

# $\epsilon$ -PPI: Searching Identity in Information Networks with Quantitative Privacy Guarantees

Yuzhe Tang <sup>†</sup>    Ling Liu <sup>†</sup>    Arun Iyengar <sup>‡</sup>    Kisung Lee <sup>†</sup>    Qi Zhang <sup>†</sup>

<sup>†</sup>Georgia Institute of Technology, GA, USA Email: {yztang, lingliu@cc., kisung.lee@cc., qzhang90}@gatech.edu

<sup>‡</sup>IBM T.J. Watson Research Center, NY, USA Email: aruni@us.ibm.com

## Abstract

In information sharing networks, having a privacy preserving index (or PPI) is critically important for providing efficient search on access controlled content across distributed providers while preserving privacy. An understudied problem for PPI techniques is how to provide controllable privacy preservation, given the innate difference of privacy of the different content and providers. In this paper we present a configurable privacy preserving index, coined  $\epsilon$ -PPI, which allows for quantitative privacy protection levels on fine-grained data units. We devise a new common-identity attack that breaks existing PPI's and propose an identity-mixing protocol against the attack in  $\epsilon$ -PPI. The proposed  $\epsilon$ -PPI construction protocol is the first without any trusted third party and/or trust relationship between providers. We have implemented our  $\epsilon$ -PPI construction protocol by using generic MPC techniques (secure multi-party computation) and optimized the performance to a practical level by minimizing the costly MPC computation part.

## I. Introduction

In information networks, autonomous service providers store private personal records on behalf of individual owners and enable information sharing under strict enforcement of access control rules. Such information networks have the following salient features: 1) providers, each under a different administrative domain, do not mutually trust each other, 2) providers have the responsibility of protecting owners' privacy and 3) it is crucial to share information between providers from an application perspective.

An example of the information network is the emerging Nationwide eHealthcare Information Network (NHIN [1] and Healthcare software CONNECT [2]), in which patients delegate their personal medical records to the hospitals where they have visited. Different hospitals may compete for customer patients and have conflicting economic interests, which renders it hard to build full trust relationships between them. Hospitals are

responsible for protecting patient privacy, which is regulated by Federal laws (e.g. HiPPA [3]). On the other hand, to provide immediate and accurate medical diagnosis and treatment, it is important to have an information sharing service; for example, when a patient who is unconscious is sent to a hospital, such information sharing can help the doctor be able to retrieve the patients' medical history that involves multiple (untrusted) hospitals. Another example is cross-university online course management systems (e.g. StudIP [4], Swiki [5]). Such online systems allow sharing of access-controlled materials within groups of students and teachers; while the information sharing is crucial for collaborative learning and improved learning efficiency, it may pose a threat to student privacy; on the other hand, protection of student privacy is required by Federal laws (e.g. FERPA [6]). In these untrusted networks with needs of cross-domain collaborations, it calls for a global mechanism to protect privacy of a patient or a student, while enabling effective information sharing.

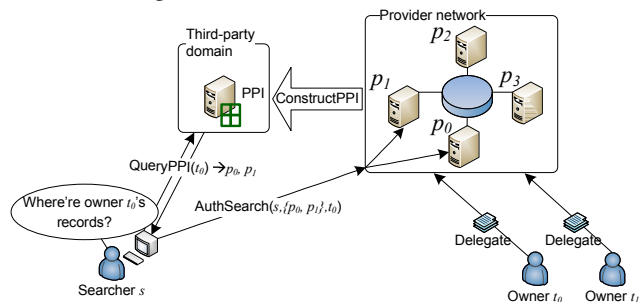


Fig. 1: The system of PPI and the provider network

To support and promote information sharing among mutually untrusted providers, privacy preserving index or PPI [7], [8], [9], [10] is proposed. The PPI aims at supporting a global search facility hosted on a third-party entity; the design of PPI should respect the providers' complete access control on personal records and protect their privacy. The typical working of a PPI, as will be elaborated in Section II, is a two-phase search procedure. As in Figure 1, a searcher, in the hope of finding certain owners' records, first queries the PPI, and obtains a list of providers that may or may not have the records of interest. Then for each provider in the list, the searcher attempts

to get authenticated and authorized before she can locally search the private records there. A PPI is usually hosted on a third-party and untrusted entity, mainly because of the difficulty to find a global entity trusted by all (mutually untrusted) providers; for example, consider the U.S. government as a candidate; various scandals including the recent PRISM program [11] have made the government lose the public trust. Hence, it is desirable to construct the PPI in such a way that the data owners’ privacy is protected.

## Quantitatively Differentiated Privacy Preservation

While existing PPI’s have addressed privacy preservation, none of these approaches recognize the needs of differentiating privacy preservation for different owners and providers. To be specific, the privacy goals addressed by a PPI system are about a private fact whether “an owner  $t_j$  has his/her record stored on provider  $p_i$ ”. It is evident that disclosing the private fact regarding different owners and providers causes different levels of privacy concerns. For example, a woman may consider her visit to a women’s health center (e.g., for an abortion) much more sensitive than her visit to a general hospital (e.g., for cough treatment). Similarly, different owners may have different levels of concerns regarding their privacy: while an average person may not care too much about their visit to a hospital, a celebrity may be more concerned about it, because even a small private matter of a celebrity can be publicized by the media (e.g., by paparazzi). To address the innate privacy difference of owners, it is therefore critical to differentiate privacy protection to address the innate different privacy concerns in a PPI system. That being said, using existing PPI approaches can not provide quantitative guarantees on the privacy preservation degree, let alone on a fine-grained per-owner basis. The cause, largely due to the degree-agnostic way of constructing PPI systems, is analyzed in Appendix B.

In this paper, we propose a new PPI abstraction for differentiated and quantitative privacy control, coined  $\epsilon$ -PPI. Here,  $\epsilon$  is a privacy aware knob that allows each owner to mark a desired privacy level when delegating data to the providers. Specifically,  $\epsilon_j$  is a value in a spectrum from 0 to 1, where value 0 is for the least privacy concern (in this case, the PPI returns the list of exactly those “true positive” providers who truly have the records of interest) and value 1 for the best privacy preservation (in this case, PPI returns all providers, and a search essentially goes to the whole network). By this means, an attacker observing the PPI search result can only have a bounded confidence by  $\epsilon$  in successfully identifying a true positive (and thus vulnerable) provider from the obscured provider list.

*Challenges:* To construct the new  $\epsilon$ -PPI abstraction, it poses challenges. On one hand, achieving quantitative privacy guarantee typically requires the index construction to know more about providers’ data (e.g. owner identity distribution across the provider network), which entails information exchange between providers. On the other hand, providers do not trust each other and may feel reluctant to disclose too detailed information. Therefore, it is essential to draw a clear

line between what is private information and what is not during the index construction, and to fully utilize the non-private information to provide as much quantitative guarantee as possible. In our proposed construction protocol, we utilize the owner frequency (i.e. the number of providers an owner has delegated her records to). Our unique insight in protocol design is that the owner frequency is private only when the value is big (i.e., when the owner’s record appears almost everywhere). This is because knowing such “common” owners would give an adversary confidence to make successful guesses regarding any provider. By only releasing the small-value frequency, we can protect providers’ privacy and deliver a quantitative guarantee.

In realizing the design protocol in a mutually untrusted network, we rely on MPC computation (i.e., secure multi-party computation [12], [13], [14], [15]) which addresses the input privacy for generic computation. However, it raises performance issues when directly applying MPC techniques in our problem setting. On one hand, current MPC computation platforms can only scale to small workloads [16]; they are practically efficient only for simple computation among few parties. On the other hand, a typical PPI construction may involve thousands of owners and tens or hundreds of providers, which entails an intensive use of bit-wise MPC computation. It is therefore critical to a practical MPC protocol to efficiently carry out the computation for  $\epsilon$ -PPI construction. In this regards, we propose to minimize the expensive MPC computation by using a parallel secure sum protocol. The secure sum can be efficiently carried out by a proposed secret sharing scheme with additive homomorphism. Based on the proposed MPC primitive, our index construction protocol protects providers’ privacy and can tolerate collusion of up to  $c$  providers ( $c$  is configurable).

The contributions of this paper can be summarized as follows,

- We propose  $\epsilon$ -PPI that differentiates the needs of privacy protection in a quantitative manner. The  $\epsilon$ -PPI exposes a new delegate operation to owners, which allows them to specify their different levels of privacy concerns. This new privacy knob, coined  $\epsilon$ , can give quantitative privacy control while enabling information sharing.
- We propose  $\epsilon$ -PPI construction protocol for an untrusted environment. As far as we know, this is the first PPI construction protocol without assumption on trusted parties or mutual trust relationships between providers. The performance of  $\epsilon$ -PPI construction protocol is extensively optimized by reducing the use of costly generic MPC and using the proposed domain-specific protocols. The proposed construction protocol is implemented and evaluated with verified performance superiority.
- We introduce a new privacy attack (called common-owner attack) that can break generic PPI systems. The new attack model targets vulnerable common owners. Our proposed  $\epsilon$ -PPI is the first to resist common-owner attacks by using a proposed term-mixing protocol.

