

# Private Set Intersection with Linear Communication from General Assumptions <sup>\*†</sup>

Brett Hemenway Falk<sup>1</sup>, Daniel Noble<sup>1</sup>, and Rafail Ostrovsky<sup>2</sup>

<sup>1</sup>University of Pennsylvania

<sup>2</sup>University of California, Los Angeles

October 2, 2019

## Abstract

This work presents a hashing-based algorithm for Private Set Intersection (PSI) in the honest-but-curious setting. The protocol is generic, modular and provides both asymptotic and concrete efficiency improvements over existing PSI protocols.

If each player has  $m$  elements, our scheme requires only  $O(m\lambda)$  communication between the parties, where  $\lambda$  is a security parameter.

Our protocol builds on the hashing-based PSI protocol of Pinkas et al. (USENIX 2014, USENIX 2015), but we replace one of the sub-protocols (handling the cuckoo “stash”) with a special-purpose PSI protocol that is optimized for comparing sets of unbalanced size. This brings the asymptotic communication complexity of the overall protocol down from  $\omega(m\lambda)$  to  $O(m\lambda)$ , and provides concrete performance improvements (10-15% reduction in communication costs) over Kolesnikov et al. (CCS 2016) under real-world parameter choices.

Our protocol is simple, generic and benefits from the permutation-hashing optimizations of Pinkas et al. (USENIX 2015) and the Batched, Relaxed Oblivious Pseudo Random Functions of Kolesnikov et al. (CCS 2016).

---

\*This research was sponsored in part by ONR grant (N00014-15-1-2750) “Syn-Crypt: Automated Synthesis of Cryptographic Constructions”, NSF grant CNS-1513671, DARPA/SPAWAR contract N66001-15-C-4065 and ODNI/IARPA contract 2019-1902070008. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the DoD, DARPA, ODNI, IARPA, or the U.S. Government.

†Please cite the conference proceedings version of this paper: Brett Hemenway Falk, Daniel Noble, and Rafail Ostrovsky. 2019. Private Set Intersection with Linear Communication from General Assumptions. In *18th Workshop on Privacy in the Electronic Society (WPES19)*, November 11, 2019, London, United Kingdom. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3338498.3358645>

# 1 Introduction

Private Set Intersection (PSI) is a secure Multi-Party Computation (MPC) protocol that allows two parties, who each hold a private set of elements from some universe,  $\mathcal{U}$ , to compute the intersection of their (private) sets, without revealing information about the elements outside the intersection. PSI is an important cryptographic tool, and is a building block for many more complex functionalities and thus has received a lot of attention from the cryptographic community. In this work, we focus on PSI in the honest-but-curious setting, and thus we focus on comparing our protocol to other protocols that target the honest-but-curious security model.

A simple, generic method for privately computing set intersection would be to securely compute all pairwise comparisons between the elements. If each player has  $m$  elements, this would require  $m^2$  comparisons. The number of comparisons can be reduced by first hashing the elements into bins. The players would agree on a hash function  $h : \mathcal{U} \rightarrow [n]$ , and then hash their elements into bins, where each bin is of size  $b$ . Then for each bin, the players can perform a secure comparison of all the elements. This basic hashing protocol requires  $nb^2$  secure comparisons. If  $n = O(m)$ , then with high probability, the maximum bin size is  $\log n / \log \log n$ , so this would require  $n^{(\log n / \log \log n)^2}$  secure comparisons. If, instead of performing all pair-wise comparisons in each bin, we recursed, hashing each bin into sub-bins, we obtain a protocol that requires  $O(n \log n)$  secure comparisons.

