

# EsPRESSo: Efficient Privacy-Preserving Evaluation of Sample Set Similarity

Carlo Blundo<sup>1</sup>, Emiliano De Cristofaro<sup>2</sup>, and Paolo Gasti<sup>3</sup>

<sup>1</sup> Università di Salerno, Italy <sup>2</sup> PARC <sup>3</sup> UC Irvine

## Abstract

In today’s digital society, electronic information is increasingly shared among different entities, and decisions are made based on common attributes. To address associated privacy concerns, the research community has begun to develop cryptographic techniques for controlled (privacy-preserving) information sharing. One interesting open problem involves two mutually distrustful parties that need to assess the *similarity* of their information sets, but cannot disclose their actual content. This paper presents the first efficient and provably-secure construction for privacy-preserving evaluation of sample set similarity, measured as the *Jaccard* similarity index. We present two protocols: the first securely computes the Jaccard index of two sets, the second approximates it, using MinHash techniques, with lower costs. We show that our novel protocols are attractive in many compelling applications, including document similarity, biometric authentication, genetic tests, multimedia file similarity. Finally, we demonstrate that our constructions are appreciably more efficient than prior work.

## 1 Introduction

The availability of electronic information is increasingly essential to the functioning of our communities and, in numerous circumstances, data needs to be shared between parties without complete mutual trust. This naturally raises commensurate privacy concerns with respect to the disclosure, and long-term safety, of sensitive contents. One interesting problem occurs whenever two or more entities need to evaluate the similarity of their information, but are reluctant to disclose data wholesale. This task faces three main technical challenges: (1) how to identify a meaningful metric to estimate similarity, (2) how to compute a measure of it such that no private information is revealed during the process, and (3) how to do so efficiently. We denote this problem as *Privacy-preserving Evaluation of Sample Set Similarity* and we motivate it below, vis-à-vis a few relevant applications.

**Document similarity:** Two parties need to estimate the similarity of their documents (or collections thereof) – in many settings, documents contain sensitive information and parties may be unwilling, or simply forbidden, to reveal their content. For instance, program chairs of a conference may want to verify that none of submitted papers is also under review in other conferences or journals, but, naturally, they are not allowed to disclose papers in submission. Likewise, two law enforcement authorities (e.g., the FBI and the local police), or two investigation teams with different clearance levels, might need to share documents pertaining suspect terrorists, but they can do so only conditioned upon a clear indication that content is relevant to the same investigation.

**Iris Matching:** Biometric identification and authentication are increasingly used due to fast and inexpensive devices that can extract biometric information from a multitude of sources, e.g., voice, fingerprints, iris, and so on. Clearly, given its utmost sensitivity, biometric data must be protected from arbitrary disclosure. Consider, for instance, an agency that needs to determine whether a given biometric appears on a government watch-list. As agencies may have different clearance levels, privacy of biometric’s owner needs to be preserved if no matches are found, but, at the same time, unrestricted access to the watch-list cannot be granted.

**Genetic Paternity/Ancestry Testing:** Advances in DNA sequencing technologies will soon yield ubiquitous and low-cost full sequencing of human genomes [16]. This will stimulate the deployment of algorithms that perform, *in computation*, various genomic tests, such as genetic *paternity* and *ancestry* testing. Since individuals tied by parent-child or ancestry relationships carry almost identical genomes, an attractive technique to establish such ties is to measure the genomes’ similarity. However, it is well-known that (human) genomic information is extremely sensitive, as a

genome not only uniquely identifies an individual, but it also reveals information about, e.g., ethnic heritage, disease predispositions, and many other phenotypic traits [21].

**Multimedia File Similarity:** Digital media, e.g., images, audio, video, are increasingly relevant in today’s computing ecosystems. Consider two parties that wish to evaluate similarity of their media files, e.g., for plagiarism detection: sensitivity of possibly unreleased material (or copyright issues) may prevent parties from revealing actual contents.

All examples above exhibit some common features: neither party can reveal its information in its entirety. What they are willing to reveal is limited to a metric that assesses their similarity. This paper presents the design of efficient cryptographic techniques for privacy-preserving evaluation of sample set similarity. Such techniques do not only appeal to examples above, but to any setting where parties need to evaluate similarity of sets, independently of their nature. It is relevant to a wide spectrum of applications, for instance, in the context of privacy-preserving sharing of information and/or recommender systems, e.g., to privately assess similarity of social network profiles, social interests, network traces, attack logs, healthcare information, and so on.

## 1.1 Technical Roadmap & Contributions

Our first step is to identify a *metric* for effectively evaluating similarity of sample sets. While several measures have been proposed, such as, cosine similarity (for vectors), Hamming and Levenshtein distances (for strings), Euclidean, Manhattan, or Minkowski similarity (for geometric objects), Pearson coefficient (for statistical data), we focus on a well-known, widely used, measure: the *Jaccard Similarity Index* [25]. It quantifies the similarity between two sets  $A$  and  $B$ , and is a rational number between 0 and 1. Experimental and analytical results indicate that high values of the Jaccard index capture well the informal notion of “roughly the same” [10] and can be used, e.g., to find near duplicate records [55] and similar documents [10], for web-page clustering [51], data mining [52], and genetics [18, 44]. As sample sets can be relatively large, in distributed settings, an approximation of the index is oftentimes preferred to its exact calculation. To this end, *MinHash* techniques [10] are used to estimate the Jaccard index, with remarkably lower computation and communication costs (see Sec. 2.1).

In this paper, we define and instantiate a cryptographic primitive geared for *privacy-preserving evaluation of sample set similarity*. We design two efficient protocols, allowing two interacting parties to compute (resp., approximate) the Jaccard index of their private sets, without reciprocally disclosing any information about their content (or, at most, an upper bound on their size). Our main cryptographic building block is Private Set Intersection Cardinality (PSI-CA) [22], introduced in Sec. 2.2. Specifically, we use PSI-CA to privately compute the magnitude of set intersection and union, and derive the value of the Jaccard index. As fast (linear-complexity) PSI-CA protocols become available, this can be done efficiently, even on large sets. Nonetheless, our work shows that, using MinHash approximations, one can obtain an estimate of the Jaccard index with remarkably increased efficiency – i.e., reducing the size of input sets, thus, the number of (costly) cryptographic operations.

Privacy-preserving evaluation of sample set similarity is appealing in many real-world scenarios. We focus on document similarity, iris matching, and multimedia file similarity, and show that privacy is attainable with low overhead. Our experiments demonstrate that our generic technique – while not bounded to any specific application – is appreciably more efficient than state-of-the-art protocols that only focus on one specific scenario, while maintaining comparable accuracy. Finally, in the process of reviewing related work, we identify limits and flaws of some prior results.

**Organization.** The rest of this paper is structured as follows. Next section introduces building blocks, then Sec. 3 presents our construction for secure computation of Jaccard index and an even more efficient technique to (privately) approximate it. Then, Sec. 4, 5, and 6 present our constructions for privacy-preserving similarity evaluation of, respectively, documents, irises, and multimedia content. The paper concludes in Sec. 7. (Also, Appendix D sketches a very fast protocol to privately approximate set intersection cardinality, that additionally hides set sizes.)

## 2 Preliminaries

### 2.1 Jaccard Similarity Index and MinHash Techniques

**Jaccard Index.** One of the most common metrics for assessing the similarity of two sets (hence, of data they represent) is the Jaccard index [25]. It measures the similarity between two sets  $A$  and  $B$  as  $J(A, B) = |A \cap B| / |A \cup B|$ . High values of the index suggest that two sets are very similar, whereas, low values indicate that  $A$  and  $B$  are almost disjoint.

The Jaccard index of  $A$  and  $B$  can be rewritten as a mere function of the *intersection*:  $J(A, B) = |A \cap B| / (|A| + |B| - |A \cap B|)$ .

**MinHash Techniques.** Computing the Jaccard index incurs a complexity linear in set sizes. Thus, in the context of a large number of big sets, its computation might be relatively expensive. In fact, for each pair of sets, the Jaccard index must be computed from scratch, i.e., no information used to calculate  $J(A, B)$  can be re-used for  $J(A, C)$ . (That is, similarity is not a transitive measure.) As a result, an *approximation* of the Jaccard index is often preferred, as it can be obtained at a significantly lower cost, e.g., using MinHash techniques [10]. Informally, MinHash techniques extract a small representation  $h_k(S)$  of a set  $S$  through deterministic (salted) sampling. This representation has a constant size  $k$ , independent from  $|S|$ , and can be used to compute an approximation of the Jaccard index. The parameter  $k$  also defines the expected error with respect to the exact Jaccard index. Intuitively, larger values of  $k$  yield smaller approximation errors. The computation of  $h_k(S)$  also incurs a complexity linear in set sizes, however, it must be performed *only once* per set, for *any* number of comparisons. Thus, with MinHash techniques, evaluating the similarity of any two sets requires only a constant number of comparisons. Similarly, the bandwidth used by two interacting parties to approximate the Jaccard index of their respective sets is also constant ( $O(k)$ ).