The rest of this paper proceeds as follows: Section II formulates the  $\epsilon$ -PPI problem. Section III and IV respectively describe the computation model and distributed implementation of the  $\epsilon$ -PPI construction protocol. Section V presents evaluation results,

and Section VI surveys the related work before the conclusion in Section VII.

## II. Problem Formulation

### A. System Model

We formally describe our system model, which involves four entities: 1) a set of  $n$  data owners, each of whom, identified by  $t_j$ , holds a set of personal records, 2) a provider network consisting of  $m$  providers in which a provider  $p_i$  is an autonomously operating entity (e.g. a hospital or a university), 3) a global PPI server in a third-party domain, 4) a data searcher who wants to find all the records of an owner of interest. The interactions between these four entities are formulated in the following four operations.

- **Delegate**( $\langle t_j, \epsilon_j \rangle, p_i$ ): A data owner  $t_j$  can delegate her records to provider  $p_i$  based on her trust relationship (e.g. such trust can be built based on her visit to a hospital). Along with the record delegation, the owner can specify her preferred privacy degree  $\epsilon_j$ . Here  $\epsilon_j$  indicates the level of privacy concerns, ranging from 0 up to 1. For example, a VIP user (e.g. a celebrity patient in the eHealthcare network) may want to set the privacy level at a high value while an average patient may set the privacy level at a medium value <sup>1</sup>.
- **ConstructPPI**( $\{\epsilon_j\}$ ): After data records are populated, all  $m$  providers in the network join a procedure **ConstructPPI** to collectively construct the privacy preserving index. The index construction should comply with owner-specified privacy degree  $\{\epsilon_j\}$ . As will be elaborated, the constructed PPI contains noises or false positives for the purpose of privacy preservation and  $\{\epsilon_j\}$  is materialized as the false positive rate of owner  $t_j$ .
- **QueryPPI**( $t_j$ )  $\rightarrow$   $\{p_i\}$ : At data serving time, a searcher  $s$ , in the hope of finding owner  $t_j$ 's records, initiates a two-phase search procedure consisting of two operations, **QueryPPI**( $t_j$ )  $\rightarrow$   $\{p_i\}$  and **AuthSearch**( $s, \{p_i\}, t_j$ ). This is illustrated in Figure 1. For the first phase, the searcher poses query request, **QueryPPI**( $t_j$ ), and the PPI server returns a list of providers  $\{p_i\}$  who may or may not have records of the requested owner  $t_j$ . The query evaluation in PPI server is trivially done since the PPI, once constructed, contains the (obscured) mapping between providers and owners.
- **AuthSearch**( $s, \{p_i\}, t_j$ ): The second phase in the search is for searcher  $s$  to contact each provider in list  $\{p_i\}$  (i.e. the result list from the first phase) and to find owner  $t_j$ 's records there. This process involves user authentication and authorization regarding searcher  $s$ ; we assume each provider has already set up its local access control subsystem for authorized access to the private personal records. Only after authorization can the searcher search the local repository on provider  $p_i$ .

<sup>1</sup>To prevent every user from setting the highest value of  $\epsilon$ , one possible way is to differentiate prices for different privacy settings. The system has incentive to do so, since high privacy settings incur more overhead in the provider networks.

We describe the internal data model in a PPI. Each personal record contains an owner identity  $t_j$ <sup>2</sup> (e.g. the person's name). As shown in Figure 2, a provider  $p_i$  summarizes its local record repository by a membership vector  $M_i(\cdot)$ ; it indicates the list of owners who have delegated their records on provider  $p_i$ . For example, provider  $p_0$  who has records of owner  $t_0$  and  $t_1$  maintains a membership vector as  $M_i = \{t_0 : 1, t_1 : 1, t_2 : 0\}$ . In our model, the same owner can have records spread across multiple providers (e.g., a patient can visit multiple hospitals). The constructed PPI maintains a mapping between providers and owners; it is essentially a combination of all provider-wise membership data, yet with noises. The PPI mapping data is an  $m \times n$  matrix  $M'(\cdot, \cdot)$ , in which each row is of an owner, each column of a provider and each cell of a Boolean value to indicate the membership/non-membership of the owner to the provider. For the purpose of privacy preservation, there are noises or false positives added in the matrix; for example, regarding provider  $p_1$  and owner  $t_0$ , value 1 in the published PPI  $M$  is a false positive in the sense that provider  $p_1$  does not have any records of owner  $t_0$  but falsely claims to do so. The false positive value is helpful for obscuring the true and private membership information.

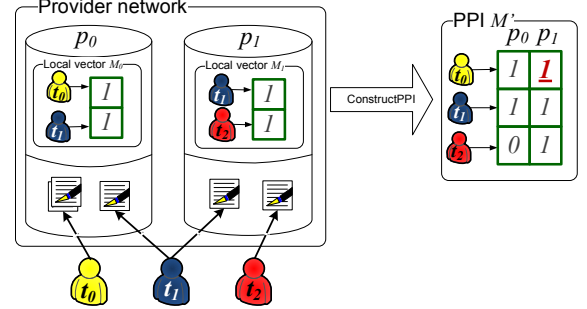


Fig. 2:  $\epsilon$ -PPI model

Table I summarizes the notations that will be used throughout the rest of the paper.

TABLE I: Notations

Symbols of system model			
$t_j$	The $j$ -th owner (identity)	$n$	Number of owners
$\epsilon_j$	Privacy degree of $t_j$		
$p_i$	The $i$ -th provider	$m$	Number of providers
$M_i(\cdot)$	Local vector of $p_i$	$M'(\cdot, \cdot)$	Data matrix in the PPI
Symbols of $\epsilon$ -PPI construction			
$\beta_j$	Publishing probability of $t_j$	$\sigma_j$	Frequency of owner $t_j$
$\lambda$	Percentage of common owners	$f p_j$	Achieved false positive rate of $t_j$

### B. Threat Model

**Privacy goals:** In our work, we are mainly concerned with the owner-membership privacy; for an owner  $t_j$ , the owner-membership privacy is about which providers the owner  $t_j$ 's records belong to, that is,  $M(i, j) = 1$ <sup>3</sup>. Knowing this information, one can develop the private personal knowledge; for example, knowing that a sport celebrity has records stored in a surgery hospital allows one to infer that he or she may have

<sup>2</sup>In this paper, we use "owner" and "identity" interchangeably.

<sup>3</sup>We use  $M(\cdot)$  and  $M'(\cdot, \cdot)$  interchangeably.

had a medical condition possibly requiring surgery and may be absent in the rest of the season. Other privacy goals related to the PPI system but not addressed in this work include searcher anonymity and record content privacy. The searcher anonymity prevents an attacker from knowing which owner(s) a searcher has searched for, which can be protected by various anonymity protocols [17]. The record content privacy [7] involves the detailed content of an owner’s record.

In order to attack the owner-membership privacy, we consider a threat model in which an attacker can exploit multiple information sources through different channels. In particular, we consider the following privacy-threatening scenarios:

- **Primary attack** : The primary attack scenario is that an attacker randomly or intentionally chooses a provider  $p_i$  and an owner  $t_j$ , and then claims that “owner  $t_j$  has delegated his/her records to provider  $p_i$ ”. To determine which providers and owners to attack, the attacker learns about the publicly available PPI data  $M'$ , and attacks only those with  $M'(i, j) = 1$ . Given an owner  $t_j$ , the attacker can randomly or intentionally (e.g. by her prior knowledge) picks a provider  $p_i$  so that  $M'(i, j) = 1$ . To further refine the attack and improve the confidence, the attacker can exploit other knowledge through various channels, such as colluding providers. Due to space limit, we focus on the attack through the public channel in this paper (the colluding attack and analysis can be found in the tech report [18]).
- **Common-identity attack** : This attack focuses on the common identity which appears in almost all providers in the network. The attacker can learn about the truthful frequency of owner identity  $\sigma_j$  from the public PPI matrix  $M'$  (as will be analyzed many PPI’s [9], [7], [8] reveals the truthful frequency) and choose the owners with high frequency. By this means, the attacker can have better confidence in succeeding an attack. For example, consider the following extreme case: by learning an owner identity is with frequency  $\sigma_j = 100\%$ , the attacker can choose any provider and be sure that the chosen provider must be a true positive (i.e.,  $M(i, j) = 1$ ).