This scheme can be improved by replacing one player’s hash table with a cuckoo hash table [PSZ14, PSSZ15, KKRT16]. In this modification, the players choose two hash functions, and Alice hashes each of her elements into two bins, and Bob uses the two hash functions to hash his elements into a cuckoo-hash table with a stash [KMW09]. Then, for each location in Bob’s cuckoo hash table, the players compute a secure comparison between the single element in that bucket and every element in Alice’s corresponding bucket. Finally, they compare every element in Bob’s stash against every element in Alice’s table. Since Bob only has a single element in each of his buckets, this method only requires  $O(m)$  comparisons to compare Bob’s cuckoo table against Alice’s table. Unfortunately, comparing Bob’s stash (of size  $s$ ) against Alice’s table still requires  $O(ms\lambda)$  communication in the protocols of [PSZ14, PSSZ15, KKRT16]. The entire protocol then requires  $O(ms\lambda)$  communication, and somewhat counterintuitively, the asymptotic communication complexity of the protocol is dominated by computing the intersection with the small ( $\omega(1)$ -sized) stash.

In this work, we give a novel PSI protocol that builds on the cuckoo-hashing based approach of [PSZ14, PSSZ15, KKRT16] by incorporating an efficient protocol for unbalanced PSI (e.g. [KLS<sup>+</sup>17, CLR17, CHLR18, RA18]) to reduce the communication complexity of comparing Bob’s stash to Alice’s set. This provides both asymptotic and practical improvements in communication complexity. Our protocol reduces the communication cost of the stash-comparison step from  $\omega(m\lambda)$  to  $O(m\lambda)$ , thus reducing the communication complexity of the entire protocol from  $\omega(m\lambda)$  to  $O(m\lambda)$ . These asymptotic efficiency gains do not

rely on concrete number-theoretic hardness assumptions, but can be realized from OT alone.

In addition to the asymptotic improvements in communication complexity, this protocol provides concrete improvements in communication complexity in set sizes of about 4,000 elements. Many practical PSI applications involve sets of thousands or millions of elements. For instance the High Country Bandits case used the intersection of three cell tower dumps, of total size about 150,000, to find a phone that had been active near three different bank robberies [SFF14]. Another example are the FBI-maintained NCIC hotlists, which combined contain over 12 million records such as license plates of stolen vehicles [FBI19].<sup>1</sup> These are routinely intersected with police agencies who record tens of thousands of license plates per day. At present the agencies are given permission to store the lists locally to perform the intersection, however, since the data is restricted to authorized users, the FBI may prefer for intersections with certain lists to be computed privately.

In Section 5.5 we calculate the actual communication cost of our protocol using different implementations of the unbalanced PSI sub-protocol, and compare these with the PSI protocol of Kolesnikov et al. (CCS 16) [KKRT16].

We also observe that in the cuckoo-hashing schemes of [PSZ14, PSSZ15, KKRT16] Bob’s hashing protocol does not need to support dynamic insertions and deletions, and thus we can replace the cuckoo-hashing scheme with a 1-out-of- $k$  hashing scheme, and compute the optimal allocation of his elements in an offline pre-processing phase. Although the size of the hash table,  $n$ , remains  $n = O(m)$ , this allows us to guarantee that an optimal allocation is found. This is described in Section 4.2.

## 2 Previous work

There have been many approaches to the problem of private set intersection (PSI), and early works focused on building custom protocols to perform PSI (some examples include [FNP04, KS05, HL10, JL09, JL10, DSMRY09, DCT10, DCT12]). Over the years, many special-purpose PSI protocols have been proposed and implemented. Many of the early protocols were designed around Oblivious Polynomial Evaluation (OPE). In this framework, Alice interpolates a polynomial with roots at her elements, and then Alice and Bob work together to privately evaluate this polynomial at points given by Bob’s (private) elements. The private evaluation can be done using any additively homomorphic cryptosystem. Some PSI protocols that fall into this framework include [FNP04, HN10, KS05, DSMRY09]. These protocols typically have good communication complexity (each side sends  $O(m)$  ciphertexts), but high computational cost (each side must do  $O(m)$  public-key operations). See [DCT12, Appendix A] for an overview of many of the special-purpose PSI protocols.

---

<sup>1</sup>Hotlist contents are not public, so while we here provide an estimate on the total number of records, we do not know the sizes of individual hotlists.

