will remove it from the directory.

B. Defence against Collusion

After computing the encrypted sum of users in her query area, the user, in cooperation with secure comparison server l , determines whether this sum is at least k . The user could simply send $\mathcal{E}_{C_l}(r + \sum_i v_i)$ and $r + k$ to server l , which would decrypt the first value, compare it to $r + k$, and inform the user of the result. Since both the sum and k are obscured with r , the server can learn neither of them. However, this solution is flawed, because it might reveal the total number of users to a secure comparison server and a location broker that are run by the same organization. Assume that the location broker is the only broker that covers the query area. Here, based on the knowledge of $\sum_i v_i$ (where the sum covers only one broker), the broker and the server can jointly determine r , which allows them to compute k . In turn, once they know a user's k , the server and the broker can infer the total number of registered

people in any query area chosen by the user, as long as the user's choice of k is static and the query area is covered by the broker. The coverage condition guarantees that the broker will be contacted by the user and hence can learn the query area. Otherwise, the server and the broker could learn only the total number of people, but not for which query area. To avoid these information leaks, we need to hide the user's input to the comparison, $r + k$, from the secure comparison server. Moreover, we need to hide the result of the comparison from the server, else the server could still infer $r + k$ in case it is found to be equal to $r + \sum_i v_i$.

To execute the comparison in this way, we exploit the Greater Than - Strong Conditional Oblivious Transfer (GT-SCOT) protocol [22]. The protocol has two participants, a receiver and a sender. The receiver and sender have private inputs x and y , respectively. The sender has two secrets, $s_0 \in D_S$ and $s_1 \in D_S$, where D_S is a subset of \mathbb{Z}_n . The sender wants to send s_0 to the receiver if $x < y$ and s_1 if $x > y$, but is oblivious about which secret is sent. In short, the sender cannot learn whether $x < y$ or $x > y$. The protocol requires a semantically secure additive homomorphic encryption scheme with large message domains, such as the Paillier scheme. In the protocol, the receiver encrypts x bit by bit with the receiver's public key and sends the vector of ciphertexts to the sender. The sender encrypts y bit by bit with the receiver's public key and finds the most significant bit that is different in the two numbers without learning its position. The sender then obliviously assigns s_0 or s_1 to that bit and randomizes all other bits. Then, the sender permutes the vector of encrypted values to prevent the receiver from learning the position of that bit and sends the vector to the receiver. The receiver decrypts the elements of the vector and stops when a value in D_S is found. For details of the GT-SCOT protocol, we refer to our technical report [20].

If s_0 and s_1 are already known to the receiver, the GT-SCOT protocol simply allows the receiver to learn whether $x < y$, $x = y$, or $x > y$. We exploit this observation in our distributed k -anonymity protocol. Here, the user sends only $\mathcal{E}_{C_l}(r + \sum_i v_i)$ to secure comparison server l . Then, the user and the server run the GT-SCOT protocol. The server uses $r + \sum_i v_i$ as the sender's input, y , and the user uses $r + k$ as the receiver's input, x . The GT-SCOT protocol guarantees that the server will not learn $r + k$ and the result of the comparison. However, it allows the user to distinguish between three cases ($\sum_i v_i < k$, $\sum_i v_i = k$, and $\sum_i v_i > k$), whereas k -anonymity does not distinguish between the equality and the greater-than case. Also, telling a user that there are precisely k people in the query area enables tracking attacks, similar to the one outlined in Section I. To avoid the equality case, we have the secure comparison server compute and compare bit-by-bit encryptions of $2 * (r + \sum_i v_i) + 1$ and $2 * (r + k)$.