There are two strategies to realize MinHashes: one employs multiple hash functions, while the other is built from a single hash function.<sup>1</sup>

**MinHash with many hash functions.** Let  $\mathcal{F}$  be a family of hash functions that map items from set  $U$  to distinct  $\tau$ -bit integers. Select  $k$  different functions  $h^{(1)}(\cdot), \dots, h^{(k)}(\cdot)$  from  $\mathcal{F}$ ; for any set  $S \subseteq U$ , let  $h_{min}^{(i)}(S)$  be the item  $s \in S$  with the smallest value  $h^{(i)}(s)$ . The MinHash representation  $h_k(S)$  of set  $S$  is a vector  $h_k(S) = \{h_{min}^{(i)}(S)\}_{i=1}^k$ . The Jaccard index  $J(A, B)$  is estimated by counting the number of indexes  $i$ -s, such that that  $h_{min}^{(i)}(A) = h_{min}^{(i)}(B)$ , and this value is then divided by  $k$ . Observe that it holds that  $h_{min}^{(i)}(A) = h_{min}^{(i)}(B)$  exactly when the minimum hash value of  $A \cup B$  lies in  $A \cap B$ .

This measure can be obtained by computing the cardinality of the intersection of  $h_k(A)$  and  $h_k(B)$ , in the following way. Each element  $a_i$  of the vector  $h_k(A)$  is encoded as  $\langle a_i, i \rangle$ . Similarly,  $\langle b_i, i \rangle$  represents the  $i$ -th element of vector  $h_k(B)$ . An unbiased estimate of the Jaccard index between  $A$  and  $B$  is given by:

$$sim(A, B) = \frac{|\{\langle a_i, i \rangle\}_{i=1}^k \cap \{\langle b_i, i \rangle\}_{i=1}^k|}{k} \quad (1)$$

As discussed in [11], if  $\mathcal{F}$  is a family of min-wise independent hash functions, then each value of a fixed set  $A$  has the same probability to be the element with the smallest hash value. Specifically, for each min-wise independent hash function  $h^{(i)}(\cdot)$  and for any set  $S$ , we have that, for any  $s_j, s_l \in S$ ,  $\Pr[s_j = h_{min}^{(i)}(S)] = \Pr[s_l = h_{min}^{(i)}(S)]$ . Thus, we also obtain that  $\Pr[h_{min}^{(i)}(A) = h_{min}^{(i)}(B)] = J(A, B)$ . In other words, if  $r$  is a random variable that is 1 when  $h_{min}^{(i)}(A) = h_{min}^{(i)}(B)$  and 0 otherwise, then  $r$  is an unbiased estimator of  $J(A, B)$ ; however, in order to reduce its variance, such random variable must be sampled several times, i.e.,  $k \gg 1$  hash values must be used. In particular, by Chernoff bounds [13], the expected error of this estimate is  $O(1/\sqrt{k})$ .

**Approximating (Jaccard) Similarity of Vectors without MinHash.** If one needs to approximate the Jaccard index of two fixed-length *vectors* (rather than sets), one could use other (more efficient) techniques similar to MinHash. Observe that a vector  $\vec{S}$  can be represented as a set  $S = \{\langle s_i, i \rangle\}$ , where  $s_i$  is simply the  $i$ -th element of  $\vec{S}$ . We now discuss an efficient strategy to approximate the Jaccard index of two vectors  $A = \{\langle a_i, i \rangle\}_{i=1}^n$ ,  $B = \{\langle b_i, i \rangle\}_{i=1}^n$  of length  $n$ , without using MinHash. The approach discussed here incurs constant ( $O(k)$ ) computational and communication complexity, i.e., it is independent from vectors' length of the vectors being compared. First, select  $k$  random values  $(r_1, \dots, r_k)$ , for  $r_i$  uniformly distributed in  $[1, n]$ , and compute  $A_k = \{\langle a_{r_i}, r_i \rangle\}_{i=1}^k$  and  $B_k = \{\langle b_{r_i}, r_i \rangle\}_{i=1}^k$ , respectively. The value  $\delta = |A_k \cap B_k|/k$  can then be used to assess the similarity of  $A$  and  $B$ . We argue that  $\delta$  is an unbiased estimate of  $J(A, B)$ : for each  $\alpha \in (A_k \cup B_k)$  we have that  $\Pr[\alpha \in (A_k \cap B_k)] = \Pr[\alpha \in (A \cap B)]$  since  $\alpha \in (A \cap B) \wedge \alpha \in (A_k \cup B_k) \Leftrightarrow \alpha \in (A_k \cap B_k)$ . We also have  $\Pr[\alpha \in (A \cap B)] = J(A, B)$ , thus,  $\delta$  is indeed an unbiased estimate of  $J(A, B)$ .

The above algorithm implements a perfect min-wise permutation for this setting: since elements  $(r_1, \dots, r_k)$  are

<sup>1</sup>This paper focuses on the former technique, thus, we defer the description of the latter to Appendix A, which also overviews possible MinHash instantiations.

uniformly distributed, for each  $i \in [1, k]$  any element in  $A$  and  $B$  has the same probability of being selected. As such, similar to MinHash with many hash function, the expected error is also  $O(1/\sqrt{k})$ .

## 2.2 Cryptography Background

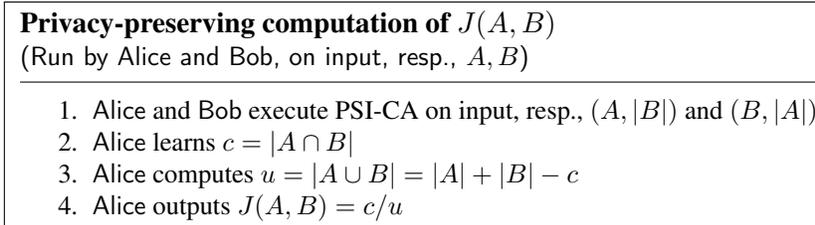
**Private Set Intersection Cardinality (PSI-CA)** is a cryptographic protocol involving two parties: Alice, on input  $A = \{a_1, \dots, a_w\}$ , and Bob, on input  $B = \{b_1, \dots, b_v\}$ , such that Alice outputs  $|A \cap B|$ , while Bob has no output. In the last few years, several PSI-CA protocols [22, 54, 32, 17] have been proposed, that are secure in different security models and under different assumptions. We choose the protocol in [17] as it achieves communication and computation complexity linear in set sizes. As a result, throughout this paper, we use the PSI-CA construction from [17], which is secure, in the presence of semi-honest adversaries, under the One-More-DH assumption [7] in the Random Oracle Model (ROM). It requires  $O(|A| + |B|)$  *offline* and  $O(|A|)$  *online* modular exponentiations in  $\mathbb{Z}_p$  with exponents from subgroup  $\mathbb{Z}_q$ . (Offline operations are computed only once, for any number of interactions and any number of interacting parties). Communication overhead amounts to  $O(|A|)$  elements in  $\mathbb{Z}_p$  and  $O(|B|)$  – in  $\mathbb{Z}_q$ . (Assuming 80-bit security parameter,  $|q| = 160$  bits and  $|p| = 1024$  bits.) PSI-CA from [17] is reviewed in Appendix B (Figure 7).

**Adversarial Model:** We use standard security models for secure two-party computation, which assume the adversary to be either semi-honest or malicious.<sup>2</sup> As per definitions in [23], protocols secure in the presence of *semi-honest adversary* assume that parties faithfully follow all protocol specifications and do not misrepresent any information related to their inputs, e.g., size and content. However, during or after protocol execution, any party might (passively) attempt to infer additional information about other party’s input. Whereas, security in the presence of *malicious parties* allows arbitrary deviations from the protocol. Security arguments in this paper are made with respect to *semi-honest* participants.

## 3 Privacy-preserving Sample Set Similarity

### 3.1 Private Computation of Jaccard Index

We now present our first construction for privacy-preserving computation of Jaccard Index. We consider two parties, Alice and Bob, on input sets  $A$  and  $B$ , respectively. We show that parties can efficiently compute  $J(A, B)$  in a privacy-preserving manner using protocol illustrated in Figure 1 below.



**Figure 1:** Proposed protocol for privacy-preserving computation of set similarity.

**Complexity.** The cost of protocol in Figure 1 is dominated by that incurred by the underlying PSI-CA protocol. As we select the PSI-CA construction of [17], which incurs linear communication and computational complexities, overall complexities of protocol in Figure 1 are also linear in the size of sets. If we were to compute the Jaccard index without privacy, asymptotic complexities would be same as our privacy-preserving protocol – i.e., linear. However, given the lack of cryptographic operations, constants hidden by the big  $O(\cdot)$  notation would be much smaller.

**Security.** Our main security claim is that, by running the protocol in Figure 1, parties do not reciprocally disclose the content of their private sets, but only learn similarity computed as the Jaccard index and the size of the the other party’s input (hence, the cardinality of set intersection and union). Therefore, security of protocol in Figure 1 relies on that of the underlying PSI-CA instantiation. In particular, Alice and Bob do not exchange any information besides messages related to the PSI-CA protocol. For this reason, a secure implementation of the underlying PSI-CA guarantees that neither Alice nor Bob learn additional information about the other party’s set (thus, we omit detailed formal proofs to ease presentation).

<sup>2</sup>Hereafter, the term *adversary* refers to protocol participants. Outside adversaries are not considered, since their actions can be mitigated via standard network security techniques.

**Performance Evaluation.** Our technique for secure Jaccard index computation can be instantiated using any PSI-CA construction. Nonetheless, to maximize efficiency, we choose the one in [17]. To assess the practicality of resulting construction, protocol in Figure 1 has been implemented in C (with OpenSSL and GMP libraries), using 160-bit random exponents and 1024-bit moduli to obtain 80-bit security.<sup>3</sup> We assume  $|A| = |B| = 1000$  and set items to be hashed using SHA-1, thus, they are 160-bit. In this setting, protocol in Figure 1 would incur (i) about 0.5s total computation time on a single Intel Xeon E5420 core running at 2.50GHz and (ii) 276KB in bandwidth. We omit running times for larger sets since, as complexities are linear, one can easily derive a meaningful estimate of time/bandwidth for any size.