This paper focuses on attacks on a single owner, while a multi-owner attack boils down to multiple single-owner attacks.

### C. Privacy Metric and Degrees

*Privacy metric:* We measure the privacy disclosure by the confidence an attacker can succeeding an attack. Formally, given an attack on an owner  $t_j$  and provider  $p_i$ , we measure the privacy disclosure by the probability that the attack can succeed, that is,  $Pr(M(i, j) = 1 | M'(i, j) = 1)$ . To measure the privacy protection level of a specific owner  $t_j$ , we use the average probability of successful attacks against all possible providers that are subject to  $M'(i, j) = 1$ . The privacy metric is formulated as following.

$$\begin{aligned} Pr(M(\cdot, j) | M'(\cdot, j)) &= \text{AVG}_{\forall i, M'(i, j)=1} (Pr(M(i, j) = 1 | M'(i, j) = 1)) \\ &= 1 - fp_j \end{aligned}$$

Here,  $fp_j$  is the false positive rate of providers in the list of providers  $M'(i, j) = 1$ . The privacy disclosure metric on owner  $t_j$  is equal to  $1 - fp_j$ , because the false positive providers determines the probability that an attack can succeed/fail. For example, if the list  $\{p_i | M(i, j) = 1\}$  is completely without any false positive providers (i.e.  $fp_j = 0\%$ ), then attacks on any provider can succeed, leading to  $100\% = 1 - fp_j$  success probability/confidence.

Based on the privacy metric, we further define four discrete privacy degrees. The definition of privacy degrees are based on an information flow model of our privacy threat model, in which an attacker obtains information from the information source through different channels.

- **UNLEAKED:** The information can not flow from the source, and the attacker can not know the information. This is the highest privacy protection level.
- **$\epsilon$ -PRIVATE:** The information can flow to attackers through the channel of public PPI data or PPI construction process. If this occurs, the PPI design protects privacy from being disclosed. The PPI can provide a quantitative guarantee on the privacy leakage. Formally, given a privacy degree  $\epsilon_j$ , this privacy degree requires the quantitative guarantee as follows.

$$Pr(M(\cdot, j) | M'(\cdot, j)) \leq 1 - \epsilon_j \quad (1)$$

In particular, when  $\epsilon = 0\%$ , the attacker might be 100% confident about success of the attack, and privacy is definitely leaked.

- **NOGUARANTEE:** The information can flow to the attacker and the PPI design can not provide any guarantee on privacy leakage. That is, the achieved value of privacy leakage metric may be unpredictable.
- **NOPROTECT:** The information can flow to the attacker and the PPI design does not address the privacy preservation. That is, the privacy is definitely leaked and the attack can succeed with 100% certainty. This is equivalent to the special case of NOGUARANTEE where  $\epsilon_j = 0\%$ . This is the lowest level of privacy preservation.

Based on our privacy model and metric, we can summarize the prior work in Table II. Due to the space limitation, we put the analysis in Appendix B.

TABLE II: Comparison of  $\epsilon$ -PPI against existing PPI’s

	Primary attack	Common-identity attack
PPI [7], [8]	NOGUARANTEE	NOGUARANTEE
SS-PPI [9]	NOGUARANTEE	NOPROTECT
$\epsilon$ -PPI	$\epsilon$ -PRIVATE	$\epsilon$ -PRIVATE

### D. Index Construction of Quantitative Privacy Preservation

In the  $\epsilon$ -PPI, we aim at achieving  $\epsilon$ -PRIVATE on a per-identity basis (i.e. differentiating privacy preservation for different owners). The formal problem that this paper address is the index construction of quantitative privacy preservation, which is stated as below.

*Proposition 2.1:* Consider a network with  $m$  providers and  $n$  owners; each provider  $p_i$  has a local Boolean vector  $M_i$  of its membership of  $n$  owners. Each owner  $t_j$  has a preferred level of privacy preservation  $\epsilon_j$ . The problem of quantitative privacy preserving index construction is to construct a PPI that can bound any attacker's confidence (measured by our per-owner privacy metric) under  $\epsilon_j$ , with regards to all attacks on owner  $t_j$  as described in our threat model.

### III. $\epsilon$ -PPI Construction: the Computation

Our  $\epsilon$ -PPI construction is based on a proposed two-phase framework in which providers first collectively calculate a global value  $\beta$ , and then each provider independently publishes its local vector randomly based on probability  $\beta$ . This framework requires complex computations. In this section, we introduce them at different levels of granularity: First we take an overview of our two-phase construction framework with emphasis on describing the second phase. We then introduce the first phase (called the  $\beta$  calculation) in details; we present the detailed calculation of  $\beta$  under two kinds of owner identities, namely the common and non-common owners. At last, we conduct the privacy analysis.

#### A. A Two-Phase Construction Framework

We propose a two-phase framework for the  $\epsilon$ -PPI construction. First, for each owner identity  $t_j$ , all  $m$  providers collectively calculate a probability value  $\beta_j$ . In the second phase, the private membership value regarding owner  $t_j$  and every provider  $p_i$  is published. In this paragraph, we assume  $\beta_j$  is already calculated and focus on describing the second phase – how to use  $\beta_j$  to publish private data. Recall that in our data model, each provider  $p_i$  has a Boolean value  $M(i, j)$  that indicates the membership of owner  $t_j$  in this provider. After knowing value of  $\beta_j$ , provider  $p_i$  starts to publish this private Boolean value by randomly flipping it at probability  $\beta_j$ . To be specific, given a membership Boolean value (i.e.  $M(i, j) = 1$ ), it is always truthfully published as 1, that is,  $M'(i, j) = 1$ . Given a non-membership value (i.e.  $M(i, j) = 0$ ), it is negated to  $M'(i, j) = 1$  at probability  $\beta_j$ . We call the negated value as the false positive in the published PPI. The following formula describes the randomized publication. Note when Boolean value  $M(i, j) = 1$ , it is not allowed to be published as  $M'(i, j) = 0$ .

$$\begin{aligned} 0 &\rightarrow \begin{cases} 1, \text{ with probability } \beta \\ 0, \text{ with probability } 1 - \beta \end{cases} \\ 1 &\rightarrow 1 \end{aligned} \quad (2)$$

The truthful publication rule (i.e.  $1 \rightarrow 1$ ) guarantees that relevant providers are always in the QueryPPI result and the 100% query recall is ensured. The false-positive publication rule (i.e.  $0 \rightarrow 1$ ) adds noises or false positives to the published PPI which can help obscure the true owner-to-provider membership and thus preserves owner-membership privacy. For multiple owners, different  $\beta$ 's are calculated and the randomized publication runs independently.

*An example:* Consider the case in Figure 2. For owner  $t_0$ , if the  $\beta_0$  is calculated to be 0.5, then provider  $p_1$  would publish its negative membership value  $M_1(0) = 0$  as value 1 with probability 0.5. In this example, it is flipped and the constructed  $\epsilon$ -PPI contains  $M'(1, 0) = 1$ . Similarly for identity  $t_2$  and provider  $p_0$ , it is also subject to flipping at probability  $\beta_2$ . In this example, it is not flipped, and the constructed  $\epsilon$ -PPI contains  $M'(0, 2) = 0$ .

#### B. The $\beta$ Calculation

In the randomized publication,  $\beta_j$  determines the amount of false positives in the published  $\epsilon$ -PPI. For quantitative privacy preservation, it is essential to calculate a  $\beta_j$  value that makes the false positive amount meet the privacy requirement regarding  $\epsilon_j$ . In this part, we focus on the calculation of  $\beta$  which serves as the first phase in  $\epsilon$ -PPI construction process. Concretely we consider two cases: the common identity case and the non-common identity case. Recall that the common identity refers to such an owner who delegates her records to almost all providers in the network. The general PPI construction is vulnerable to the common-identity attack and it needs to be specially treated.

1) *The Case of Non-common Identity:* In the case of non-common identity, negative providers suffice to meet the desired privacy degree. We consider the problem of setting value  $\beta_j$  for identity  $t_j$  in order to meet the desired  $\epsilon_j$ . Recall the randomized publication: Multiple providers independently runs an identical random process, and this can be modeled as a series of Bernoulli trials (note that the publishing probability  $\beta(t_j)$  is the same to all providers). Our goal is to achieve privacy requirement that  $fp_j \geq \epsilon_j$  with high level success rate  $p_p$ , that is,  $p_p = Pr(fp_j \geq \epsilon_j)$ . Under this model, we propose three policies to calculate  $\beta$  with different quantitative guarantees: a basic policy  $\beta_b$  that guarantees  $fp_j \geq \epsilon_j$  with 50% probability, and an incremented expectation based policy  $\beta_d$ , and a Chernoff bound based policy  $\beta_{c}$  that guarantees  $fp_j \geq \epsilon_j$  with  $\gamma$  probability where success rate  $\gamma$  can be configured.