C. Defence against Binary Search

A flaw of our protocol is that, using binary search, a user might still be able to learn the precise number of users in a cell. Namely, the user could present $\mathcal{E}_{C_l}(r + \sum_i v_i)$ multiple

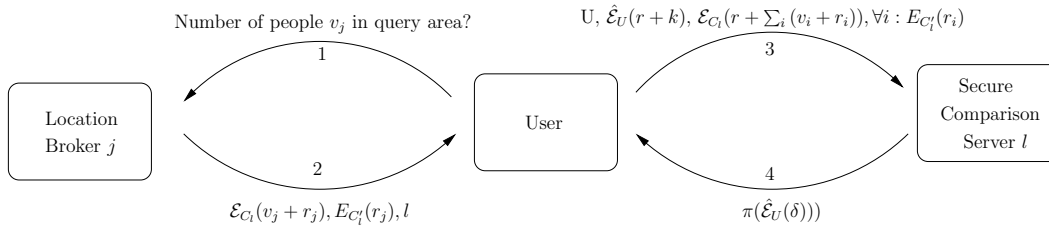


Fig. 2. Distributed k -anonymity protocol. U and C_l are the Paillier public key of the user and secure comparison server l , respectively. C_l' is the RSA public key of secure comparison server l . $\mathcal{E}_A(\cdot)$ denotes regular Paillier encryption with public key A . $\hat{\mathcal{E}}_A(\cdot)$ denotes *bit-by-bit* Paillier encryption with public key A . $E_A(\cdot)$ denotes RSA encryption with public key A . $\pi(\cdot)$ is a random permutation. δ is the result vector obliviously computed by secure comparison server l while executing the GT-SCOT protocol. The ciphertexts $E_{C_l}(r_i)$ also include an expiration date and are signed by location broker i (not shown).

times to the secure comparison server, maybe with a different value of r each time. By adjusting the value of k in each run of the GT-SCOT protocol, the user can perform a binary search for the actual value of $\sum_i v_i$.

To prevent this attack, we use expiring tickets. Instead of sending $\mathcal{E}_{C_l}(v_j)$ to the user, a location broker sends $\mathcal{E}_{C_l}(v_j + r_j)$ and a ticket that contains $E_{C_l'}(r_j)$, where r_j is a random number changing with each request and $E_{C_l'}(\cdot)$ is the RSA encryption function using RSA public key C_l' of the secure comparison server. A location broker also includes an expiration date in the ticket and signs the ticket. A secure comparison server will decrypt all $E_{C_l'}(r_j)$ and subtract $\sum_i r_i$ from $r + \sum_i (v_i + r_i)$. The server also remembers tickets till their expiration date and refuses to re-use a ticket seen previously. This way, the attack mentioned above will fail, even if r is changed. Also, the user cannot use fresh tickets with a previously presented encrypted sum, since the r_i value will be different, meaning the secure comparison server cannot compute the correct input value for the GT-SCOT protocol and the user cannot learn any useful information from this operation. Similar to traditional k -anonymity approaches based on a central trusted server, a location broker can limit the query frequency of users.

Figure 2 illustrates our distributed k -anonymity protocol based on the GT-SCOT protocol and expiring tickets.

D. Choice of Secure Comparison Server

Having multiple secure comparison servers distributes load and avoids a single point of failure. Our distributed k -anonymity protocol requires that we pick one of the secure comparison servers at the beginning of a protocol run. Because of possible collusion between a user and a secure comparison server, which would allow the user to learn the total number of registered users in her query area, we cannot let a user choose a secure comparison server. Instead, we need to choose a server such that, over time, the risk of a user working together with a colluding server is limited by k/n , where n corresponds to the number of secure comparison servers, with k of them colluding with the user. Therefore, we cannot statically assign a secure comparison server to a user, since we might be unlucky and pick a colluding one. Moreover, a user might decide not to trust servers maintained by particular organizations, and she might refrain from using our system if we forced her to use

such a server all the time.

Another strategy is to have each location broker randomly choose a secure comparison server for a query. However, this strategy has two flaws: First, our protocol requires that all brokers choose the same server, which will likely not be the case here. Second, if a user is assigned to a non-colluding server, she can repeat her query until a colluding server is chosen. To address these flaws, we need an assignment scheme that, within a particular time frame, has all brokers assign the same server to a particular user. The length of the time frame should be such that the impact of using a malicious server within the entire duration of the time frame is limited (e.g., the time frame should be shorter than a day) and such that if a user decides to perform an attack at a particular moment in time, her expected waiting time till she is being assigned a colluding server is so long that the attack environment (e.g., locations of users) will likely have significantly changed by then (e.g., the time frame should be longer than a minute).