We also implement an optimized prototype that further improves total running time by (1) pipelining computation and transmission and (2) parallelizing computation on two cores. We test the prototype by running Alice and Bob on two PCs equipped with 2 quad-core Intel Xeon E5420 processors running at 2.50GHz, however, we always use (at most) 2 cores. On a conservative stance, we do not allow parties to perform any pre-computation offline. We simulate a 9Mbps link, since, according to [41], it represents the current average Internet bandwidth in US and Canada.

In this setting, and again considering  $|A| = |B| = 1000$ , total running time of protocol in Figure 1 amounts to 0.23s. Whereas, the computation of Jaccard index *without privacy* takes 0.018s. Therefore, we conclude that privacy protection, in our experiments, only introduces a 12-fold slowdown, independently from set sizes.

**Comparison to prior work.** Performance evaluation above does not include any prior solutions, since, to the best of our knowledge, there is no comparable cryptographic primitive for privacy-preserving computation of the Jaccard index. The work in [48] is somewhat related: it targets private computation of the Jaccard index using Private Equality Testing (PET) [34] and deterministic encryption, however, it introduces the need for a non-colluding semi-honest *third party*, which violates our design model. Also, it incurs an impractical number of public-key operations, i.e., *quadratic* in the size of sample sets (as opposed to linear in our case). Finally, additional (only vaguely) related techniques include: (i) work on privately approximating dot product of two vectors, such as, [45, 30], and (ii) probabilistic/approximated private set operations based on Bloom filters, such as, [30, 31]. (None of this techniques, however, can be used to solve problems considered in this paper.)

### 3.2 Private Estimation of Jaccard Index based on MinHash

The computation of the Jaccard index, with or without privacy, can be relatively expensive when (1) sample sets are very large, or (2) each set must be compared with a large number of other sets. Thus, MinHash techniques, introduced in Sec. 2.1, are often used to estimate the Jaccard index, trading off an expected error with appreciably faster computation. We now show how to privately approximate the similarity of two sample sets combining MinHash and PSI-CA. Our construction is general and does not assume any specific instantiation, given that it is a multi-hash MinHash.

Recall, from Sec. 2.1, that Jaccard index can also be approximated as  $sim(A, B) = |\{\langle a_i, i \rangle\}_{i=1}^k \cap \{\langle b_i, i \rangle\}_{i=1}^k| / k$ , where  $a_i = h_{min}^{(i)}(A)$  and  $b_i = h_{min}^{(i)}(B)$ . Therefore, privacy-preserving estimation of the Jaccard index, using multi-hash MinHash, can be reduced to securely computing cardinality of set intersection above. The resulting protocol is presented in Figure 2 below.

<p><b>Private Jaccard index estimation</b> <math>sim(A, B)</math> (Run by Alice and Bob, on input, resp., <math>A, B</math>)</p> <ol style="list-style-type: none"> <li>1. Alice and Bob compute, <math>\{\langle a_i, i \rangle\}_{i=1}^k</math> and <math>\{\langle b_i, i \rangle\}_{i=1}^k</math>, resp., using multi-hash MinHash</li> <li>2. Alice and Bob execute PSI-CA on input, resp., <math>(\{\langle a_i, i \rangle\}_{i=1}^k, k)</math> and <math>(\{\langle b_i, i \rangle\}_{i=1}^k, k)</math></li> <li>3. Alice learns <math>\delta =  \{\langle a_i, i \rangle\}_{i=1}^k \cap \{\langle b_i, i \rangle\}_{i=1}^k </math></li> <li>4. Alice outputs <math>sim(A, B) = \delta / k</math></li> </ol>
---

**Figure 2:** Proposed protocol for privacy-preserving approximation of set similarity.

It is easy to observe that, compared to the Jaccard index computation (Sec. 3), the use of MinHash leads to executing PSI-CA on smaller sets, as it holds  $k \ll \min(|A|, |B|)$ . Thus, communication and computation overhead depends on  $k$ , since inputs to PSI-CA are now sets of  $k$  items, independently from the size of  $A$  and  $B$ .

**Security and Extensions.** Similar to the protocol in Sec. 3.1, the security of protocol in Figure 2 relies on the security of the underlying PSI-CA construction.

<sup>3</sup>Source-code implementation of all protocols and experiments will be released along with the final version of the paper.

In this protocol, the Alice learns some additional information compared to the protocol in Sec. 3.1. In particular, rather than computing the similarity – and therefore the size of the intersection – of sets  $A$  and  $B$ , she determines how many elements from a particular subset of  $A$  (constructed using minhash) also appear in the subset selected from  $B$ .

To fix this issue, Alice and Bob can construct their input sets (Step 1 in Figure 2) using a set of OPRFs rather than a set of hash functions: Alice and Bob engage in a multi-party protocol where Alice inputs her set  $A = \{a_1, \dots, a_v\}$  and learns a random permutation of  $\text{OPRF}_{key_j}(a_1), \dots, \text{OPRF}_{key_j}(a_v)$  for random keys  $key_j$ ,  $1 \leq j \leq k$ . Alice constructs her input selecting the smallest value  $\text{OPRF}_{key_j}(a_i)$  for each  $j$ . Bob constructs his input without interacting with Alice. While the cost of this protocol is linear in the size of the input sets, it is significantly higher than that of protocol Sec. 3.1.

**Performance Evaluation.** We also tested the performance of our construction for privacy-preserving approximation of Jaccard similarity, again using the PSI-CA from [17]. We used the same setting of Sec. 3.1, i.e., we selected sets with 1000 items, 1024-bit moduli and 160-bit random exponents, and ran experiments on two PCs with 2.5GHz CPU and a 9Mbps link. Selecting  $k = 400$ , thus, bounding the error to 5%, the total running time of protocol in Figure 2 amounts to 0.09s – less than half compared to the one in Figure 1. Whereas, in the same setting, the approximation of Jaccard index *without privacy* takes 0.007s. Thus, the slow-down factor introduced by the privacy-protecting layer (similar to the protocol proposed in Sec. 3.1) is 12-fold. Again, note that times for different set sizes can be easily estimated since the complexity of the protocol is linear.

**Prior Work.** The estimation of set similarity through MinHash – whether privacy-preserving or not – requires counting the number of times for which it holds that  $h_{min}^{(i)}(A) = h_{min}^{(i)}(B)$ , with  $i = 1, \dots, k$ . We have denoted this number as  $\delta$ . Protocol in Figure 2 above attains secure computation of  $\delta$  through privacy-preserving set intersection cardinality. However, it appears somewhat more intuitive to do so by using the approach proposed by [4] in the context of social-network friends discovery. Specifically, in [4], Alice and Bob compute, resp.,  $\{a_i\}_{i=1}^k$  and  $\{b_i\}_{i=1}^k$ , just like in our protocol. Then, Alice generates a public-private keypair  $(pk, sk)$  for Paillier’s additively homomorphic encryption cryptosystem [43] and sends Bob  $\{z_i = \text{Enc}_{pk}(a_i)\}_{i=1}^k$ . Bob computes  $\{(z_i \cdot \text{Enc}_{pk}(-b_i))^{r_i}\}_{i=1}^k$  for random  $r_i$ ’s and returns the resulting vector of ciphertexts after shuffling it. Upon decryption, Alice learns  $\delta$  by counting the number of 0’s. Nonetheless, the technique proposed by [4] actually incurs an increased complexity, compared to our protocol in Figure 2 (instantiated with PSI-CA from [17]). Assuming 80-bit security parameters, thus, 1024-bit moduli and 160-bit subgroups, and 2048-bit Paillier moduli, and using  $m$  to denote a multiplication of 1024-bit numbers, multiplications of 2048-bit numbers count for  $4m$ . Using square-and-multiply, exponentiations with  $q$ -bit exponents modulo 1024-bit count for  $(1.5|q|m)$ . In [4], Alice performs  $k$  Paillier encryptions (i.e.,  $2k$  exponentiations and  $k$  multiplications) and  $k$  decryptions (i.e.,  $k$  exponentiations and multiplications), while Bob computes  $k$  exponentiations and multiplications. Therefore, the total computation complexity amounts to  $(4 \cdot 4 \cdot 1.5 \cdot 1024 + 3 \cdot 4)km = 24,588km$ . Whereas, our approach (even without pre-computation) requires both Alice and Bob to perform  $4k$  exponentiations of 160-bit numbers modulo 1024-moduli and  $2k$  multiplications, i.e.,  $(4 \cdot 1.5 \cdot 160 + 2)km = 962km$ , thus, our protocol achieves a 25-fold efficiency improvement. Communication overhead is also higher in [4]: it amounts to  $(2 \cdot 2048)k$  bits; whereas, using PSI-CA, we need to transfer  $(2 \cdot 1024 + 160)k$  bits, i.e., slightly more than half the traffic.

## 4 Privacy-Preserving Document Similarity

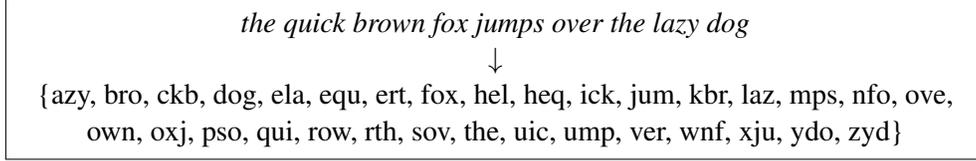
After building efficient (linear-complexity) primitives for privacy-preserving computation/approximation of Jaccard index, we now explore their applications to a few compelling problems. We start with evaluating the similarity of two documents, which is relevant in many common applications, including copyright protection, file management, plagiarism prevention, duplicate submission detection, law enforcement. In last few years, the security community has started investigating privacy-preserving techniques to enable detection of similar documents without disclosing documents’ actual contents. Below, we first review prior work and, then, present our technique for efficient privacy-preserving document similarity.