*Basic policy:* The basic policy sets the  $\beta$  value so that the expected amount of false positives among  $m$  providers can reach a desired level, which is,  $\epsilon_j \cdot m(1 - \sigma_j)$ . We can have the following,

$$\begin{aligned} \epsilon_j &= \frac{(1 - \sigma_j) \cdot \beta_b(t_j)}{(1 - \sigma_j) \cdot \beta_b(t_j) + \sigma_j} \\ \Rightarrow \beta_b(t_j) &= [(\sigma_j^{-1} - 1)(\epsilon_j^{-1} - 1)]^{-1} \end{aligned} \quad (3)$$

The basic policy has poor quality in attaining the desired privacy preservation; the actual value  $fp_j$  is bigger than  $\epsilon_j$  with only 50% success rate.

*Incremented expectation-based policy:* The incremented expectation-based approach is to increase the expectation-based  $\beta_b(t_j)$  by a constant value, that is,

$$\beta_d(t_j) = \beta_b(t_j) + \Delta \quad (4)$$

Incremental  $\Delta$  can be configurable based on the quality requirement; the bigger the value is, the higher success rate  $p_p$  is expected to attain. However, there is no direct connection between the configured value of  $\Delta$  and the success rate  $p_p$  that

can be achieved, leaving it a hard task to figure out the right value of  $\Delta$  based on desired  $p_p$ .

*Chernoff bound-based policy:* Toward an effective policy to calculate  $\beta$ , we apply the Chernoff bounds to the Bernoulli trial model of the randomized publication process. This policy allows direct control of the success rate. Formally, it has the property described in Theorem 3.1 (with the proof in Appendix A-A).

*Theorem 3.1:* Given desired success rate  $\gamma > 50\%$ , let  $G_j = \frac{\ln \frac{1}{1-\gamma}}{(1-\sigma_j)m}$  and

$$\beta_c(t_j) \geq \beta_b(t_j) + G_j + \sqrt{G_j^2 + 2\beta_b(t_j)G_j} \quad (5)$$

Then, randomized publication with  $\beta(t_j) = \beta_c(t_j)$  statistically guarantees the published  $\epsilon$ -PPI with privacy requirement  $fp_j \geq \epsilon_j$  with success rate larger than  $\gamma$ .

2) *The Case of Common Identities:* With the above  $\beta$  calculation for non-common identities, the constructed  $\epsilon$ -PPI is vulnerable to the common-identity attack. Because the  $\beta_*$ <sup>4</sup> bears information of identity frequency  $\sigma_j$ , and during our index construction framework,  $\beta$  needs to be released to all participating providers. A colluding provider would release such information to the attacker who can easily obtain the truthful identity frequency  $\sigma$  (e.g., from Equation 3 assuming  $\epsilon_j$  is publicly known) and effectively formulates the common-identity attack.

To defend against the common-identity attack,  $\epsilon$ -PPI construction employs an identity-mixing technique for common identities. The idea is to mix common identities with certain non-common identities by exaggerating the calculated  $\beta_j$  (i.e. falsely increasing certain  $\beta_j$  to 100%) from which one can not distinguish common identities from the rest. To be specific, for a non-common identity  $t_j$ , we allow its  $\beta_j$  to be exaggerated to 100% with probability  $\lambda$ , that is,

$$\beta = \begin{cases} \begin{cases} \beta_*, & 1 - \lambda \\ 1, & \lambda \end{cases}, & \beta_* < 1 \\ 1, & \beta_* \geq 1 \end{cases} \quad (6)$$

Given a set of common identities, we need to determine how many non-common identities should be chosen for mixing, in other words, to determine the value of  $\lambda$ . While a big value of  $\lambda$  can hide common identities among the non-common ones, it incurs unnecessarily high search cost. On the other hand, a value of  $\lambda$  which is too small would leave common identities unprotected and vulnerable. In  $\epsilon$ -PPI, we use the following heuristic-based policy to calculate  $\lambda$ .

- In the set of mixed identities, the percentage of non-common identities should be no smaller than  $\xi$ . Since there are  $\sum_{\beta_* \geq 1} 1$  common identities and thus  $\sum_{\beta_* < 1} \lambda$  non-common identities in the set, we have the following

formula.

$$\begin{aligned} \xi &\leq \frac{\sum_{\beta_* < 1} \lambda}{\sum_{\beta_* \geq 1} 1 + \sum_{\beta_* < 1} \lambda} \\ \Rightarrow \lambda &\geq \frac{\xi}{1 - \xi} \cdot \frac{\sum_{\beta_* \geq 1} 1}{n - \sum_{\beta_* \geq 1} 1} \end{aligned} \quad (7)$$

3)  *$\beta$  Calculation: Putting It Together:* We summarize the  $\beta$  calculation in the  $\epsilon$ -PPI construction. For each identity  $t_j$ ,  $\beta(t_j)$  is calculated based on Equation 6, which follows the computation flows as below. The underline symbol indicates the variable is private and  $\Rightarrow$  indicates the computation is fairly complex (e.g. involving square root when calculating  $\beta_*$ ).

$$\begin{aligned} \text{Frequency } \underline{\sigma} &\Rightarrow \text{Raw probability } \underline{\beta_*} \rightarrow \\ &\rightarrow \sum_{\underline{\beta_*} \geq 1} 1 \rightarrow \text{Common id percentage } \lambda \rightarrow \text{Final probability } \underline{\beta} \end{aligned} \quad (8)$$

### C. Privacy Analysis of Constructed $\epsilon$ -PPI

We present the privacy analysis of the constructed  $\epsilon$ -PPI under our threat model.

*Privacy under primary attack:* The property of the three policies of calculating  $\beta_*$  suggests that the false positive rate in the published  $\epsilon$ -PPI should be no smaller than  $\epsilon_j$  in a statistical sense. Recall that the false positive rate bounds the attacker's confidence; it implies that  $\epsilon$ -PPI achieves an  $\epsilon$ -PRIVATE degree against the primary attack. It is noteworthy that our  $\epsilon$ -PPI is fully resistant to repeated attacks against the same identity over time, because the  $\epsilon$ -PPI is static; once constructed and having privacy protected, it stays the same.

*Privacy under common-identity attack:* For common-identity attack, the attacker's confidence in choosing a true common identity depends on the percentage of true common identities among the (mixed) common identities in the published  $\epsilon$ -PPI. Therefore the privacy preservation degree is bounded by the percentage of false positives (in this case, it depends on the percentage of the non-common identities which is mixed and published as common identities in the published  $\epsilon$ -PPI), which equals  $\xi$ . By properly setting  $\lambda$ , we can have  $\xi = \max_{\forall t_j \in \{\text{common identities}\}} \epsilon_j$ . By this way, it is guaranteed to achieve the per-identity  $\epsilon$ -PRIVATE degree against the common-identity attack.

## IV. $\epsilon$ -PPI Construction: Realization

The information network lacks mutual trusts between providers, which poses new challenges when putting the  $\epsilon$ -PPI construction in practice. This section describes the design and implementation of a distributed and secure protocol that realizes the computation of  $\epsilon$ -PPI construction described in the previous section.

### A. Challenge and Design

The goal of our protocol is to efficiently and securely compute the publishing probability  $\{\beta_j\}$  among a set of mutually untrusted providers who are reluctant to exchange the private

<sup>4</sup>We use  $\beta_*$  to denote the probability value calculated by any of the three policies for non-common identities.



membership vector with others. On one hand, the secure computation would require multi-party computation (or MPC) which respects the per-party input privacy. Current techniques for MPC only support small computation workloads [16]. On the other hand, the computation required in  $\epsilon$ -PPI construction is big and complex; the computation model involves large number of identities and providers; even for a single identity involves fairly complex computation (e.g., square root and logarithm as in Equation 5). This poses a huge challenge to design a practical protocol for secure  $\epsilon$ -PPI construction.