We now present our algorithm for choosing a secure comparison server. There is a sign-up server that randomly assigns each user to one out of n groups when she signs up to our system, where n corresponds to the number of secure comparison servers. (The assignment can expire to deal with new secure comparison servers.) The sign-up server could be identical to the directory server. To determine the secure comparison server to be used for encrypting the response to a user's query, each location broker executes a well-known, deterministic function that maps the identity of the user's group to a secure comparison server. (Directly mapping the user's identity, instead of the identity of the user's group, is not an option because letting a broker know of a user's identity would make the user trackable and identifiable by a location broker.) This mapping function is bijective and depends on the current time. For example, the function can be defined as $f = (ep + id) \bmod n$. Here, time is split into epochs, with ep indicating the current epoch ($ep \geq 0$). Our suggested duration of an epoch is one hour. id is the identity of a user's group ($0 \leq id < n$), as reported by the user when sending her query to the location broker. We can also design more complex mapping functions, such as a function that changes the order in which a group is mapped to a server over time.

To ensure that the group identity reported by a user is accurate, we take advantage of cryptographic group signatures [23].

In a group signature scheme, any member of a group can sign a message without a signature verifier being able to infer which member generated the signature. Only the group manager, that is, the entity generating and distributing signing keys to group members, can trace a signature to its issuer. In our scheme, the group manager corresponds to the sign-up server. When querying a location broker, a user creates a group signature to prove membership in her group to the broker. This allows the user to remain anonymous within her group while not being able to claim membership in other groups.

In general, to minimize the risk of collusion, we do not let random people deploy a secure comparison server. Instead, the directory server should vet a server before listing it, similar to the limited vetting done by a directory server in Tor.

V. LOCATION REGISTRATION

As mentioned in Section III-B, consistent with earlier work, our threat model assumes that a location broker can detect attempts by a user to register with the broker multiple times in parallel. Having multiple location brokers, as it is the case in our solution, introduces a new vulnerability. Namely, a user could register multiple times, but each time with a *different* location broker. This way, other users might be wrongly told that their k -anonymity preference is satisfied. There are both technical and non-technical controls for this vulnerability. Charging money is an example of a non-technical control. Namely, if location brokers are maintained by operators of cellphone networks as a service to customers, a user would have to buy multiple cellphones and plans to register in parallel with multiple brokers, which makes the attack expensive. In the remainder, we present a technical control that does not make any assumptions about the underlying communication technology. Since we control the vulnerability, our threat model for user behaviour can remain identical to the threat model in the earlier work.

There are two naïve approaches to prevent a user from registering with different location brokers concurrently. In the first one, a location broker contacts the other location brokers whenever a user registers and inquires whether the user has already registered with one of them. This approach raises privacy concerns and is expensive in terms of performance. The second approach has each broker keep records of its registered users. Periodically, the brokers compare records and try to detect misbehaving users. The main problems of this approach are the privacy concerns raised by the record keeping and by the comparison and that it is retroactive.

Our solution gets around these problems. It is based on e-cash [24] and is outlined in Figure 3. E-cash allows a player to withdraw a coin from the bank and to spend it with a second player. The second player deposits the received coin with the bank. In our solution, a user gets one (and only one) coin from the bank. The role of the bank can be assumed by the directory server. When a user registers with a location broker, she spends her coin at the broker. Since the user has only one coin, registering with another broker amounts to double spending of the coin. The other broker detects this

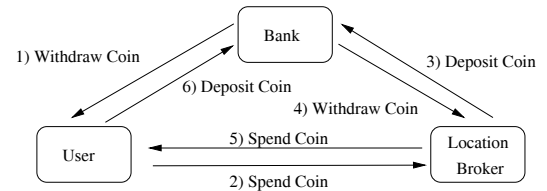


Fig. 3. Defence against multiple registrations based on e-cash. By being given only one coin, a user can register only once. The user is returned her coin when de-registering.