Another class of PSI protocols use Oblivious Pseudo-Random Functions (OPRFs) [FIPR05]. An Oblivious PRF is a protocol where Alice holds a PRF key,  $\kappa$ , and Bob holds an input,  $x$ , and the OPRF protocol allows Bob to learn  $\text{PRF}_\kappa(x)$  while Alice learns nothing about  $x$ .

OPRFs provide a natural method for a linear-communication PSI protocol in the semi-honest model. If Alice has a set,  $X$ , and Bob has a set,  $Y$ , Alice will generate a key,  $\kappa$ , for an Oblivious PRF, and for every  $y \in Y$ , they will use the OPRF protocol to give Bob  $\text{PRF}_\kappa(y)$ . Then Alice will locally evaluate  $\text{PRF}_\kappa(x)$  for all  $x \in X$ , and send these evaluations to Bob. Bob can then locally compute the intersection by comparing his evaluations to those received from Alice.

OPRFs can be implemented generically, using secure Multiparty Computation to compute an ordinary PRF, where Alice’s private input is the PRF key, and Bob’s private input is his evaluation point. This approach was taken in [PSSW09] where they used garbled circuits to obliviously compute the AES-based PRF. There have also been many special-purpose constructions of Oblivious PRFs designed specifically for set intersection protocols. The work of [HL10] shows how to instantiate an oblivious version of the Naor-Reingold PRF (based on the DDH assumption), and make it secure against malicious adversaries. The works [DCT10, DCT12] use the one-more RSA assumption in Random Oracle Model to build an OPRF-based linear-time PSI protocol, and the work of [JL10] uses an OPRF-based PSI protocol to provide security against malicious adversaries based on the one-more gap Diffie-Hellman problem in the Random Oracle Model. The work of [JL09] builds an OPRF secure in the standard model under the Decisional Composite Residuosity assumption.

In [KLS<sup>+</sup>17] it was observed that the OPRF-based PSI protocols are well-suited to applications where the parties hold sets of unequal size. In particular, if Bob’s set  $Y$  is much smaller than Alice’s set  $X$ , then using an OPRF-based PSI protocol, they only need  $|Y|$  OPRF calls, followed by  $|X|$  communication. In particular, they show that the natural approach of using garbled circuits to implement an AES-based PRF is extremely efficient when Bob’s set is sufficiently small. The recent work of [CLR17] uses leveled fully homomorphic encryption to create a PSI protocol for unbalanced set sizes, and this work was later extended to the malicious setting [CHLR18]. The work of [RA18] uses cuckoo filters and a pairing-based public-key cryptosystem design and implements an efficient PSI protocol for unbalanced sets.

Unfortunately, when the set sizes are roughly balanced, the generic OPRF protocols that use general-purpose MPC machinery to obliviously evaluate a PRF are not as efficient as the custom-PSI protocols, whereas the custom OPRF-based PSI protocols like [DCT10, DCT12, JL10] achieve linear complexity and practical efficiency, but under strong and non-standard cryptographic assumptions.

In the face of the plethora of custom PSI protocols, [HEK12] proposed the idea, that circuit-based PSI protocols built on general-purpose MPC have many advantages, most notably that they were easy to implement and integrate into other, more complex secure computation protocols. The work of [HEK12] identified three natural PSI protocols that could easily be implemented by an off-