To address the above challenge, we propose an efficient and secure construction protocol by following the design principle of *minimizing the secure computation*. Given a computation flow in Equation 8, our secure protocol design has three salient features: 1) It separates the secure and non-secure computations by the last appearance of private variables in the flow (note that the computation flows from the private data input to the end of non-private result). 2) It reorders the computation to minimize the expensive secure computation. The idea is to push down complex computation towards the non-private end. To be specific, instead of first carrying out complex floating point computations for raw probability  $\beta$ , as in Formula 8, we push such computations down through the flow and pull up the obscuring computations for private input, as in Formula 9. 3) To scale to a large number of providers, we propose an efficient protocol for calculating secure sum, and use it to reduce the “core” of the MPC part in  $\epsilon$ -PPI construction.

$$\underline{\sigma} \rightarrow \sum_{\underline{\sigma} < \sigma'} 1 \rightarrow \lambda \rightarrow \begin{cases} \rightarrow \beta = 1 \\ \Rightarrow \beta = \beta_* \end{cases} \quad (9)$$

## B. The Distributed Algorithm

Following our design, we propose a practical distributed algorithm to run the two-phase  $\epsilon$ -PPI construction. The overall workflow is illustrated in Figure 3. For simplicity, we focus on phase 1 for  $\beta$  calculation. The  $\beta$  calculation is realized in two stages by itself: As illustrated in Algorithm 1, the first stage is a SecSumShare protocol which, given  $m$  input Boolean from the providers, outputs  $c$  secret shares whose sum is equal to the sum of these  $m$  Boolean. Here,  $c$  is the number of shares that can be configurable based on the tolerance on provider collusion. The output  $c$  shares have the security property that a party knowing  $x < c$  shares can not deduce any information about the sensitive sum of  $m$  Boolean. For different identities, the SecSumShare protocol runs multiple instances independently and in parallel, which collectively produce  $c$  vectors of shares, denoted by  $s(i, \cdot)$ , where  $i \in [0, c-1]$ . The  $c$  vectors are distributed to  $c$  coordinate providers (for simplicity we assume they are providers  $p_0, \dots, p_{c-1}$ ) on which the second-stage protocol, CountBelow, is run. As shown by Algorithm 2, given  $c$  vectors  $s(0, \cdot), \dots, s(c-1, \cdot)$  and a threshold  $t$ , the CountBelow algorithm sums them to vector  $\sum_i s(i, \cdot)$  and counts the number of elements that are bigger than  $t$ .

1) *Distributed Algorithm for SecSumShare*: We use an example in the top box in Figure 3 to illustrate the distributed algorithm of SecSumShare. In the example  $c = 3$  and there

TABLE III: Distributed algorithms for  $\epsilon$ -PPI construction

**Algorithm 1** calculate-beta( $M_0, \dots, M_{n-1}$ )

- 1:  $\{s(0, \cdot), \dots, s(c-1, \cdot)\} \leftarrow \text{SecSumShare}(M_0, \dots, M_{n-1})$
- 2:  $\sigma'(\cdot)$  is calculated under condition  $\beta_* = 1$ , by either Equation 3, or 4 or 5.
- 3:  $\sum_{\sigma \geq \sigma'} 1 \leftarrow \text{CountBelow}(s(0, \cdot), \dots, s(c-1, \cdot), \sigma'(\cdot) \cdot m)$
- 4:  $\{\beta_0, \dots, \beta_{m-1}\} \leftarrow \sum_{\sigma \geq \sigma'} 1 \quad \triangleright \text{By Equation 9}$

**Algorithm 2** CountBelow( $s(0, \cdot), \dots, s(c-1, \cdot)$ , threshold  $t$ )

- 1: count  $\leftarrow 0$
- 2: **for**  $\forall j \in [0, m-1]$  **do**
- 3:  $S[j] \leftarrow \sum_i s(i, j)$
- 4: **if**  $S[j] < t$  **then**
- 5: count++
- 6: **end if**
- 7: **end for**
- 8: **return** count

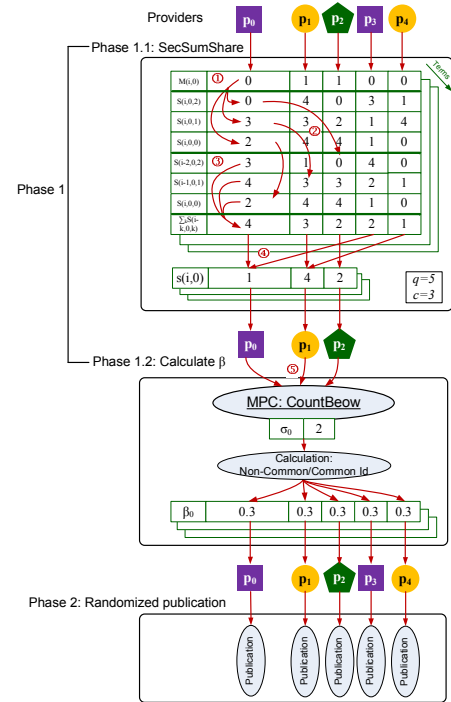


Fig. 3: An example of  $\epsilon$ -PPI construction algorithm

are five providers  $p_0, \dots, p_4$ . The example focuses on a single identity case for  $t_j$  (e.g.  $j = 0$ ). Out of the 5 providers,  $p_1$  and  $p_2$  have records of owner  $t_0$  (i.e.,  $M(1, 0) = M(2, 0) = 1$ ). SecSumShare requires modular operations; in this example, the modulus divisor is  $q = 5$ . It runs in the following 4 steps.

- 1 Generating shares: each provider  $p_i$  decomposes its private input Boolean  $M(i, j)$  into  $c$  shares, denoted by  $\{S(i, j, k)\}$ , with  $k \in [0, c-1]$ . The first  $c-1$  shares are randomly picked from interval  $[0, q]$  and the last share is deterministically chosen so that the sum of all shares equals the input Boolean  $M(i, j)$  in modulo  $q$ . That is,  $(\sum_{k \in [0, c]} S(i, j, k)) \bmod q = M(i, j)$ . In Figure 3, as depicted by arrows ①,  $p_0$ 's input  $M(0, 0)$  is decomposed to  $c = 3$  shares,  $\{S(0, 0, k)\}_{k=0}^2 = \{2, 3, 0\}$ . It ensures  $(2 + 3 + 0) \bmod 5 = 0$ .

- 2 Distributing shares: each provider  $p_i$  then distributes her shares to the next  $c - 1$  neighbor providers;  $k$ -th shares  $S(i, j, k)$  will be sent out to  $k$ -th successor of provider  $p_i$ , that is,  $p_{(i+k) \bmod m}$ . As shown by arrows ② in Figure 3,  $p_0$  keeps the first share 2 locally, sends her second share 3 to her successor  $p_1$  and the third share 0 to 2-hop successor  $p_2$ .
- 3 Summing shares: each provider then sums up all shares she has received in the previous step to obtain the *super-share*. In Figure 3, after the step of share distribution, provider  $p_0$  receives 3 from  $p_3$ , 4 from  $p_4$  and 2 from herself. As depicted by arrows ③, the super-share is calculated to be  $3 + 4 + 2 \bmod 5 = 4$ .
- 4 Aggregating super-shares: each provider sends her super-share to a set of  $c$  coordinators. These coordinators receiving super-shares then sum the received shares up and output the summed vector  $s(i, \cdot)$  to the next-stage CountBelow protocol. In Figure 3, provider  $p_0, p_1, p_2$  are chosen as coordinators and arrow ④ shows that the sum of super-shares on provider  $p_0$  is  $s(0, 0) = (4 + 2) \bmod 5 = 1$ . The sum of all the values on coordinators should be equal to the number of total appearances of identity  $t_0$ . That is,  $1 + 4 + 2 \bmod 5 = 2$ . Note two providers have identity  $t_0$ . This total appearance number or identity frequency may be sensitive (in the case of common identity) and can not be disclosed immediately, which is why we need the second stage protocol, CountBelow.

2) *Implementation of CountBelow computation:* The secure computation of CountBelow (in Algorithm 2) is implemented by using a generic MPC protocol. Each party corresponds to a coordinate provider in the  $\epsilon$ -PPI system. Specifically, we choose a Boolean-circuit based MPC protocol FairplayMP [13] for implementation. The reason is that compared to an arithmetic-circuit based protocol, it lends itself to the computation of comparison required in Algorithm 2 (i.e., in Line 4). In particular for  $c = 2$ , the computation in CountBelow essentially boils down to a comparison operation (i.e.,  $s(0, i) > t - s(1, i)$ ), and the problem is reduced to a Millionaire problem [19]. The distributed algorithm to carry out MPC (and thus our MPC-based CountBelow computation) can be found in [12], [13]. Since Algorithm 2 is implemented by expensive MPC it normally becomes the bottleneck of the system; in practice,  $c \ll m$  and thus the network can scale to large number of providers  $m$  while the MPC is still limited to small subset of the network.