double spending when depositing the coin, either immediately in case of an online e-cash scheme [24] or in a delayed way in case of an offline scheme [25]. In the former case, the broker will deny registration. In the later case, the bank will learn the user’s identity and will ban her from the system (see Section VI-C). Whenever a location broker deposits a user’s (valid) coin, the broker also withdraws a fresh coin from the bank. (Alternatively, since location brokers are not malicious, the bank can periodically give a set of coins to the broker for increased performance.) When the user wants to de-register, she asks the broker to spend this coin by giving it to the user. The user then deposits the coin and withdraws a fresh coin from the bank, which she can later spend at another broker.

The benefits of our solution are that it does not require all location brokers to be contacted for a registration and that location brokers do not need to keep records of registered users after their de-registration. Moreover, since e-cash is anonymous, the bank cannot learn which user registers with which location broker, and a location broker cannot learn a user’s identity and where a user has registered previously.

In case a location broker crashes, a registered user will not be able to register with a new broker. Here, we have the user contact the bank with a proof of registration issued by the broker. The bank will then issue a new coin to the user. When the broker comes back up, it will re-synchronize with the bank.

VI. SECURITY ANALYSIS

In this section, we first review how our protocol defends against the threats listed in Section III-B, where we assume that location brokers and secure comparison servers are honest-but-curious. In the remainder, we discuss how our architecture can be extended to defend against malicious parties.

A. Threat Analysis

In our protocol, location brokers do not interact with each other, so they cannot learn the location of users in the query area who are registered with other brokers, not even their total number. Our architecture is based on e-cash and group signatures, which allows users to remain anonymous to a location broker. A user should not directly connect to a broker during a registration or query operation. Information about this connection (e.g., the user’s IP address) might allow the broker to re-identify the user and to learn her location from her query area. Instead, the user should communicate through a trusted proxy or an anonymous communication channel (e.g., Tor). In

practice, the location broker that a user registers with might already know the user’s identity and might forbid anonymous registrations. (E.g., an operator of a cellphone network often knows the identity of its customers.) In this setup, the location broker that a user registers with can serve as the trusted proxy for contacting the other brokers during a query operation. Despite the communication between these brokers and a user being proxied, the brokers might be able to re-identify the user based on her query area if the area is associated with the user or under physical surveillance by a broker. To address this threat, the user should query the brokers only if she knows that there are at least $k - 1$ other people in the query area, that is, we end up with a chicken-egg problem. Therefore, as stated in Section III-A, we choose the width of a cell in our architecture such that there is a realistic chance that multiple users can be located in the cell, which makes this attack hard.

A secure comparison server learns no useful information during a comparison operation, not even its outcome. A server also gains no benefit from colluding with other secure comparison servers. To remain anonymous to a location broker, a user should not directly connect to a secure comparison server, since the tickets issued by the broker allow a secure comparison server that is run by the same organization to link a query and a comparison operation.

A user learns only whether the total number of users in her query area is at least k . The expiring tickets prevent her from learning the actual number of users with a binary search.

The directory server cannot learn any location information, because users do not retrieve individual records for their current cell from the server. The published directory is signed, which prevents the directory server from misbehaving.

B. Malicious Servers or Brokers

Assume that a malicious secure comparison server fails to correctly execute some of the steps in the GT-SCOT protocol. While it is not possible for the server to learn the total number of users, due to the randomness added by the user, the server could misbehave with the intent to give the user the impression that there are at least k registered users in an area, even if this is not the case, or vice versa. We could address this concern by adding zero-knowledge proofs to each step of the protocol, proving that the step was executed faithfully. For example, Groth [26] proposes an efficient scheme for proving in zero-knowledge the correctness of a permutation of homomorphic encryptions. As it turns out, this scheme requires three additional rounds of interactions between the prover and the verifier, which makes it expensive for mobile devices. Therefore, in our scheme, we choose a retroactive approach. We have a secure comparison server log the random values used in its encryption and permutation operations. Furthermore, the server has to sign all its generated messages to achieve non-repudiation. If users suspect misbehaviour, they, likely in collaboration with the directory server, can force the secure comparison server to reveal the logged values and its private key and can validate the server’s computations.