the-shelf MPC protocol. If the universe  $\mathcal{U}$  of elements is known in advance, and is not too large, the players can simply encode their sets as characteristic vectors in  $\{0, 1\}^{|\mathcal{U}|}$ , and then perform  $|\mathcal{U}|$  secure bit-wise AND operations to compute their intersection (called the Bit-Wise And (BWA) protocol). When the universe is not known in advance (or is too large), the players can securely perform all pairwise comparisons (this requires  $m^2$  secure comparisons to intersect two sets of size  $m$ ), this is called the Pair-Wise Comparison (PWC) protocol, and is included only as a baseline, or straw-man protocol. Finally, they introduce the Sort-Compare-Shuffle (SCS) paradigm, where each player locally sorts their sets, then they engage in a secure computation to securely merge their sorted sets (using the bitonic sorting network). After the joint multi-set is sorted, all elements in the intersection will occur twice in two adjacent positions. Thus the intersection can be computed using  $2m - 1$  secure comparisons (comparing element  $i$  and  $i + 1$  for  $i = 1, \dots, 2m - 1$ ). Finally, before the intersection can be revealed, it must be randomly shuffled (using the Waksman permutation network [Wak68]) to hide information carried by the position of the intersected elements. The bitonic sorting network<sup>2</sup> requires  $O(m \log m)$  comparisons to merge two sorted lists of length  $m$ , the Waksman permutation network requires  $O(m \log m)$  gates (each of which could be implemented using a comparison) to randomly permute  $m$  elements and thus the total number of comparisons required by the SCS approach is  $O(m \log m)$ . Another class of protocols uses hashing. If the players agree on a hash function (or hash functions), they can use the hash functions to locally sort their elements into bins, and then perform pairwise comparisons on the bins. In its simplest form, the players agree on a hash function,  $h : \mathcal{U} \rightarrow [n]$ , and some bucket size  $b$ . Then they locally hash their  $m$  elements into  $n$  buckets of size  $b$ . If any bucket receives more than  $b$  elements, the protocol will fail, so  $b$  must be chosen to be large enough so that this probability is sufficiently small. Then for each bucket, the players will engage in a secure computation to compare all elements within that bucket. If they use brute-force comparison within the bucket, this requires  $b^2$  comparisons, and the entire protocol requires  $nb^2$  comparisons to compute the intersection. If the players use the SCS method (described above) within each bucket, the number of comparisons drops to  $O(nb \log b)$ . If  $n = O(m)$ , then  $b$  must be  $O(\log m / \log \log m)$ , and the entire protocol is  $O(m \log m)$ .

The works of [PSZ14, PSSZ15, PSZ16] outline an optimization of this approach, where Alice uses a  $k$ -out-of- $k$  hashing, while Bob uses cuckoo hashing. To do this, Alice and Bob agree on  $k$  hash functions  $h_1, \dots, h_k$ , and Alice hashes each of her elements into the  $k$  buckets defined by these hash functions. Alice's buckets will be sized to store as many elements as necessary. Bob, on the other hand, uses  $h_1, \dots, h_k$  to build a cuckoo hash table (with a stash), and hashes each element into this cuckoo hash table. For each of the  $n$  bins, Alice and Bob engage in a secure computation to compare the single element in that bin in

---

<sup>2</sup>More precisely, this only uses the final “layer” of a bitonic sorting network, which merges two sorted lists using  $O(m \log(m))$  comparisons. A full bitonic sorting network sorts an arbitrarily permuted list and requires  $O(m \log^2(m))$  comparisons.

Bob’s cuckoo hash table to the  $b$  elements Alice has in her corresponding bin. Finally, they compare each element in Bob’s stash to every one of Alice’s elements. If the stash has size  $s$ , this requires  $nb + ms$  secure comparisons. Using cuckoo hashing, we can set  $n = O(m)$ ,  $s = \omega(1)$ , and as above  $b = O(\log(m))$ , and thus the protocol uses  $O(m \log m)$  secure comparisons. The protocols of [PSZ14, PSSZ15] use an OT-masking protocol to replace the  $(b + s)n$  secure comparisons to simply sending  $(1 + s)n$  pseudo random masks. This reduces the communication complexity of these protocols to  $\omega(n\lambda) = \omega(m\lambda)$ . The primary improvement introduced in [PSSZ15] is the notion of *permutation-based hashing* which reduces the complexity of each secure comparison (but not the number of secure comparisons). Permutation-based hashing can be used to improve the performance of all the hashing-based PSI protocols (including ours), and we review the details of permutation-based hashing in Section 4.1. The performance of [PSSZ15] can be further improved by viewing the OT-based solution as a special OPRF-based solution, and instantiating it with novel, special-purpose OPRFs [KKRT16]. It appears that the communication of these OT-masking protocols can be improved by the use of “silent-OT” extension [BCG<sup>+</sup>19], but this would come at a cost in terms of computation, and to the best of our knowledge this modification has never been implemented. We remark, however, that our techniques could still be applied to further improve the communication cost if the traditional OT extension were replaced by silent OT extension.