### C. Privacy Analysis of Constructing $\epsilon$ -PPI

We analyze the privacy preservation of  $\epsilon$ -PPI construction process. We mainly consider a semi-honest model, which is consistent with the existing MPC work [13]. The privacy analysis is conducted from three aspects: 1) The privacy guarantee of SecSumShare protocol. It guarantees: 1.1)  $(2c - 3)$ -secrecy of input privacy [9]: With less than  $c$  providers in collusion, none of any private input can be learned by providers other than its owner. 1.2)  $c$ -secrecy of output privacy: the private sum can only be reconstructed when all  $c$  shares are used. With less than

$c$  shares, one can learn nothing regarding the private sum. The output privacy is formally presented in Theorem 4.1 with proof in Appendix A-B. 2) The security and privacy of CountBelow relies on that of the MPC used in implementation. The generic MPC technique can provide information confidentiality against colluding providers on  $c$  participating parties [13]. 3) The final output  $\beta$  does not carry any private information, and is safe to be released to the (potentially untrusted) providers for randomized publication.

*Theorem 4.1:* The SecSumShare's output is a  $(c, c)$  secret sharing scheme. Specifically, for an owner  $t_j$ , SecSumShare protocol outputs  $c$  shares,  $\{s(i, j) | \forall i \in [0, c - 1]\}$ , whose sum is the secret  $v_j$ . The  $c$  shares have the following properties.

- *Recoverability:* Given  $c$  output shares, the secret value  $v_j$  (i.e. the sum) can be easily reconstructed.
- *Secrecy:* Given any  $c - 1$  or fewer output shares, one can learn nothing about the secret value, in the sense that the conditional distribution given the known shares is the same as the prior distribution,

$$\forall x \in \mathbb{Z}_q, Pr(v_j = x) = Pr(v_j = x | V \subset \{s(i, j)\})$$

where  $V$  is any proper subset of  $\{s(i, j)\}$ .

## V. Experiments

To evaluate the proposed  $\epsilon$ -PPI, we have done two set of experiments: The first set, with simulation-based experiments, evaluates how effective the  $\epsilon$ -PPI can be in terms of delivering quantitative privacy protection, and the second set evaluates the performance of our index construction protocol. For realistic performance study, we have implemented a functioning prototype for  $\epsilon$ -PPI construction.

### A. Effectiveness of Privacy Preservation

*Experimental setup:* To simulate the information provider network, we used a distributed document dataset [20] of 2,500 – 25,000 small digital libraries, each of which simulates a provider in our problem setting. To be specific, this dataset defines a “collection” table, which maintains the mapping from the documents to collections. The documents are further derived from NIST’s publicly available TREC-WT10g dataset [21]. To adapt to our problem setting, each collection is treated as a provider and the source web URLs (as defined in TREC-WT10g dataset) of the documents are treated as owner’s identity. If not otherwise specified, we use no more than 10,000 providers in the experiments. Using the collection table, it also allows us to emulate the membership matrix  $M$ . The dataset does not have a privacy metric for the query phrase. In our experiment, we randomly generate the privacy degree  $\epsilon$  in the domain  $[0, 1]$ . We use a metric, success ratio, to measure the effectiveness. The success rate is the percentage of identities whose false positive rates in the constructed PPI are no smaller than the desired rate  $\epsilon_j$ . Due to space limit, the experiment results of different  $\beta$  calculation policies can be found in Appendix C-A.



1)  $\epsilon$ -PPI versus Existing Grouping-based PPI's: The experiments compare the  $\epsilon$ -PPI with existing PPI's. The existing PPI's [7], [8], [9] are based on a grouping abstraction; providers are organized into disjoint privacy groups so that different providers from the same group are indistinguishable from the searchers. By contrast,  $\epsilon$ -PPI does not utilize grouping technique and is referred to in this section as a non-grouping approach. In the experiment, we measure the success rate of privacy preservation and search performance. Grouping PPI's are tested under different group sizes. Given a network of fixed providers, we use the group number to change average group size. We test grouping PPI with the Chernoff bound-based and the incremented expectation-based policy under the default setting. The expected false positive rate is configured at 0.8, and the number of providers is 10,000. We uniformly sample 20 times and report the average results.

Results are illustrated in Figure 4. Non-grouping PPI generally performs much better and is more stable than the grouping approach in terms of success ratio. With proper configuration (e.g.  $\Delta = 0.01$  for incremental expectation-based policy and  $\gamma = 0.9$  for Chernoff policy), the non-grouping  $\epsilon$ -PPI always achieves near-1.0 success ratio. By contrast, the grouping PPI's display instability in their success ratio. For example, as shown by the "Grouping (#groups 2000)" series in Figure 4a, the success rate fluctuates between 0.0 and 1.0, which renders it difficult to provide a guarantee to the system and owners. The reason is that with 2000 groups, sample space in each group is too small (i.e., with 50 providers) to hold a stable result for success ratio. When varying  $\epsilon$ , similar behavior is shown in Figure 4b; the success rate of grouping PPI's quickly degrades to 0, leading to unacceptable privacy quality. This is due to the owner agnostic design in grouping PPI. This set of experiments shows that the privacy degree of non-grouping PPI's can be effectively tuned in practice, implying the ease of configuration and more control exposed to applications.

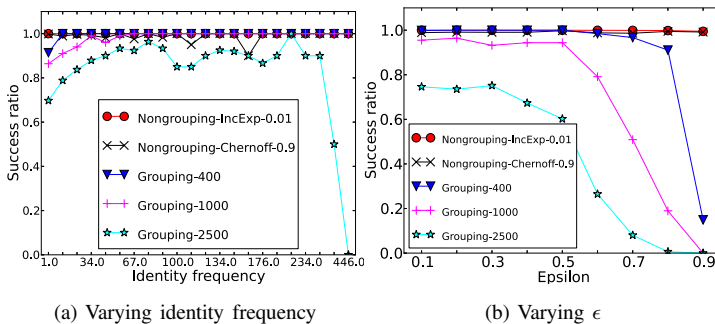


Fig. 4: Comparing non-grouping and grouping

## B. Performance of Index Construction

*Experimental setup:* We evaluate the performance of our distributed  $\epsilon$ -PPI construction protocol. Towards that, we have implemented a functioning prototype. The CountBelow is implemented by using an MPC software, FairplayMP [13], which is based on Boolean circuits. The implemented CountBelow

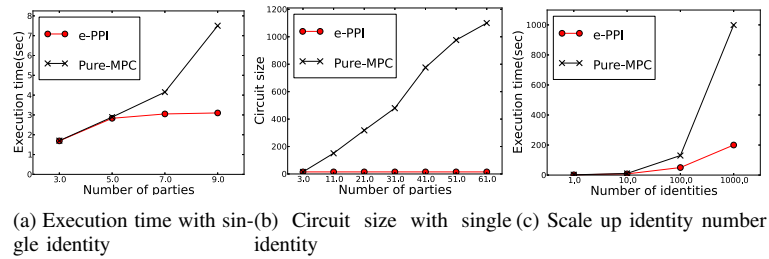


Fig. 5: Performance of index construction protocol

protocol is written in SFDL, a secure function definition language exposed by FairplayMP, and is compiled by the FairplayMP runtime to Java code, which embodies the generated circuit for secure computation. We implement the SecSumShare protocol in Java. In particular, we use a third-party library Netty [22] for network communication and Google's protocol buffer [23] for object serialization. We conduct experiments on a number of machines in Emulab [24], [25], each equipped with a 2.4 GHz 64-bit Quad Core Xeon processor and 12 GB RAM. In the experiments, the number of machines tested is varied from 3 to 9 (due to limited resource at hand). For each experiment, the protocol is compiled to and run on the same number of parties. Each party is mapped to one dedicated physical machine. The experiment uses a configuration of  $c = 3$ .

To justify the standpoint of our design that MPC is expensive, we compare our reduced-MPC approach as in the  $\epsilon$ -PPI construction protocol against a pure MPC approach. The pure MPC approach does not make use of SecSumShare protocol to reduce the number of parties involved in the generic MPC part and directly accepts inputs from the  $m$  providers. The metric used in the experiment is the start-to-end execution time, which is the time duration from when the protocol starts to run to when the last machine reports to finish. The result is shown as in Figure 5a. It can be seen that the pure MPC approach generally incurs longer execution time than our reduced-MPC approach (used in  $\epsilon$ -PPI construction): As the information network grows large, while the execution time of pure MPC approach increases super-linearly, that of reduced-MPC approach increases slowly. This difference is due to the fact that the MPC computation in our reduced-MPC approach is fixed to  $c$  parties and does not change as the number of providers  $m$  grows. And the parallel SecSumShare in reduced-MPC approach is scalable in  $m$  as well, since each party runs in constant rounds, and each round sends a constant number (at most  $c - 1$ ) of messages to its neighbors. For scaling with more parties, we use the metric of circuit size, which is the size of the compiled MPC program. As a valid metric, the circuit size determines the execution time<sup>5</sup> in real runs. By this means, we can show the scalability result of up to 60 parties as in Figure 5b. Similar performance improvement can be observed except that the circuit size grows linearly with number of parties involved. Finally, we also study the scalability from running the protocol with multiple identities in a three-party network. The result in Figure 5c shows that  $\epsilon$ -PPI construction grows with the number of identities at a much

<sup>5</sup>Regarding the detailed definition of circuit size and the exact correlation between circuit size and execution time, it can be found in FairplayMP [13].

slower rate than that of the pure MPC approach.