Similar to the secure comparison server, a malicious location broker can misbehave while executing our protocol. In particular, a broker can encrypt a value that is different from the actual number of users registered in an area. It is possible to ask a broker to keep a record of all its actions. However, this record would have to include location registrations of users, which is problematic in terms of privacy. We prefer a less invasive approach. If users suspect misbehaviour by a location broker (and misbehaviour by a secure comparison server can be excluded, based on the above mechanism), they report the set of location brokers from which they retrieved information to the directory server. Over time, this will allow the directory server to single out a particular location broker.

C. Malicious Users

Malicious users could report wrong locations to a location broker during registration. As it turns out, a complete defence against this attack is likely impossible. A determined attacker can give her mobile device to another user or simply tamper with the location reporting mechanism on her mobile device. A user could also acquire multiple devices, maybe under different identities, and use them to register multiple times. As mentioned in our threat model (see Section III-B), these threats are not new to our system; they also arise in previous schemes. Let us outline some mechanisms that make these attacks harder.

A location broker might be able to detect wrongly reported locations. For example, if a broker is controlled by the operator of a WiFi network, the operator can ensure that a reported location is close to the WiFi access point from which the registration request was sent. An operator of a cellphone network can verify whether the reporting device is close to a particular cellphone tower.

We can also exploit the sign-up process, as introduced in Section V, to defend against malicious users. The sign-up server can require physical identification, which reduces the danger of a user signing up multiple times. However, this approach makes the system more difficult to use. An alternative is to ask the user for a credit card number, including her name and billing address. This option becomes especially attractive if the system charges its users in the first place. Billing for the usage of our system itself can become a mechanism for reducing misbehaviour, because an attacker might not have the necessary resources for a large-scale attack.

VII. EVALUATION

In this section, we evaluate our distributed k -anonymity protocol. We first examine the cost of contacting a location broker, followed by the cost of contacting a secure comparison server. In our evaluation, we focus on the cost of the homomorphic encryption operations and of the GT-SCOT protocol. To the best of our knowledge, no measurement-based evaluation of this protocol has been published, whereas there are such evaluations of the other two cryptographic protocols used in our architecture. For example, Belenkiy et al. [27] evaluate e-cash and Cornelius et al. [28] evaluate group signatures.

We implemented our protocol using the OpenSSL and NTL [29] libraries. The key size for RSA and Paillier is 1024 bits. We deploy a location broker and a secure comparison server on a 2.4 GHz Intel Xeon Dual Core running Linux 2.6.24. The user has a slow laptop (a ThinkPad T43 with a 2 GHz Intel Pentium M running Linux 2.6.22) to approximate the capabilities of a modern smartphone. Communication runs over WiFi and is protected against eavesdroppers with TLS using AES128 in CBC mode with an ephemeral Diffie-Hellman key exchange for forward secrecy.

A. Location Broker

We examine the performance of querying a location broker for the number of people in the query area and of adding this number to an existing encrypted sum. In the experiment, when a user connects to a location broker, the location broker sends back a Paillier encrypted random value. The user then performs a homomorphic addition. We repeat the experiment ten times and report mean and standard deviation.

The overall delay experienced by the user is 39.9 ± 0.7 ms. It takes 32.5 ± 0.7 ms to set up a TLS connection, which includes client and server authentication. The server takes 7.4 ± 0.0 ms to Paillier encrypt a random value. The cost of the homomorphic addition operation by the user is negligible. In summary, setting up the TLS connection is about four times as expensive as the Paillier encryption operation. As mentioned in Section VI-A, to hide her identity, a user might not directly connect to a location broker. Here, the cost of encryption in relation to the cost of connection setup becomes even smaller.

In practice, the user will likely contact multiple location brokers. Apart from the addition operation, whose cost is negligible, the brokers can be contacted in parallel. If this is not feasible for the user's device, the overall delay will be linear in the number of location brokers. We envision that this number is small (5-10 brokers) in most scenarios. This number reflects the number of cellphone and WiFi network operators providing coverage for the query area, which tends to be small. In addition, there might be a small number of independently operated location brokers.