### 4.1 Related Work

The work in [28] (later extended in [39]) is the first to realize *privacy-preserving document similarity*. It realizes secure computation of the *cosine similarity* of vectors representing the documents, i.e., each document is represented as the list of words appearing in it, along with the normalized number of occurrences. Recently, Jiang and Samanthula [29] have proposed a novel technique relying on the Jaccard index and  $N$ -gram based document representation [38]. (Given

any string, an  $N$ -gram is a substring of size  $N$ ). According to [29], the  $N$ -gram based technique presents several advantages over cosine similarity: (1) it improves on finding *local similarity*, e.g., overlapping of pieces of texts, (2) it is language-independent, (3) it requires a much simpler representation, and (4) it is less sensitive to document modification. We overview it below.

**Documents as sets of  $N$ -grams.** A document can be represented as a set of  $N$ -grams contained in it. To obtain such a representation, one needs to remove spaces and punctuation and build the set of successive  $N$ -grams in the document. An example of a sentence, along with its  $N$ -gram representation (for  $N = 3$ ), is illustrated in Figure 3. The similarity of two documents can then be estimated as the *Jaccard index* of the two corresponding sets of  $N$ -grams. In the context of document similarity, experts point out that 3 results as a good choice of  $N$  [10].



**Figure 3:** Tri-gram representation.

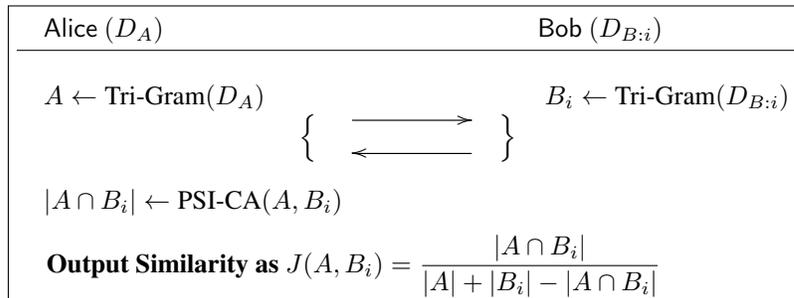
To enable privacy-preserving computation of Jaccard index, and therefore estimation of document similarity, Jiang and Samanthula [29] propose a two-stage protocol based on *Paillier's* additively homomorphic encryption [43]. Suppose Alice wants to privately evaluate the similarity of her document  $D_A$  against a list of  $n$  documents held by Bob, i.e.,  $D_{B:1}, \dots, D_{B:n}$ . First, Bob generates a global space,  $|S|$ , of tri-grams based on his document collection. This way,  $D_A$  as well as each of Bob's document,  $D_{B:i}$ , can be represented as binary vectors in the global space of tri-grams: each component is 1 if the corresponding tri-gram is included in the document and 0 otherwise. We denote with  $A$  the representation of  $D_A$  and with  $B_i$  that of  $D_{B:i}$ . Then, Alice and Bob respectively compute random shares  $a$  and  $b_i$  such that  $a + b_i = |A \cap B_i|$ . Next, they set  $c = |A| - a$  and  $d_i = |B_i| - b_i$ . Finally, Alice and Bob, on input  $(a, c)$  and  $(b_i, d_i)$ , resp., execute a Secure Division protocol (e.g., [12, 1]) to obtain  $(a + b_i)/(c + d_i) = |A \cap B_i|/|A \cup B_i| = J(A, B_i)$ .

The computational complexity of the protocol in [29] amounts to  $O(|S|)$  Paillier encryptions performed by Alice, and  $O(n \cdot |S|)$  modular multiplications – by Bob. Whereas, communication overhead amounts to  $O(n \cdot |S|)$  Paillier ciphertexts.

**Flaw in [29].** Unfortunately, protocol in [29] *is not secure*, since Bob has to disclose his global space of tri-grams (i.e., the set of all tri-grams appearing in his document collection). Therefore, Alice can passively check whether or not a word appears in Bob's document collection. Actually, Alice can learn much more, as we show in Appendix C. We argue that this flaw could be fixed by considering the global space of tri-grams as the set of all possible tri-grams, thus, avoiding the disclosure of Bob's tri-grams set. Assuming that documents are stripped of any symbol and contain only lower-cased letters and digits, we obtain  $S = \{a, b, \dots, z, 0, 1, \dots, 9\}^3$ . Unfortunately, this modification would tremendously increase computation and communication overhead.

## 4.2 Our Construction

As discussed in Sec. 3, we can realize privacy-preserving computation of the Jaccard index using PSI-CA. To privately evaluate the similarity of documents  $D_A$  and any document  $D_{B:i}$ , Alice and Bob execute protocol in Figure 4. Function  $\text{Tri-Gram}(\cdot)$  denotes the representation of a document as the set of tri-grams appearing in it.



**Figure 4:** Privacy-preserving *evaluation* of document similarity of documents  $D_A$  and  $D_{B:i}$ .

**Complexity.** Complexity of protocol in Figure 4 is bounded by that of the underlying PSI-CA construction. Using the technique in [17], computational complexity amounts to  $O(|A| + |B_i|)$  modular exponentiations, whereas, communication overhead – to  $O(|A| + |B_i|)$ . Observe that, in the setting where Alice holds one documents and Bob a collection of  $n$  documents, complexities should be amended to  $O(n|A| + \sum_{i=1}^n |B_i|)$ . However, due to the nature of protocol in [17], Bob can perform  $O(\sum_{i=1}^n |B_i|)$  computation *off-line*, ahead of time. Hence, total *online* computation amounts to  $O(n|A|)$ .

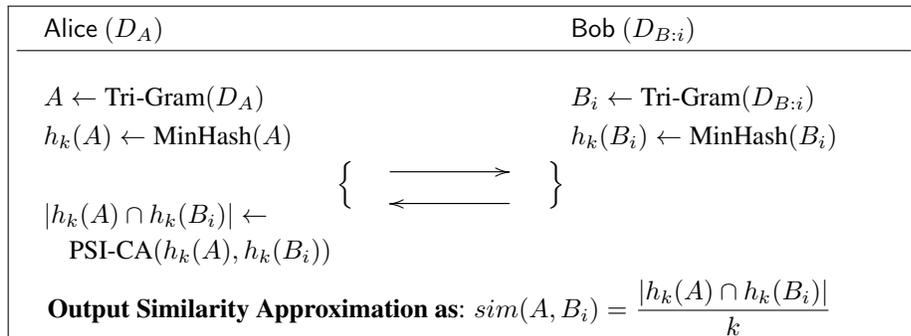
**More efficient computation using MinHash.** As discussed in Sec. 2.1, one can approximate the Jaccard index by using MinHash techniques, thus, trading off accuracy with significant improvement in protocol complexity. The resulting construction is similar to the one presented above and is illustrated in Figure 5. It adds an intermediate step between the tri-gram representation and the execution of PSI-CA: Alice and Bob apply MinHash to sets  $A$  and  $B_i$ , respectively, and obtain  $h_k(A)$  and  $h_k(B_i)$ . The main advantage results from the fact that PSI-CA is now executed on smaller sets, of constant size  $k$ , thus, achieving significantly improved communication and computational complexities. Again, note that the error is bounded by  $O(1/\sqrt{k})$ .

**Performance Evaluation.** We now compare the performance of our constructions to the most efficient prior technique, i.e., the protocol in [29] (that, unfortunately, is insecure). We consider the setting of [29], where Bob maintains a collection of  $n$  documents. Recall that our constructions use the PSI-CA in [17]. Assuming 80-bit security parameters, we select 1024-bit moduli and 160-bit random exponents. As [29] relies on Paillier encryption, it uses 2048-bit moduli and 1024-bit exponents. In the following, let  $m$  denote a multiplication of 1024-bit numbers. Multiplications of 2048-bit numbers count for  $4m$ . Modular exponentiations with  $q$ -bit exponents modulo 1024-bit count for  $(1.5|q|)m$ . The protocol in [29] requires  $O(|S|)$  Paillier encryptions and  $O(n \cdot |S|)$  modular multiplications. As pointed above, we need  $|S| = 36^3 = 46,656$ . Therefore, the total complexity amounts to  $(4 \cdot 1.5 \cdot 1024 + 4n)|S|m = (6144 + 4n)|S|m \approx (2.9 \cdot 10^8 + 1.9 \cdot 10^5 n)m$ .

Our construction above requires  $(2 \cdot 1.5 \cdot 160n|A|)m$  for the computation of Jaccard index similarity and  $(1.5 \cdot 160nk)m$  for its approximation. Thus, to compare performance of our protocol to that of [29], we need to take into account the dimensions of  $A$ ,  $B_i$ , as well as  $n$  and  $k$ . To this end, we collected 393 scientific papers from the KDDcup dataset of scientific papers published in ArXiv between 1996 and 2003 [15]. The average number of different tri-grams appearing in each paper is 1307. Therefore, cost of our two techniques can be estimated as  $(2 \cdot 1.5 \cdot 160 \cdot 1307n)m$  and  $(1.5 \cdot 160 \cdot nk)m$ , respectively. Thus, our technique for privacy-preserving document similarity is faster than [29] for  $n < 2000$ . Furthermore, using MinHash techniques, complexity is always faster (and of at least one order of magnitude), using both  $k = 40$  and  $k = 100$ . Also, recall that, as opposed to ours, the protocol in [29] is *not* secure.