The recent work of [PRTY19] improves the general cuckoo-hashing technique of [PSZ14, PSSZ15, PSZ16], by replacing the OPRF with a *multi-point* OPRF. In [PSZ14, PSSZ15, PSZ16] Bob uses cuckoo hashing to ensure that each of his hash buckets contains at most one element. This allows them to use a “one-time” OPRF for each bucket, where Bob learns the OPRF value calculated on the single element in his bucket, while Alice learns the OPRF key, and compute the values on all elements in her bucket. In [PRTY19], Alice and Bob use a one-time multi-point OPRF, which allows Bob to hash to buckets that can hold more than one element, and this drives down both the asymptotic communication cost (to  $\Theta(m\lambda)$ ) and the concrete cost.

Bloom filters provide a natural, generic method for computing set intersections. If each participant inserts their  $m$  elements into a Bloom filter of size  $n$ , then they can use a secure bitwise-AND calculation to calculate the intersection of their Bloom filters. The players can locally query this “intersected” Bloom filter on each of their elements to find the intersection. This approach was taken in [MBD12]. It is straightforward to check that if an element appears in the intersection, it will also show up in the intersected Bloom filter. It is not too hard to see that the “intersected” Bloom filter created in this way may have extra ones that would not appear in a fresh Bloom filter created by inserting only the elements in the intersection of the two private sets. These extra ones, have the potential to leak information about the underlying sets, and thus this simple method of computing a set intersection using Bloom filters cannot be made to meet the security definitions of PSI.

Nevertheless, Bloom filters can be used to perform PSI. The protocol of [DCW13] introduces the notion of a garbled Bloom filter. In a traditional Bloom

filter, an empty set is represented by the all-zero vector. An element  $x$  is inserted by setting each of the  $k$  locations defined by the hash functions  $h_1, \dots, h_k$  to 1. In a garbled Bloom filter, each entry holds a string (rather than a single bit), and an element  $x$  is inserted by secret-sharing  $x$  using a  $k$ -out-of- $k$  secret sharing scheme, ( $x = s_1 + \dots + s_k$ ) and inserting the shares  $s_i$  into the location determined by  $h_i$ . If the slot  $h_i(x)$  is occupied, we *re-use* the existing share in that location. As long as one of the  $k$  slots is unoccupied, there will be enough freedom to make the  $s_i$  sum to  $x$ . If all  $k$  slots are occupied, then the insertion fails (just as in a regular Bloom filter). This approach can also be made secure against malicious adversaries [RR17].

The garbled Bloom filter can be made into a PSI protocol as follows. Alice will create a standard Bloom filter encoding her set, while Bob will create a garbled Bloom filter encoding his set. Denote these Bloom filters  $A \in \{0, 1\}^n$  and  $B \in \left(\{0, 1\}^\lambda\right)^n$ . Then for each entry  $i \in [n]$ , if  $A[i] = 0$ , then set  $C[i] \stackrel{\$}{\leftarrow} \{0, 1\}^\lambda$ . If  $A[i] = 1$ , then  $C[i] = B[i]$ . It is not hard to check that this is a garbled Bloom filter that encodes all elements in the intersection. The somewhat surprising result from [DCW13] is that this resulting garbled Bloom filter has exactly the same distribution as a garbled Bloom filter created by the intersection, and thus it leaks no information beyond the intersection. Implementing this PSI protocol using MPC requires  $n$  secure (single-bit) comparisons. Note that in the semi-honest setting, Alice can generate all the random values (for when  $A[i] = 0$ ) and then Bob can receive the garbled Bloom filter, compute the intersection, and send the intersection to Alice, and thus there is no need to generate the random values within the MPC protocol. For a false-positive rate  $\epsilon$ , a Bloom filter holding  $m$  elements needs to be of size  $O(m \log(1/\epsilon))$ , thus when  $\epsilon = O(m^{-1})$ ,  $n = O(m \log m)$ .