The user also needs to Paillier encrypt the random value that she will add to the encrypted sum of users reported by the location brokers. This encryption takes 67.2 ± 0.5 ms. As expected, the encryption operation is slower on the laptop than on the server. However, as opposed to the other operations, this encryption can occur offline. Moreover, the user can use an encrypted value multiple times for a secure comparison server.

B. Secure Comparison Server

We evaluate the performance of the GT-SCOT protocol for different bit lengths of the bit-by-bit homomorphic encryption operation. We vary the bit length between 4 and 16 and perform fifty runs for each configuration.

We present our results in Figure 4. The bottom graph shows the cost of the sender side of the GT-SCOT protocol, which varies between 66.3 ± 0.3 ms for a bit length of 4 and 247.9 ± 0.9 ms for a bit length of 16. The middle graph corresponds to

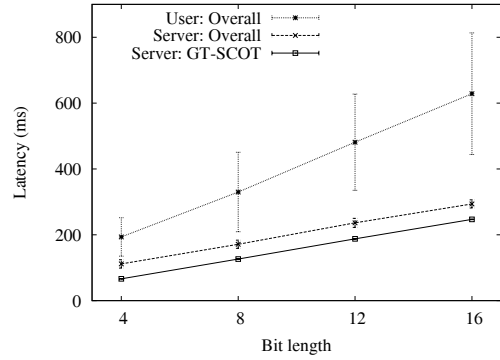


Fig. 4. Latency experienced by the secure comparison server and the user in relationship to the bit length used in the GT-SCOT protocol. We show mean and standard deviation.

the overall cost by the server. In addition to the sender side of the GT-SCOT protocol, it also includes the cost of setting up a TLS connection and Paillier decryption of the total number of users. Finally, the top graph shows the overall latency, as experienced by the user. It varies between 193.5 ± 58.4 ms for a bit length of 4 and 628.6 ± 184.7 ms for a bit length of 16. The overall latency corresponds to the overall cost by the server plus the cost of the receiver side of the GT-SCOT protocol, which takes 77.5 ± 47.7 ms for a bit length of 4 and 330.4 ± 179.6 ms for a bit length of 16. The standard deviation is large, because the user decrypts the permuted result vector received from the secure comparison server element by element and stops as soon as she finds the server's answer.

In our implementation, we let a user choose the bit length. In practice, we expect that bit lengths between 8 and 12 will be used mainly, depending on the number of location brokers covering the query area and the maximum number of reported users for the query area (which is different from the number of registered users in the query area). If there are a bits in total, we can support up to 2^c location brokers and up to $2^b - 1$ reported users per location broker per query area, where $a = b + c + 2$. (This implies $0 \leq k \leq 2^c * (2^b - 1)$.) A user informs a broker of her choice of b ; if there are more than $2^b - 1$ registered users in the query area, the broker simply reports $2^b - 1$ users. (As stated in Section IV-A, each broker adds a random value to its reported number. We can ignore these values here, since they are subtracted by the secure comparison server before running the GT-SCOT protocol. This will also revert any wrap-arounds that might have occurred due to choosing a large random value.) The two remaining bits leave space for adding the random number r ($0 \leq r \leq 2^c * (2^b - 1)$)¹ chosen by the user to the total number of users, which requires at most $b + c + 1$ bits, and for allowing the secure comparison server to double the resulting sum to avoid the equality case. For example, for a bit length of 8, we can support up to 8

¹If $r = 0$ and $\forall i : v_i = 0$, the secure comparison server can infer these values. Similar for the maximum case. To lower the probability for this scenario, the user can dynamically increase (decrease) her lower (upper) bound for r and keep her choice secret.

location brokers and 7 reported users per broker per query area. Here, the overall latency experienced by the user is 330.0 ± 120.7 ms. For a bit length of 12, we can support up to 16 location brokers and 63 reported users per broker per query area. Here, the overall latency experienced by the user is 481.2 ± 146.2 ms. In short, we expect the overall latency to be noticeable, but tolerable.