Assuming that it takes about  $1\mu s$  to perform modular multiplications of 1024-bit integers (as per our experiments on a single Intel Xeon E5420 core running at 2.50GHz), we report estimated running times in Table 1 for increasing values of  $n$  (i.e., the number of Bob’s documents).

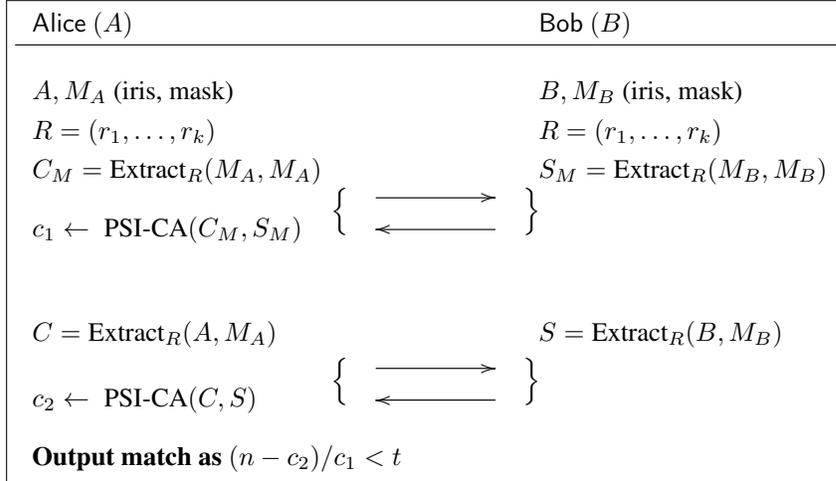
We performed some statistical analysis to determine the real magnitude of the error introduced by MinHash, when compared to the Jaccard index without MinHash. Our analysis is based on the trigrams from documents in the KDDcup dataset [15], and confirms that the average error is within the expected bounds: for  $k = 40$ , we obtained an average error of 14%, while for  $k = 100$  the average error was 9%. This is acceptable, considering that the Jaccard index actually provides a normalized *estimate* of the similarity between two sets, not a definite metric.



**Figure 5:** Privacy-preserving *approximation* of document similarity of documents  $D_A$  and  $D_{B:i}$ .

$n$	[29]	Figure 4	Figure 5	
			$k = 100$	$k = 40$
10	9.5 mins	6.3 secs	0.05 secs	0.05 secs
$10^2$	9.9 mins	63 secs	1.9 secs	1.9 secs
$10^3$	12.7 mins	10.4 mins	48 secs	19.2 secs
$10^4$	40.7 mins	1.74 hours	8 mins	3.2 mins
$10^5$	5.3 hours	17.4 hours	1.2 hours	32 mins

**Table 1:** Computation time of privacy-preserving document similarity.



**Figure 6:** Privacy-preserving iris matching of biometric  $A$  and  $B$ .

## 5 Privacy-Preserving Iris Matching

Advances in biometric recognition enable the use of biometric data as a practical mean of authentication and identification. Today, several governmental agencies around the world perform large-scale collections of different biometric features. As an example, the US Department of Homeland Security (DHS) collects face, fingerprint and iris images, from visitors within its US-VISIT program [53]. Iris images are also collected from all foreigners, by the United Arab Emirates (UAE) Ministry of Interior, as well as fingerprints and photographs from certain types of travelers [24].

While biometry serves as an excellent mechanism for identification of individuals, biometric data is, undeniably, extremely sensitive and must be subject to minimal exposure. As a result, such data cannot be disclosed arbitrarily. Nonetheless, there are many legitimate scenarios where biometric data should be shared, in a controlled way, between different entities. For instance, an agency may need to determine whether a given biometric appears on a government watch-list. As agencies may have different clearance levels, privacy of biometric’s owner should be preserved if no matches are found, but, at the same time, unrestricted access to the watch-list cannot be granted.

### 5.1 Related Work

As biometric identification techniques are increasingly employed, related privacy concerns have been investigated by the research community. A number of recent results address the problem of privacy-preserving face recognition. The work in [19] is the first to present a secure protocol, based on Eigenfaces, later improved by [46]. Next, [42] designs a new privacy-preserving face recognition algorithm, called SCiFI. Furthermore, the protocol in [6] realizes privacy-preserving fingerprint identification, using FingerCodes [26]. FingerCodes use texture information from a fingerprint to compare two biometrics. The algorithm is not as discriminative as traditional fingerprint matching techniques based on location of minutiae points, but it is adopted in [6] given its suitability for efficient *privacy-preserving* realization. Among all biometric techniques, this paper focuses on *iris-based* identification. The problem of privacy-preserving iris matching has been introduced by Blanton and Gasti in [8], who propose an approach based on a combination of garbled circuits [56] and homomorphic encryption.

### 5.2 Our Construction

A (human) iris can be digitalized as an  $n$ -bit string  $S = s_1 s_2 \dots s_n$  with an  $n$ -bit mask  $M_S = m_{s_1} m_{s_2} \dots m_{s_n}$ . The mask indicates which bits of  $S$  have been read reliably. In particular, the  $i$ -th bit of  $S$  should be used for matching

only if the  $i$ -th bit of  $M_S$  is set to 1. A common value for  $n$ , which we use in our experiments, is 2048. As, during an iris scan, the subject may rotate its head, a right or left shift can occur in the iris representation, depending on the direction of the rotation. Therefore, the distance between two irises  $A$  and  $B$  is computed as the minimum distance between all rotations of  $A$  and the representation of  $B$ . In practice, it is reasonable to assume that the rotation is limited to a shift of at most 5 positions towards left/right [8].

The matching between two irises,  $A$  and  $B$ , is computed via the *Weighted Hamming Distance (WHD)* of the samples. Let  $M = (M_A \wedge M_B)$ ;  $WHD$  is computed as:

$$WHD(A, B, M) = \frac{HD(A \wedge M, B \wedge M)}{\|M\|} \quad (2)$$

where  $\|\cdot\|$  denotes hamming weight, i.e. the number of string bits set to 1. Given a threshold  $t$ , if  $WHD(A, B, M) < t$ , we say that irises  $A$  and  $B$  are *matching*. (Assuming a maximum rotation of 5 positions, the distance must be computed 11 times.)

In the following, we propose a probabilistic technique for privately estimating of  $WHD(A, B, M)$ , that relies on the construction for privacy-preserving estimation of Jaccard index based on MinHash (introduced in Sec. 3.2). The error on the approximation is bounded by the MinHash parameter  $k$ .

Proposed protocol is illustrated in Figure 6. Given any two  $n$ -bit strings  $X$  and  $Y$  and any list of  $k$  values  $R = (r_1, \dots, r_k)$ , with  $r_i \in [1, n]$ , we define:

$$\text{Extract}_R(X, Y) = \{w_{r_1}, \dots, w_{r_k}\}, \text{ where } w_{r_i} = \begin{cases} \langle x_{r_i}, r_i \rangle & \text{if } y_{r_i} = 1 \\ r \leftarrow \{0, 1\}^\tau & \text{otherwise} \end{cases}$$

Given  $A, M_A, B, M_B$ , Alice and Bob privately determine  $WHD(A, B, M_A \wedge M_B)$ :

- Alice and Bob negotiate  $k$  random values  $R = (r_1, \dots, r_k)$ , with  $r_i \in [1, n]$ .
- Alice computes  $C_M = \text{Extract}_R(M_A, M_A)$ ; Bob computes  $S_M = \text{Extract}_R(M_B, M_B)$ .
- Alice and Bob engage in a PSI-CA protocol where their inputs are  $C_M$  and  $S_M$  and Alice learns the output  $c_1$  of PSI-CA.
- Alice computes  $C = \text{Extract}_R(A, M_A)$ ; similarly, Bob computes  $S = \text{Extract}_R(B, M_B)$ .
- Alice and Bob interact in a PSI-CA protocol with input  $C$  and  $S$  respectively; at the end of the protocol, Alice learns  $c_2$ , i.e. the cardinality of the intersection.
- Biometric  $A$  matches  $B$  iff  $(n - c_2)/c_1 < t$ .

### 5.3 Comparison to prior work

We now compare our technique for privacy-preserving iris matching to prior work, namely the technique in [8]. First, observe that protocol in Figure 6 estimates the Weighted Hamming Distance with bounded error, whereas, construction in [8] yields its exact computation. However, as we discuss below, the error incurred by our technique is low enough to be used in practice and achieves reduced computational complexity. In fact, our probabilistic protocol could be used to perform a fast, preliminary test: if differences between two irises are significant, then there is no need for further tests. Otherwise, the two parties can engage in the protocol in [8] to obtain (in a privacy-preserving way) a precise result. Next, as opposed to the technique in [8], Alice also learns an estimate on the number of bits set to 1 in the combined mask  $M_A \wedge M_B$ , but not their position. However, this information is not sensitive, thus, it does not leak any information about the iris sampled by Alice or Bob.

**Optimization.** As discussed above, it is reasonable to assume that “Bob” (e.g., DHS) holds a database with a large number of biometric samples, whereas, “Alice” (e.g., TSA) has only one or few samples that she is searching in Bob’s database. To this end, we now show how the protocol in Figure 6 can be optimized, by pre-computing several expensive operations offline, for such a scenario.