Several recent works [CO18, PSWW18, PSTY19] considered the problem of designing PSI protocols that do not reveal the intersection itself, but instead allow the players to compute *secret shares* of the intersection which could then be used in future secure computation protocols. Both these works build on the hashing-based PSI protocols of [PSZ14, PSSZ15, PSZ16]. In these hashing-based schemes, Bob uses cuckoo-hashing to hash each of his elements into one of  $k$  possible buckets, whereas Alice hashes each of her elements into all  $k$  potential buckets. Then, for each bucket, they must compare whether Bob's element matches one of Alice's elements in the corresponding bucket. Thus, for each bucket the players need a method for securely computing a *Private Set Membership* (PSM), where one party has a single element, and the other party has a large set. Prior works used OT to compute a one-time OPRF for each bucket [PSZ14, PSSZ15, PSZ16], but the hashing-based PSI protocols can thus be instantiated with any secure PSM protocol in place of the one-time OPRF. If the PSM protocol can be made to output secret shares, rather than membership status in the clear, then the entire hashing-based PSI protocol can be made to output secret-shares, rather than the intersection in the clear. A generic MPC calculation of set membership would require a number of equality tests equal to the size of the set, and would thus result in a PSI protocol that is fairly

inefficient. The primary technical contribution of [CO18] is the development of a novel, efficient PSM protocol that outputs secret shares (or encryptions) of the membership query. At a high-level, their PSM protocol works as follows: the sender constructs a binary tree, each node of which contains a key for a symmetric key cryptosystem. If the tree has depth  $t$  (i.e., the sender’s elements are bit-strings of length  $t$ ), then the sender and receiver engage in  $t$   $\binom{2}{1}$ -OTs, where the receiver’s inputs are the bits of her element. This allows the receiver to traverse the tree obliviously, and learn a single value, corresponding to whether her element was in the sender’s set. Using OT-extension, this PSM protocol has comparable efficiency to the one-time OPRF-based schemes of [PSZ14, PSSZ15, PSZ16]. Asymptotically, however, the protocol still requires  $\omega(m\lambda)$  communication.

The work of [PSWW18] takes a different approach, and instead modifies the hashing structure, introducing the notion of two-dimensional cuckoo hashing. In the basic cuckoo-hashing scheme, Bob maps each of his elements into one of  $k$  buckets, and Alice maps her elements to  $k$ -out-of- $k$  buckets, thus ensuring that if there is an overlap in their sets, it will result in an overlap in exactly one of the buckets. The reason this approach is not amenable to a generic circuit-based protocol is that Alice’s buckets may have many (about  $\log m / \log \log m$ ) elements. The primary contribution of [PSWW18] is to introduce a new type of hashing scheme, where there are two tables and Alice hashes each of her elements into *one* of the two tables using a 2-out-of-2 hashing scheme, and Bob hashes each of his elements into *both* of the two tables using a 1-out-of-2 (cuckoo) hashing scheme. Then a bucket-by-bucket equality test can be instantiated using any generic MPC protocol. The overall communication complexity of this approach is then  $\omega(m\lambda)$ . As in the cuckoo hashing protocols of [PSZ14, PSSZ15, PSZ16] the protocol is  $\omega(m\lambda)$  instead of  $O(m\lambda)$  because asymptotically Bob’s “stash” needs to be of size  $\omega(1)$ . In practice, however, they make do with a constant-sized stash, and this results in an extremely efficient circuit-based PSI that can then be implemented using any generic circuit-based MPC protocol. Note, however, that this “MPC-friendly” PSI protocol is significantly more communication intensive than existing “cleartext” PSI protocols (e.g., [PSZ14, PSSZ15, PSZ16]).