The user also needs to perform a bit-by-bit encryption of the sum of her privacy preference and of her chosen random value. The cost of this operation varies between 210.4 ± 0.8 ms for a bit length of 3 and 864.4 ± 115.0 ms for a bit length of 15. (The large variation is an artifact of using a bit length of 15. The variation is small for a bit length of 16, which has a larger mean. We are investigating this behaviour.) However, as opposed to the other operations, this encryption can be done offline. Moreover, the user can use an encrypted value multiple times for a secure comparison server.

VIII. CONCLUSIONS AND FUTURE WORK

We have presented a protocol for location privacy based on k -anonymity that needs neither a single trusted server nor users to trust each other. Our sample implementation and its evaluation have shown that the protocol is efficient.

In terms of future work, we are integrating our protocol into a platform for location-based services, which will allow us to gather more insights about the protocol's usability in practice.

ACKNOWLEDGEMENTS

We thank the anonymous reviewers for their comments. This work is supported by the Natural Sciences and Engineering Research Council of Canada.

REFERENCES

- [1] M. Gruteser and D. Grunwald, "Anonymous Usage of Location-Based Services Through Spatial and Temporal Cloaking," in *Proceedings of 1st International Conference on Mobile Systems, Applications, and Services (MobiSys 2003)*, May 2003, pp. 31–42.
- [2] B. Bamba, L. Liu, P. Pesti, and T. Wang, "Supporting Anonymous Location Queries in Mobile Environments with PrivacyGrid," in *Proceedings of 17th International World Wide Web Conference (WWW2008)*, April 2008, pp. 237–248.
- [3] A. R. Beresford, "Location privacy in ubiquitous computing," Computer Laboratory, University of Cambridge, Tech. Rep. 612, January 2005.
- [4] C. Bettini, S. Mascetti, X. S. Wang, and S. Jajodia, "Anonymity in Location-based Services: Towards a General Framework," in *Proceedings of 8th International Conference on Mobile Data Management (MDM 2007)*, May 2007, pp. 67–79.
- [5] C.-Y. Chow, M. F. Mokbel, and X. Liu, "A Peer-to-Peer Spatial Cloaking Algorithm for Anonymous Location-based Services," in *Proceedings of 14th ACM International Symposium on Advances in Geographic Information Systems (ACM-GIS'06)*, November 2006, pp. 171–178.
- [6] B. Gedik and L. Liu, "Location Privacy in Mobile Systems: A Personalized Anonymization Model," in *Proceedings of 25th International Conference on Distributed Computing Systems (ICDCS 2005)*, June 2005, pp. 620–629.
- [7] G. Ghinita, P. Kalnis, and S. Skiadopoulos, "PRIVÉ: Anonymous Location-Based Queries in Distributed Mobile Systems," in *Proceedings of 16th International World Wide Web Conference (WWW2007)*, May 2007, pp. 371–380.
- [8] —, "MobiHide: A Mobile Peer-to-Peer System for Anonymous Location-Based Queries," in *Proceedings of Proceedings of 10th International Symposium on Spatial and Temporal Databases (SSTD 2007)*, July 2007, pp. 221–238.
- [9] P. Kalnis, G. Ghinita, K. Mouratidis, and D. Papadias, "Preserving Anonymity in Location Based Services," School of Computing, The National University of Singapore, Tech. Rep. TRB6/06, 2006.
- [10] A. Kapadia, N. Triandopoulos, C. Cornelius, D. Peebles, and D. Kotz, "AnonySense: Opportunistic and Privacy-Preserving Context Collection," in *Proceedings of 6th International Conference on Pervasive Computing (Pervasive 2008)*, May 2008, pp. 280–297.