## VI. Related Work

This section surveys related work on indexing support on untrusted servers. We focus on information confidentiality or privacy on secure index design, and do not survey the issues of integrity and authenticity.

*Non-encryption based privacy preserving index:* PPI is designed to index access controlled contents scattered across multiple content providers. While being stored on an untrusted server, PPI aims at preserving the content privacy of all participant providers. Inspired by the privacy definition of  $k$ -anonymity [26], existing PPI work [7], [8], [9] follows the *grouping-based* approach; it organizes providers into disjoint privacy groups, such that providers from the same group are indistinguishable to the searchers. To construct such indexes, many existing approaches [7], [8], [27] assume providers are willing to disclose their private local indexes, an unrealistic assumption when there is a lack of mutual trust between providers. SS-PPI [9] is proposed with resistance against colluding attacks. While most existing grouping PPI's utilize a randomized approach to form groups, its weakness is studied in SS-PPI but without a viable solution. Though the group size can be used to configure grouping-based PPI's, it lacks per-owner concerns and quantitative privacy guarantees. Moreover, organizing providers in groups usually leads to query broadcasting (e.g. with positive providers scattered in all groups), rendering search performance inefficient. By contrast,  $\epsilon$ -PPI is a brand new PPI abstraction without grouping (i.e. non-grouping PPI as mentioned before), which provides quantitative privacy control on a per-owner basis.

*Secure index and search-able encryption:* Building searchable indexes over encrypted data has been widely studied in the context of both symmetric key cryptography [28] and public key cryptography [29], [30], [31]. In this architecture, content providers build their local indices and encrypt all the data and indices before submitting them to the untrusted server. During query time, the searcher first gets authenticated and authorized by the corresponding content provider; the searcher then contacts the untrusted server and searches against the encrypted index. This system architecture makes the assumption that a searcher already knows which provider possesses the data of her interest, which is unrealistic in the PPI scenario. Besides, unlike the encryption-based system, performance is a motivating factor behind the design of our PPI, by making no use of encryption during the query serving time.

## VII. Conclusion

In this paper, we propose  $\epsilon$ -PPI for quantitative privacy control in information networks. The privacy of our  $\epsilon$ -PPI can be controlled by each individual in a quantitative fashion. We identify a vulnerability of generic PPI on protecting common owner identities and address this vulnerability in our  $\epsilon$ -PPI design by proposing an identity mixing technique. We

have implemented the index construction protocol without any trusted party and applied a performance-optimization design that minimizes the amount of secure computation. We have built a generic privacy threat model and performed security analysis which shows the advantages of  $\epsilon$ -PPI over other PPI system in terms of privacy preservation quality.

## Acknowledgement

This research is partially supported by grants from NSF CISE NetSE program, SaTC program, I/UCRC, an IBM faculty award and a grant from Intel ICST on Cloud Computing.

## References

- [1] "Nhin: <http://www.hhs.gov/healthit/healthnetwork/>."
- [2] "Nhin connect, <http://www.connectopensource.org/>."
- [3] "Hippa, <http://www.cms.hhs.gov/hipaaeninfo/>."
- [4] "Studip, <http://www.studip.de/>."
- [5] "Swiki, <http://en.wikipedia.org/wiki/swiki>."
- [6] "Ferpa, <http://www2.ed.gov/ferpa/>."
- [7] M. Bawa, R. J. B. Jr., and R. Agrawal, "Privacy-preserving indexing of documents on the network," in *VLDB*, 2003, pp. 922–933.
- [8] M. Bawa, R. J. Bayardo, Jr, R. Agrawal, and J. Vaidya, "Privacy-preserving indexing of documents on the network," *The VLDB Journal*, vol. 18, no. 4, 2009.
- [9] Y. Tang, T. Wang, and L. Liu, "Privacy preserving indexing for health information networks," in *CIKM*, 2011, pp. 905–914.
- [10] S. Zerr, E. Demidova, D. Olmedilla, W. Nejdl, M. Winslett, and S. Mitra, "Zerber: r-confidential indexing for distributed documents," in *EDBT*, 2008, pp. 287–298.
- [11] "Prism, [http://en.wikipedia.org/wiki/prism\\_\(surveillance\\_program\)](http://en.wikipedia.org/wiki/prism_(surveillance_program))."
- [12] D. Malkhi, N. Nisan, B. Pinkas, and Y. Sella, "Fairplay - secure two-party computation system," in *USENIX Security Symposium*, 2004, pp. 287–302.
- [13] A. Ben-David, N. Nisan, and B. Pinkas, "Fairplaymp: a system for secure multi-party computation," in *ACM Conference on Computer and Communications Security*, 2008, pp. 257–266.
- [14] W. Henecka, S. Kögl, A.-R. Sadeghi, T. Schneider, and I. Wehrenberg, "Tasty: tool for automating secure two-party computations," in *ACM CCS*, 2010, pp. 451–462.
- [15] I. Damgård, M. Geisler, M. Krøigaard, and J. B. Nielsen, "Asynchronous multiparty computation: Theory and implementation," in *Public Key Cryptography*, 2009, pp. 160–179.
- [16] A. Narayan and A. Haeberlen, "DJoin: differentially private join queries over distributed databases," in *OSDI*, Oct. 2012.
- [17] M. Wright, M. Adler, B. N. Levine, and C. Shields, "An analysis of the degradation of anonymous protocols," in *NDSS*, 2002.
- [18] Y. Tang and L. Liu, "Searching information networks with quantitative privacy guarantee," *Georgia Tech Technical Report 2012*, <http://www.cc.gatech.edu/~ytang36/docs/techreport-12.pdf>.
- [19] A. C.-C. Yao, "Protocols for secure computations (extended abstract)," in *FOCS*, 1982, pp. 160–164.
- [20] J. Lu and J. P. Callan, "Content-based retrieval in hybrid peer-to-peer networks," in *CIKM*, 2003, pp. 199–206.
- [21] D. Hawking, "Overview of the trec-9 web track," in *TREC*, 2000.
- [22] "Netty: <http://netty.io/>."
- [23] "Protobuf: <http://code.google.com/p/protobuf/>."
- [24] "<http://www.emulab.net/>."
- [25] B. White, J. Lepreau, L. Stoller, R. Ricci, S. Guruprasad, M. Newbold, M. Hibler, C. Barb, and A. Joglekar, "An integrated experimental environment for distributed systems and networks," in *OSDI*, 2002.
- [26] L. Sweeney, "k-anonymity: A model for protecting privacy," *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 10, no. 5, pp. 557–570, 2002.
- [27] M. Bawa, R. J. B. Jr., S. Rajagopalan, and E. J. Shekita, "Make it fresh, make it quick: searching a network of personal webservers," in *WWW*, 2003, pp. 577–586.
- [28] D. X. Song, D. Wagner, and A. Perrig, "Practical techniques for searches on encrypted data," in *IEEE SSP*, 2000, pp. 44–55.

- [29] C. Wang, N. Cao, J. Li, K. Ren, and W. Lou, "Secure ranked keyword search over encrypted cloud data," in *ICDCS*, 2010, pp. 253–262.
- [30] M. Li, S. Yu, N. Cao, and W. Lou, "Authorized private keyword search over encrypted data in cloud computing," in *ICDCS*, 2011, pp. 383–392.
- [31] N. Cao, C. Wang, M. Li, K. Ren, and W. Lou, "Privacy-preserving multi-keyword ranked search over encrypted cloud data," in *INFOCOM*. IEEE, 2011, pp. 829–837.
- [32] M. Mitzenmacher and E. Upfal, *Probability and computing - randomized algorithms and probabilistic analysis*. Cambridge University Press, 2005.