Note that Bob can perform the offline phase of PSI-CA protocol of [17] (see Figure 7) on all bits of his biometric samples: *unlike the protocol in [8], this is required only once, independently on the number of interactions between Bob and any user.*

Protocol in Figure 6				Protocol in [8]	
		Offline	Online	Offline	Online
Bob	± 5-bit rot.	0.13 ms + 5.8 s/rec	71.5 ms/rec	2.8 s + 71.55 ms/rec	97.2 ms + 134.28 ms/rec
	no rot.	0.13 ms + 530 ms/rec	6.5 ms/rec	2.6 s + 6.48 ms/rec	97.2 ms + 12.33 ms/rec
Alice	± 5-bit rot.	71.63 ms	71.5 ms/rec	12.2 s + 3 ms/rec	20.34 ms/rec
	no rot.	6.63 ms	6.5 ms/rec	11.7 s + 0.27 ms/rec	1.8 ms/rec

**Table 2:** Computation overhead of our randomized iris matching technique in Figure 6 and that of [8]. Experiments are performed with 5-bit left/right rotation and with no rotation of the iris sample. “Rot” abbreviates “rotation” and “rec” – “record”.

After negotiating with Bob the values  $R = (r_1, \dots, r_k)$ , and before receiving her input, Alice pre-computes  $k$  pairs  $\langle \alpha_{0,i} = H(\langle 0, r_i \rangle)^{R_c}, \alpha_{1,i} = H(\langle 1, r_i \rangle)^{R_c} \rangle$ . (This assumes the use of the PSI-CA in [17].) Once Alice’s mask has been revealed, she constructs the corresponding encrypted representation by simply selecting the appropriate element from each pair. Similarly, she computes  $k$  triples  $(\alpha_{0,i}, \alpha_{1,i}, \alpha_{\rho,i})$  where  $\alpha_{0,i}, \alpha_{1,i}, \alpha_{\rho,i}$  represent 0, 1 and a random element in  $\{0, 1\}^\tau$ , respectively. Alice later uses such triples to represent each bit  $\beta_i$  of her iris sample as  $\alpha_i = \alpha_{\beta,i}$  if the corresponding bit in the mask is 1, else, as  $\alpha_i = \alpha_{\rho,i}$ . During the online phase, Alice selects the appropriate pre-computed values to match the mask and the iris bits. Similarly, Bob inputs the selected bits of each record’s mask and iris into the PSI-CA protocol.

**Performance Comparison.** In Table 2, we report running times from implementations of, respectively, our protocol in Figure 6 and technique in [8]. We assume that about 75% of the bits in the mask are set to 1 (like in [8]). We set the length of each iris and mask to 2048 bits and the database size to 320 irises, which is the number used in prior work. All tests are run on a single Intel Xeon E5420 core running at 2.50GHz. We set  $k = 25$ , thus, obtaining an expected error in the order of  $1/\sqrt{25}$ , i.e., 20%.

Observe that *online* cost incurred by Bob with our technique is significantly lower compared to that of protocol in [8]. Whereas, it is higher for Alice. Nonetheless, summing up the computation overhead incurred by both Alice and Bob, our protocol always results faster than the one in [8] for the online computation.

The *offline* cost imposed on Bob is about twice as high as its counterpart in protocol from [8]. However, in our protocol, the offline part is done once, for all possible interactions, independently from their number. Whereas, in [8], the offline computation needs to be performed anew, for *every* interaction. In settings where Bob interacts frequently with multiple entities, this may significantly effect protocol’s overall efficiency. Furthermore, the offline cost imposed on Alice is markedly lower (several orders of magnitudes) using our technique.

We conclude that the protocol in Figure 6 improves, in many settings, overall efficiency compared to state of the art. However, it introduces a maximum error of about 20%, whereas, the scheme in [8] compute the exact – rather than approximate – outcome of an iris comparison. Thus, a good practice is to use the scheme in Figure 6 to perform an initial selection of relevant biometric samples, using a threshold  $t' > t$ , in order to compensate for the error. The final matching on selected samples can then be done, in a privacy-preserving manner, using the protocol in [8].

## 6 Privacy-Preserving Multimedia File Similarity

Amid widespread availability of digital cameras, digital audio recorders, and media-enabled smartphones, users generate a staggering amount of *multimedia* content. As a result, secure online storage (and management) of large volumes of multimedia data becomes increasingly desirable. According to [57], YouTube received more than 13 million hours of video in 2010, and 48 hours are uploaded every minute (i.e., 8 years of content each day). On a similar note, Flickr users upload about 60 photos every second.

Such an enormous amount of multimedia data requires automated analysis tools, e.g., Content-Based Image Retrieval (CBIR) [49], that is used in the context of images. There are several available techniques to implement CBIR, including search techniques based on color histograms [50], bin similarity coefficients [40], texture for image characterization [37], shape features [47], edge directions [27], and matching of shape components such as corners, line segments or circular arcs [14].

**Related Work.** Motivated by the potential sensitivity of multimedia data, the research community has begun to develop mechanisms for secure signal processing. For instance, authors in [20] are the first to investigate secure signal processing related to multimedia documents. Then, the work in [35, 36] introduces two protocols to search over encrypted multimedia databases. Specifically, it extracts 256 visual features from each image. Then, files are encrypted in

		Offline	Online
Bob	Exact	0.13 ms + 33.28 ms/record	33.28 ms/record
	Approximate	0.13 ms + 13 ms/record	13 ms/rec
Alice	Exact	66.69 ms	33.28 ms/record
	Approximate	26.13 ms	13 ms/record

**Table 3:** Computation cost of our multimedia documents similarity protocol.

a distance-preserving fashion, so that encrypted features can be directly compared for similarity evaluation. Similarity is computed using the Jaccard index between the visual features of searched image and those of images in a database. However, the security of the scheme relies on order-preserving encryption (used to mask frequencies of recurring visual features), which is known to provide only a limited level of security [9].

**Our Approach.** We use the Jaccard index to assess the similarity of multimedia files. As showed in Sec. 3, we can do so, in a privacy-preserving way, using protocol in Figure 1. Observe that our approach is independent of the underlying feature extraction algorithm, even though protocol accuracy naturally relies on the quality of the feature extraction phase. Once features have been extracted, our privacy-preserving protocols only reveal their similarity, thus, without disclosing the features themselves. As an example, we instantiate our techniques for privacy-preserving *image* similarity. Our approach for feature extraction is based on [36], since its accuracy is reasonable enough for real-world use, using color histograms in the color space of Hue, Saturation and Value (HSV). Thus, our scheme achieves the same accuracy of [36], in terms of precision and recall. Once again, to obtain improved efficiency, similarity can be approximated using Min-Hash techniques, as per protocol in Figure 2. In this case, both performance and accuracy depend on the MinHash parameter  $k$ .

**Performance Evaluation.** We test our technique with the same dataset used by [36], i.e., 1000 images from the standard Corel dataset. We extract 256 features from each image, for a total of 256,000 features for the whole database. We envision a user, Alice, willing to assess similarity of an image against an image database, held by Bob. We run our protocol for privately computing the Jaccard index (“Exact” rows in Table 3) and for estimated similarity, using MinHash with  $k = 100$  (“Approximate” row). Table 3 summarizes our experiments. All tests are run on a single Intel Xeon E5420 core running at 2.50GHz and show that privacy protection is attainable at a very limited cost.

Remark that a thorough performance comparison between our protocol and related work is out of the scope of this paper, since the main effort of prior work has been achieving high accuracy in similarity detection, rather improving efficiency. Thus, we defer it to future work. Nonetheless, the authors of [36] report that the running time of their protocol is in the order of 1 second per image, on a hardware comparable to our testbed (a dual-core 3GHz PC with 4GB of RAM). Therefore, it is safe to assume that our protocol for privacy-preserving multimedia file similarity is about one order of magnitude faster than available techniques, even without considering pre-computation.

## 7 Conclusion

This paper introduced the first efficient construction for privacy-preserving evaluation of sample set similarity, relying on the Jaccard index measure. We also presented an efficient randomized protocol that approximates, with bounded error, this similarity index. Our techniques are generic and practical enough to be used as a basic building block for a wide array of different privacy-preserving functionalities, including document and multimedia file similarity, biometric matching, genomic testing, similarity of social profiles, and so on. Experimental analyses support our efficiency claims and demonstrate improvements over prior results. Source-code implementation of all proposed protocols and experiments is available upon request and will be released along with the final version of the paper.

Naturally, our work does not end here: additional applications and extensions require further investigation. Also, as part of future work, we plan to investigate privacy-preserving computation of other similarity measures.

## References

- [1] M. J. Atallah, M. Bykova, J. Li, K. B. Frikken, and M. Topkara. Private collaborative forecasting and benchmarking. In *WPES*, 2004.
- [2] G. Ateniese, M. Blanton, and J. Kirsch. Secret handshakes with dynamic and fuzzy matching. In *NDSS*, 2007.
- [3] G. Ateniese, E. De Cristofaro, and G. Tsudik. (If) size matters: Size-Hiding Private Set Intersection. In *PKC*, 2011.