- [11] S. Mascetti and C. Bettini, "A Comparison of Spatial Generalization Algorithms for LBS Privacy Preservation," in *Proceedings of International Workshop on Privacy-Aware Location-based Mobile Services (PALMS)*, May 2007.
- [12] M. F. Mokbel, C.-Y. Chow, and W. G. Aref, "The New Casper: Query Processing for Location Services without Compromising Privacy," in *Proceedings of 32nd International Conference on Very Large Data Bases (VLDB 2006)*, September 2006, pp. 763–774.
- [13] H. Shin, V. Atluri, and J. Vaidya, "A Profile Anonymization Model for Privacy in a Personalized Location Based Service Environment," in *Proceedings of 9th International Conference on Mobile Data Management (MDM 2008)*, April 2008, pp. 73–80.
- [14] G. Zhong and U. Hengartner, "Toward a Distributed k -Anonymity Protocol for Location Privacy," in *Proceedings of 7th Workshop on Privacy in the Electronic Society (WPES 2008)*, October 2008, pp. 33–37.
- [15] P. Samarati and L. Sweeney, "Protecting Privacy when Disclosing Information: k -Anonymity and Its Enforcement through Generalization and Suppression," SRI International, Tech. Rep. SRI-CSL-98-04, 1998.
- [16] S. Zhong, Z. Yang, and R. N. Wright, "Privacy-Enhancing k -Anonymization of Customer Data," in *Proceedings of 24th Symposium on Principles of Database Systems (PODS 2005)*, June 2005, pp. 139–147.
- [17] A. R. Beresford and F. Stajano, "Location Privacy in Pervasive Computing," *IEEE Pervasive Computing*, vol. 2, no. 1, pp. 46–55, 2003.
- [18] T. Jiang, H. J. Wang, and Y.-C. Hu, "Preserving Location Privacy in Wireless LANs," in *Proceedings of 5th International Conference on Mobile Systems, Applications, and Services (MobiSys 2007)*, June 2007, pp. 246–257.
- [19] R. Dingedine, N. Mathewson, and P. Syverson, "Tor: The Second-Generation Onion Router," in *Proceedings of 13th USENIX Security Symposium*, August 2004, pp. 303–319.
- [20] G. Zhong and U. Hengartner, "A Distributed k -Anonymity Protocol for Location Privacy," Centre for Applied Cryptographic Research, University of Waterloo, Tech. Rep. CACR 2008-17, September 2008.
- [21] P. Paillier, "Public-Key Cryptosystems Based on Composite Degree Residuosity Classes," in *Proceedings of EUROCRYPT '99*, May 1999, pp. 223–238.
- [22] I. F. Blake and V. Kolesnikov, "Strong Conditional Oblivious Transfer and Computing on Intervals," in *Proceedings of ASIACRYPT 2004*, December 2004, pp. 515–529.
- [23] D. Chaum and E. van Heyst, "Group Signatures," in *Proceedings of EUROCRYPT '91*, April 1991, pp. 257–265.
- [24] D. Chaum, "Blind Signatures for Untraceable Payments," in *Proceedings of CRYPTO '82*, August 1982, pp. 199–203.
- [25] J. Camenisch, S. Hohenberger, and A. Lysyanskaya, "Compact E-Cash," in *Proceedings of EUROCRYPT 2005*, May 2005, pp. 302–321.
- [26] J. Groth, "A Verifiable Secret Shuffle of Homomorphic Encryptions," in *Proceedings of 6th International Workshop on Practice and Theory in Public Key Cryptography*, January 2003, pp. 145–160.
- [27] M. Belenkiy, M. Chase, C. C. Erway, J. Janotti, A. Küpçü, A. Lysyanskaya, and E. Rachlin, "Making P2P Accountable without Losing Privacy," in *Proceedings of 6th Workshop on Privacy in the Electronic Society (WPES 2007)*, October 2007, pp. 31–40.
- [28] C. Cornelius, A. Kapadia, D. Kotz, D. Peebles, M. Shin, and N. Triandopoulos, "AnonySense: Privacy-Aware People-Centric Sensing," in *Proceedings of 6th International Conference on Mobile Systems, Applications, and Services (MobiSys 2008)*, June 2008, pp. 211–224.
- [29] V. Shoup, "NTL: A Library for doing Number Theory," <http://www.shoup.net/ntl>, accessed December 2008.