## Appendix A

### Proof of theorems

#### A. Proof of Theorem 3.1

*Proof:* We model the problem as Bernoulli trials and prove the theorem by applying Chernoff bound. For a term  $t_j$ , the total number of false positive providers is modeled as sum of  $T = m(1 - \sigma_j)$  Bernoulli trials, because there are  $m(1 - \sigma_j)$  negative providers for term  $t_j$  and each negative provider independently and randomly publishes its own bit, a process that can be modeled as a single Bernoulli trials. In the trial, when the negative provider becomes a false positive (i.e.,  $0 \rightarrow 1$ ) which occurs at probability  $\beta(t_j)$ , the Bernoulli random variable, denoted by  $X$ , takes on value 1. Otherwise, it takes the value 0. Let  $E(X)$  be the expectation of variable  $X$ , which in our case is,

$$E(X) = m(1 - \sigma_j) \cdot \beta(t_j) \quad (10)$$

We can apply the Chernoff bound for the sum of Bernoulli trials,  $Pr(X \leq (1 - \delta)E(X)) \leq e^{-\delta^2 E(X)/2}$  [32], where  $\delta > 0$  is any positive number. For term  $t_j$ , the expected success rate, denoted by  $p_p(t_j)$ , is equal to the probability of a publication success, that is,  $p_p(t_j) = Pr(fp_j > \epsilon_j)$ . Note  $fp_j = \frac{X}{X + \sigma_j m}$ , we have,

$$\begin{aligned} p_p(t_j) &= 1 - Pr(fp_j \leq \epsilon_j) \\ &= 1 - Pr\left(X \leq m \frac{\sigma_j}{\epsilon_j^{-1} - 1}\right) \\ &\geq 1 - e^{-\delta_j^2 m(1 - \sigma_j)\beta(t_j)/2} \end{aligned} \quad (11)$$

In here,  $\delta_j = 1 - \frac{1}{(\epsilon_j^{-1} - 1)(\sigma_j^{-1} - 1)} \cdot \frac{1}{\beta(t_j)} = 1 - \frac{\beta_b(t_j)}{\beta(t_j)}$ . Recall that  $\gamma$  is the required minimal success rate. If we can have

$$1 - e^{-\delta_j^2 m(1 - \sigma_j)\beta(t_j)/2} \geq \gamma \quad (12)$$

for all indexed terms, then  $\forall j, p_p(t_j) \geq \gamma$ . This means in the case of large number of terms, the percentage of successfully published terms or  $p_p$  is expected to be larger than or equal to  $\gamma$ , i.e.,  $p_p \geq \gamma$ , which is the proposition. Hence, by plugging  $\delta_j$  in Equation 12, we can derive,

$$(\beta_c(t_j))^2 - 2\left(\beta_b(t_j) + \frac{\ln \frac{1}{1-\gamma}}{(1 - \sigma_j)m}\right)\beta_c(t_j) + (\beta_b(t_j))^2 \geq 0$$

Note  $\frac{\ln \frac{1}{1-\gamma}}{(1 - \sigma_j)m} = G_j$ , and  $\beta_c(t_j)$  should be bigger than  $\beta_b(t_j)$  since success ratio is larger than 50%. Solving the inequality and taking only the solution that satisfies  $\beta_c(t_j) > \beta_b(t_j)$ , we have,

$$\beta_c(t_j) \geq \beta_b(t_j) + G_j + \sqrt{G_j^2 + 2\beta_b(t_j)G_j} \quad \blacksquare$$

## B. Proof of Theorem 4.1

*Proof:* Recoverability can be trivially proved based on the fact that  $\sum_{\forall i \in [0, c-1]} s(i, j) = v_j$ .

To prove secrecy, we examine the process of generating super-shares  $s(i, j)$ . It is easy to see that the SecSumShare protocol uses a  $(c, c)$  secret sharing to split each private input  $M(i, j)$ . The generated  $c$  shares for each input value are distributed to  $c$  different output super-shares. For each private input  $M(i, j)$ , an output super share  $s(i, j)$  has included *one and only one* share from it. Therefore, when an adversary knows at most  $c-1$  outputs, at least one share of each private input is still unknown to her. This leaves the value of any input completely undetermined to this adversary, thus the secret or the sum of input values completely undetermined.  $\blacksquare$

## Appendix B

### Analysis of Conventional PPIs

We analyze the privacy of existing PPI work and compare it with that of  $\epsilon$ -PPI. Here, we consider the primary attack and the common-term attack. Before that, we briefly introduce the construction protocol of existing PPI. To be consistent with terminology, we use term to refer to owner's identity in this section, for example, the common-identity attack is referred to as the common-term attack.

*Grouping PPI:* Inspired by  $k$ -anonymity [26], existing PPI work [7], [8], [9] constructs its index by using a grouping approach. The idea is to assign the providers into disjoint privacy groups, so that true positive providers are mixed with the false positives in the same group and are made indistinguishable. Then, a group reports binary value 1 on a term  $t_j$  as long as there is at least one provider in this group who possesses the term. For example, consider terms are distributed in a raw matrix  $M$  as in Figure 2. If providers  $p_2$  and  $p_3$  are assigned to the same group, say  $g_1$ , then in the published PPI group  $g_1$  would report to have term  $t_0$  and  $t_2$  but not  $t_1$ , because both  $p_2$  and  $p_3$  do not have term  $t_1$ .

a) *Privacy under primary attack:* To form privacy groups, existing PPIs randomly assign providers to groups. By this means, the false positive rate resulted in the PPI varies non-deterministically. Furthermore, grouping based approach is fundamentally difficult to achieve per-term privacy degree. Because different terms share the same group assignment, even if one can tune grouping strategy (instead of doing it randomly) to meet privacy requirement for one or few terms, it would be extremely hard, if not impossible, to meet the privacy requirement for thousands of terms. For primary attack, the privacy leakage depends on the false positive rate of row at term  $t_j$  in PPI  $M'$ . This way, the grouping based PPI can at best provide a privacy level at NOGUARANTEE for primary attacks. Our experiments in Section V-A1 confirms our analysis as well.

b) *Privacy under common-term attack:* The grouping based PPI work may disclose the truthful term-to-provider distribution and thus the identity of common terms. We use a specific example to demonstrate this vulnerability.

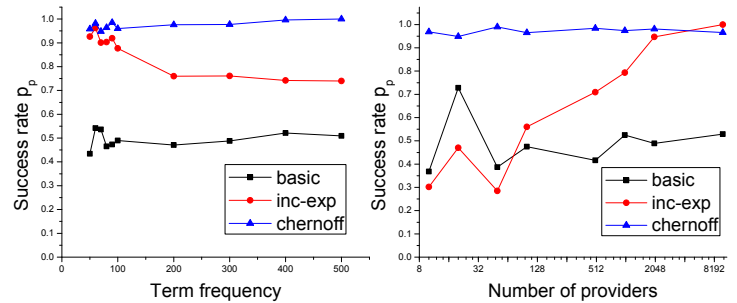
**Example** In an extreme scenario, one common term is with 100% frequency and all other terms show up in only one provider. For group assignment, as long as there are more than two groups, the rare terms can only show up in one group. In this case, the only common term in  $M'$  is the true one in  $M$ , in spite of the grouping strategy. This allows the attacker to be able to identify the true common terms in  $M$  and mount an attack against it with 100% confidence.

Given information of term distribution, one can fully exploit the vulnerability to amount common-term attacks. And the privacy degree depends on availability of term distribution information. For certain existing PPI [9], it directly leaks the sensitive common term's frequency  $\sigma_j$  to providers during index construction, leading to a NOPROTECT privacy level. Other PPI work, which does not leak exact term distribution information, still suffers from data-dependent privacy protection, resulting in a NOGUARANTEE privacy level.

## Appendix C Extra Experiment Results

### A. Effectiveness of different $\beta$ -calculation policies

We evaluate the effectiveness of three  $\beta$ -calculation policies with  $\epsilon$ -PPI, and the result shows the advantages of Chernoff bound-based policy in meeting desired privacy requirements. In the experiments, we have tested various parameter settings. We show representative results with the following values:  $\Delta = 0.02$  in incremented expectation-based policy and expected success rate  $\gamma = 0.9$  in the Chernoff bound based policy. The default false positive rate is set at  $\epsilon = 0.5$ . The experiment results are reported in Figure 6; we consider success rate as the metric. In Figure 6a, we vary identity frequency from near 0 to about 500 providers with the number of providers fixed at 10,000, and in Figure 6b we vary the number of providers with the identity frequency held constant at 0.1. It can be seen from the results that while Chernoff bound-based policy (with  $\gamma = 0.9$ ) always achieves near-optimal success rate (i.e., close to 1.0), the other two policies fall short in certain situations; the expectation-based policy is not configurable and constantly has its success rate to be around 0.5. This is expected because the expectation-based approach works on an average sense. For the incremented expectation-based policy, its success ratio, though approaching 1.0 in some cases, is unsatisfactory for common identities with high frequency (as in Figure 6a) and in the relatively small network of few providers (as in Figure 6b). On the other hand, the high-level privacy preservation of the Chernoff bound policy comes with reasonable extra search overhead. Due to space limit, this part of experiments can be found in technical report [18].



(a) Varying frequency under 10,000 providers (b) Varying provider numbers under frequency 0.1

Fig. 6: Quality of privacy preservation