- [4] E. Baglioni, L. Becchetti, L. Bergamini, U. Colesanti, L. Filippini, A. Vitaletti, and G. Persiano. A lightweight privacy preserving sms-based recommendation system for mobile users. In *RecSys*, 2010.
- [5] P. Baldi, R. Baronio, E. De Cristofaro, P. Gasti, and G. Tsudik. Countering gattaca: efficient and secure testing of fully-sequenced human genomes. In *CCS*, 2011.
- [6] M. Barni, T. Bianchi, D. Catalano, M. Di Raimondo, R. Labati, P. Failla, D. Fiore, R. Lazzeretti, V. Piuri, F. Scotti, and A. Piva. Privacy-preserving fingerprint authentication. In *MM&Sec*, 2010.
- [7] M. Bellare, C. Namprempe, D. Pointcheval, and M. Semanko. The one-more-RSA-inversion problems and the security of Chaum’s blind signature scheme. *Journal of Cryptology*, 2003.
- [8] M. Blanton and P. Gasti. Secure and Efficient Protocols for Iris and Fingerprint Identification. In *ESORICS*, 2011.
- [9] A. Boldyreva, N. Chenette, Y. Lee, and A. O’Neill. Order-preserving symmetric encryption. In *EUROCRYPT*, 2009.
- [10] A. Broder. On the resemblance and containment of documents. In *Compression and Complexity of Sequences*, 1997.
- [11] A. Broder, M. Charikar, A. Frieze, and M. Mitzenmacher. Min-wise independent permutations. In *STOC*, 1998.
- [12] P. Bunn and R. Ostrovsky. Secure two-party k-means clustering. In *CCS*, 2007.
- [13] H. Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics*, 1952.
- [14] S. D. Cohen and L. J. Guibas. Shape-based image retrieval using geometric hashing. In *ARPA Image Understanding Workshop*, 1997.
- [15] Cornell Univ. KDDCUP Dataset. <http://www.cs.cornell.edu/projects/kddcup/datasets.html>.
- [16] K. Davies. *The \$1,000 genome: the revolution in DNA sequencing and the new era of Personalized Medicine*. Free Press, 2010.
- [17] E. De Cristofaro, P. Gasti, and G. Tsudik. Fast and Private Computation of Set Intersection Cardinality. <http://eprint.iacr.org/2011/141>, 2012.
- [18] P. Dombek, L. Johnson, S. Zimmerley, and M. Sadowsky. Use of repetitive DNA sequences and the PCR to differentiate *Escherichia coli* isolates from human and animal sources. *Applied and Environmental Microbiology*, 66(6), 2000.
- [19] Z. Erkin, M. Franz, J. Guajardo, S. Katzenbeisser, I. Lagendijk, and T. Toft. Privacy-preserving face recognition. In *PETS*, 2009.
- [20] Z. Erkin, A. Piva, S. Katzenbeisser, R. L. Lagendijk, J. Shokrollahi, G. Neven, and M. Barni. Protection and retrieval of encrypted multimedia content: When cryptography meets signal processing. *EURASIP J. Information Security*, 2007, 2007.
- [21] J. Fowler, J. Settle, and N. Christakis. Correlated genotypes in friendship networks. *Proceedings of the National Academy of Sciences*, 108(5), 2011.
- [22] M. Freedman, K. Nissim, and B. Pinkas. Efficient private matching and set intersection. In *EUROCRYPT*, 2004.
- [23] O. Goldreich. *Foundations of cryptography*. Cambridge Univ Press, 2004.
- [24] IrisGuard Press Release. <http://cl.ly/3KIB>.
- [25] P. Jaccard. Etude comparative de la distribution florale dans une portion des Alpes et du Jura, 1901.
- [26] A. Jain, S. Prabhakar, L. Hong, and S. Pankanti. Filterbank-based fingerprint matching. *IEEE Transactions on Image Processing*, 9(5), 2000.
- [27] A. K. Jain and A. Vailaya. Image retrieval using color and shape. *Pattern Recognition*, 1996.
- [28] W. Jiang, M. Murugesan, C. Clifton, and L. Si. Similar document detection with limited information disclosure. In *ICDE*, 2008.
- [29] W. Jiang and B. K. Samanthula. N-gram based secure similar document detection. In *DBSec 2011*, Lecture Notes in Computer Science. Springer, 2011.
- [30] M. Kantarcioglu, R. Nix, and J. Vaidya. An efficient approximate protocol for privacy-preserving association rule mining. In *KDD*, 2009.
- [31] F. Kerschbaum. Outsourced Private Set Intersection Using Homomorphic Encryption. In *AsiaCCS*, 2012.
- [32] L. Kissner and D. Song. Privacy-preserving set operations. In *CRYPTO*, 2005.

- [33] M. Li, N. Cao, S. Yu, and W. Lou. Findu: Privacy-preserving personal profile matching in mobile social networks. In *INFOCOM*, 2011.
- [34] H. Lipmaa. Verifiable Homomorphic Oblivious Transfer and Private Equality Test. In *ASIACRYPT*, 2003.
- [35] W. Lu, A. Swaminathan, A. L. Varna, and M. Wu. Enabling search over encrypted multimedia databases. In *Media Forensics and Security*, 2009.
- [36] W. Lu, A. L. Varna, A. Swaminathan, and M. Wu. Secure image retrieval through feature protection. In *ICASSP*, 2009.
- [37] W.-Y. Ma and B. S. Manjunath. Texture features and learning similarity. In *CVPR*, 1996.
- [38] U. Manber. Finding similar files in a large file system. In *USENIX*, 1994.
- [39] M. Murugesan, W. Jiang, C. Clifton, L. Si, and J. Vaidya. Efficient privacy-preserving similar document detection. *The VLDB Journal*, 19, August 2010.
- [40] W. Niblack, R. Barber, W. Equitz, M. Flickner, E. H. Glasman, D. Petkovic, P. Yanker, C. Faloutsos, and G. Taubin. The QBIC Project: Querying Images by Content, Using Color, Texture, and Shape. In *SPIE*, 1993.
- [41] Ookla Net Metrics. Canada and US Source Data. <http://www.netindex.com/source-data/>, 2011.
- [42] M. Osadchy, B. Pinkas, A. Jarrous, and B. Moskovich. SCiFi – A system for secure face identification. In *IEEE S&P*, 2010.
- [43] P. Paillier. Public-key cryptosystems based on composite degree residuosity classes. In *EUROCRYPT*, 1999.
- [44] M. Popescu, J. Keller, and J. Mitchell. Fuzzy measures on the gene ontology for gene product similarity. *Transactions on Computational Biology and Bioinformatics*, 2006.
- [45] P. Ravikumar, W. Cohen, and S. Fienberg. A secure protocol for computing string distance metrics. In *PSDM*, 2004.
- [46] A.-R. Sadeghi, T. Schneider, and I. Wehrenberg. Efficient privacy-preserving face recognition. In *ICISC*, 2009.
- [47] S. Sclaroff. Deformable prototypes for encoding shape categories in image databases. *Pattern Recognition*, 30(4), 1997.
- [48] M. Singh, P. Krishna, and A. Saxena. A privacy preserving jaccard similarity function for mining encrypted data. In *TENCON*, 2009.
- [49] A. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain. Content-based image retrieval at the end of the early years. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(12), 2000.
- [50] J. R. Smith and S.-F. Chang. Tools and techniques for color image retrieval. In *SPIE*, 1996.
- [51] A. Strehl, J. Ghosh, and R. Mooney. Impact of similarity measures on web-page clustering. In *AAAI*, 2000.
- [52] P. Tan, M. Steinbach, V. Kumar, et al. *Introduction to data mining*. Pearson, 2006.
- [53] U.S. Department of Homeland Security. <http://www.dhs.gov/files/programs/usv.shtm>.
- [54] J. Vaidya and C. Clifton. Secure set intersection cardinality with application to association rule mining. *Journal of Computer Security*, 13(4), 2005.
- [55] C. Xiao, W. Wang, X. Lin, and J. Yu. Efficient similarity joins for near duplicate detection. In *WWW*, 2008.
- [56] A. Yao. How to generate and exchange secrets. In *FOCS*, 1986.
- [57] Youtube press release. [http://www.youtube.com/t/press\\_statistics](http://www.youtube.com/t/press_statistics).
- [58] Y. Zhang, W. Liu, and W. Lou. Anonymous communications in mobile ad hoc networks. In *INFOCOM*, volume 3, 2005.

## A Additional Details on MinHash Techniques

**Single-Hash MinHashes.** Besides the multiple-hash technique presented in Sec. 2, another approach for approximating the Jaccard index using MinHash employs a single hash function. In this case, rather than selecting one value per hash function, one selects the  $k$  values from set  $A$  that hash the to smallest integers. Specifically, let  $h(\cdot)$  be a hash function, and  $k$  a fixed parameter; for any set  $S$ , define  $h_k(S) \subset S$  as the set of the pre-images of the  $k$  smallest hash values of elements of  $S$ . Consider:

$$\text{sim}(A, B) = \frac{|(h_k(h_k(A) \cup h_k(B))) \cap (h_k(A) \cap h_k(B))|}{|h_k(h_k(A) \cup h_k(B))|} \quad (3)$$

It holds that  $\text{sim}(A, B)$  is an unbiased estimate of the Jaccard index of  $A$  and  $B$ . Again, by standard Chernoff bounds [13], the expected error is  $O(1/\sqrt{k})$ .

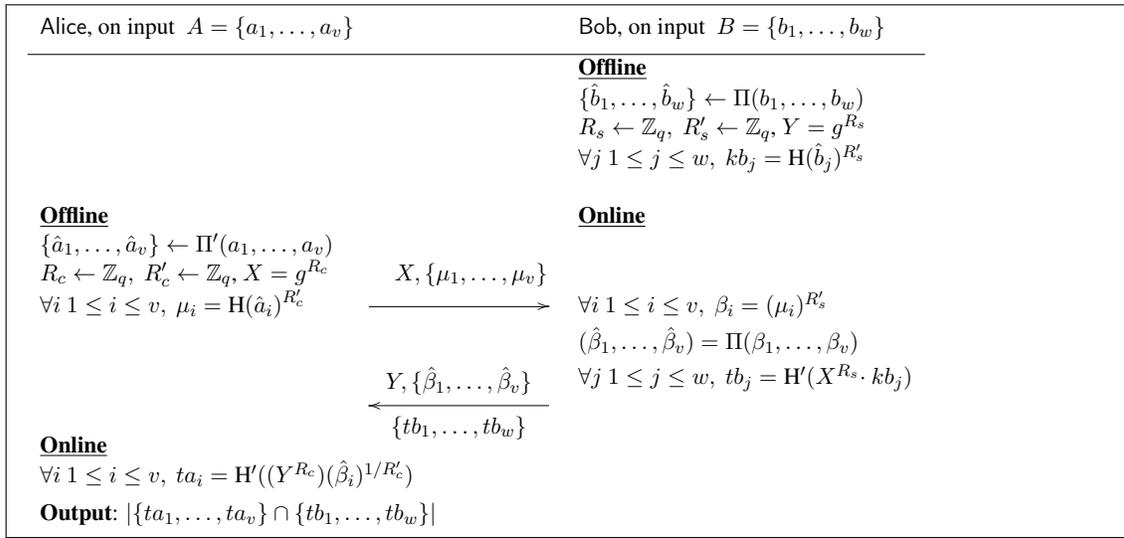
**MinHash Instantiations.** In order to implement the MinHash schemes described in this paper, the hash function should be defined by a random permutation over the set  $A \cup B$ . Assuming  $m = |A \cup B|$ , then one would need  $\Omega(m \log m)$  bits to specify a truly random permutation, thus, yielding an infeasible overhead even for small values of  $m$ . Broder, et al. [11] point out that one can obviate to this problem by using *Min-wise Independent Permutation* (MWIP) families rather than random permutations. Using MWIPs, for any subset of the domain, any element is equally likely to be the minimum, but the number of bits to specify such a permutation is showed to be still be relatively large, i.e.,  $\Omega(m)$ . In practice, however, one can allow certain relaxations. To this end, [11] introduces *approximate* MWIPs, by accepting a small error  $\varepsilon$ . The authors require all items in a set  $S$  to have only a (almost equal) chance to become the minimum element of  $A$ 's image under the permutation. Thus, for any approximate MWIP, implemented using  $h_{min}^{(i)}(\cdot)$  as defined above, for an expected relative error  $\varepsilon$ , it holds:

$$\left| \Pr \left[ h_{min}^{(i)}(S) = s \right] - \frac{1}{|S|} \right| \leq \frac{\varepsilon}{|S|}$$

A class of permutations often used in practice is one based on *linear transformations*. It is assumed that the universe is  $\mathbb{Z}_p$  for some prime  $p$  and the family of permutation is constructed using a hash function computed as  $h(x) = ax + b \pmod p$ , where  $(a, b) \in (\mathbb{Z}_p^*, \mathbb{Z}_p)$ . Such a linear transformations are easy to represent and efficiently calculable.

## B PSI-CA Protocol from [17]

In Figure 7, we review the Private Set Intersection Cardinality (PSI-CA) protocol presented in [17]. Notation is slightly modified to ease presentation.



**Figure 7:** PSI-CA protocol from [17]. It executes on common input of two primes  $p$  and  $q$  (such that  $q|p-1$ ), a generator  $g$  of a subgroup of size  $q$  and two hash functions,  $H$  and  $H'$ , modeled as random oracles.  $\Pi(\cdot)$  and  $\Pi'(\cdot)$  denote random permutations. All computation is mod  $p$ .

**Correctness.** It is easy to see that, for any  $a^* \in A$ , and  $b^* \in B$  s.t.  $a^* = b^*$ , then  $\exists (ta_i, tb_j)$  s.t.:

$$tb_j = H'(X^{R_s} \cdot H((b^*)^{R'_s})) = H'(g^{R_c R_s} \cdot H((b^*)^{R'_s})) = H'(Y^{R_c} \cdot H((a^*)^{R'_s})) = ta_i$$

## C Flaw in Private Document Similarity in [29]

In this section, we show that the protocol in [29] is not privacy-preserving (even in semi-honest model). In fact, Bob, in order to participate in the protocol, must disclose his global space of tri-grams. Given this information, Alice can efficiently check whether a word, e.g.,  $w$ , appears in Bob's document collection. Indeed, Alice computes  $w$ 's tri-gram based representation, then she checks whether all such tri-grams appear in Bob's public global space. If so, Alice learns that  $w$  appears in a document held by Bob with some non-zero probability. Technically, this probability is not 1 because Alice could have a *false positive*, i.e.  $w$  may not be in Bob's documents even though  $w$ 's trigrams are in Bob's public global space. On the other hand, if at least one of the tri-grams of  $w$  is not in Bob's public global space, Alice learns that

**Privacy-preserving approximation of  $|A \cap B|$** Run by Alice and Bob on input, resp.,  $A$  and  $B$ 

1. Alice and Bob compute,  $\{\langle a_i, i \rangle\}_{i=1}^k$  and  $\{\langle b_i, i \rangle\}_{i=1}^k$ , resp., using multi-hash MinHash where:  $a_i \stackrel{\text{def}}{=} h_{\min}^{(i)}(A)$  and  $b_i \stackrel{\text{def}}{=} h_{\min}^{(i)}(B)$
2. Alice and Bob execute PSI-CA on input, resp.,  $(\{\langle a_i, i \rangle\}_{i=1}^k, k)$  and  $(\{\langle b_i, i \rangle\}_{i=1}^k, k)$
3. Alice learns  $\delta = |\{\langle a_i, i \rangle\}_{i=1}^k \cap \{\langle b_i, i \rangle\}_{i=1}^k|$
4. Bob sends  $w$  to Alice
5. Alice outputs  $\delta \cdot (v + w)/(1 + \delta)$

**Figure 8:** Our technique for Approximated Private Set Intersection Cardinality.

Bob’s documents do not contain  $w$ . This, obviously, violates privacy requirements. If Alice and Bob include punctuation and spaces in their tri-grams representation of their documents, the probability of *false positive* becomes negligible. We do not exploit “relations” between consecutive meaningful words in the sentence, which could potentially (further) aggravate information leakage about Bob’s documents.

We now show yet another attack that lets Alice learn even more, since the N-grams representation *embeds* document’s structure. From the global space of tri-grams  $\mathcal{GS}$ , we can construct a directed graph  $G(V, E)$  representing relations between tri-grams in Bob’s document collection. Any path in such a graph will lead to a textual fragment contained in some document held by Bob. A vertex in the graph represents a tri-gram; whereas, an edge between two vertices implies that the two corresponding tri-grams are consecutive tri-grams in a word. Given a trigram  $x \in \mathcal{GS}$ , with  $x^{(i)}$  we denote the  $i$ -th letter in  $x$ . The directed graph  $G(V, E)$  is constructed as follows. The vertex set is  $V = \{V_x \mid x \in \mathcal{GS}\}$  and the edge set is  $E = \{\langle V_x, V_y \rangle \mid x^{(2)} = y^{(1)} \wedge x^{(3)} = y^{(2)}\}$ . A path  $V_{x_1}, \dots, V_{x_n}$  in  $G$ , will correspond to the string  $x_1^{(1)} x_1^{(2)} x_2^{(3)} x_3^{(3)} \dots x_n^{(3)}$ . Such a string (or some of its substring) appears in some document in Bob’s collection. By using algorithms based on Deep First Search visit of a graph, a vocabulary, and syntactic rules, we could extract large document’s chunks. We did not explore further other techniques to extract “information” from the global space of tri-grams as we consider them to be out of the scope of this paper.

## D Faster and Size-Hiding (Approximated) PSI-CA

Privacy-preserving computation of set intersection cardinality has been investigated quite extensively by the research community [22, 32, 54, 17], motivated by several interesting applications, including: privacy-preserving authentication and key exchange protocols [2], data and association rule mining [54], genomic applications [5], healthcare [32], policy-based information sharing [17], anonymous routing [58], and – as argued by this paper – sample set similarity.

In many of the application scenarios, however, it may be enough to obtain an estimation, rather than the exact measure, of set intersection cardinality. For instance, if PSI-CA is used to privately quantify the number of common social-network friends (e.g., to assess profile similarity) [33], then one may want to trade off a bounded accuracy loss with a significant improvement in protocol overhead (and without sacrificing the level of attained privacy protection). Naturally, such an improved construction is particularly appealing whenever participants’ input sets are very large.

Using MinHash techniques, we realize privacy-preserving estimation of set intersection cardinality with (constant) computation and communication complexities that only depend on the MinHash parameter – i.e.,  $O(k)$ . Proposed construction is illustrated in Figure 8. Observe that while we tolerate a bounded accuracy loss – again, depending on MinHash’s parameter  $k$ , i.e.,  $O(\sqrt{k})$  – our protocol achieves the same, provably-secure, privacy guarantees as if we ran PSI-CA on whole sets.

**Size-Hiding.** Another factor motivating the use of MinHash techniques for PSI-CA is related to input size secrecy. Available PSI-CA protocols always disclose, from the execution, at least an upper bound on input set sizes. Whereas, protocol in Figure 8 conceals – unconditionally – Alice’s set size, thus, achieving *Size-Hiding* Private Set Intersection Cardinality. Considering recent results motivating the need for size-hiding features in private set operations (see [3]), this additional feature is particularly valuable.