The recent work of [PSTY19] also computes *secret shares* of the intersection, but has lower communication cost than [CO18, PSWW18] both asymptotically ( $\Theta(m)$  cost) and concretely. The protocol of [PSTY19] uses the same hashing techniques as [PSZ14, PSSZ15, PSZ16], where Alice hashes her set into buckets using  $k$ -out-of- $k$  hashing, and Bob uses 1-out-of- $k$  (cuckoo) hashing to hash his elements into buckets, where each bucket holds a single element. Then, instead of doing one OPRF evaluation per bucket, they use a *programmable* OPRF, where, for each bucket, Alice programs the PRF to take the same value on all elements in her bucket. This programming can be easily achieved XOR-ing a traditional PRF output with a carefully interpolated polynomial. Then, they use a traditional MPC to do a single secure equality test for each bin (comparing only the PRF outputs). In this setting, where the results are secret shared, the stash can be handled, by simply reversing the roles of Alice and Bob, and having



Bob input only the elements he previously mapped to the stash.

### 3 Preliminaries

In this section, we review notation and some of the basic functionalities needed for our constructions. Our primitives are standard; we assume the reader has some familiarity with them, and thus we only provide brief reviews. Formal definitions can be found in cryptographic textbooks, e.g., [Gol01, Gol04].

#### 3.1 Common Notation

Below we show notation that is used repeatedly in the paper.

- $\mathcal{U}$ - The universe of possible set elements
- $m$  - Number of elements in each set
- $n$  - Number of bins
- $b$  - Size of bins
- $s$  - Size of stash
- $h$  - Hash function mapping elements to bins,  $h : \mathcal{U} \rightarrow [n]$
- $\lambda$  - A security parameter
- $\sigma$  - A parameter determining the failure probability
- $k$  - Number of hash functions used
- $\kappa$ - Key of a PRF, OPRF, or some variant thereof
- $X$  - Alice's set  $|X| = m$ ,  $X \subseteq \mathcal{U}$
- $Y$  - Bob's set  $|Y| = m$ ,  $Y \subseteq \mathcal{U}$
- $F$  - A PRF, OPRF, or some variant thereof
- $A[i][j]$  - The  $j^{\text{th}}$  element in Alice's  $i^{\text{th}}$  bin
- $B[i]$  - The (unique) element in Bob's  $i^{\text{th}}$  bin
- $\text{Stash}_B[i]$  - The  $i^{\text{th}}$  element of Bob's stash
- $d$  - Cost of an OPRF (bits)
- $d'$  - Cost of a one-time OPRF (bits, amortized)
- $t$  - Length of element bit-representation, typically  $t = \lceil \log_2 |\mathcal{U}| \rceil$

#### 3.2 Pseudorandom Functions

A pseudorandom function (PRF) is an efficiently computable, deterministic, keyed function with the property that for any adversary without the key, the outputs of the function  $F_\kappa(\cdot)$  are indistinguishable from independent uniformly random values.

**Definition 1** (PRF). *A deterministic function  $F : \{0, 1\}^\lambda \times \{0, 1\}^n \rightarrow \{0, 1\}^m$ , is called a pseudorandom function (PRF) if it satisfies the security properties captured by the game below.*

- The challenger generates a key,  $\kappa \xleftarrow{\$} \{0, 1\}^\lambda$

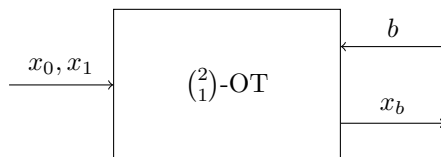


Figure 1: Oblivious Transfer. The sender provides two strings,  $x_0, x_1$ , and the receiver has a selection bit,  $b$ . The receiver receives a string  $x_b$ . The sender receives nothing.

- The challenger uniformly chooses a bit  $b \xleftarrow{\$} \{0, 1\}$ , and initializes a set  $X = \emptyset$ .
- The challenger and adversary engage in the following protocol where the adversary sends the challenger an input  $x_i \in \{0, 1\}^m$ , and receives a response  $y_i \in \{0, 1\}^m$  until the adversary outputs a guess  $b'$ . The challenger generates its responses as follows
  - If  $b = 0$ , the challenger sets  $y_i = F_{\kappa}(x_i)$
  - If  $b = 1$ , the challenger checks if there is a pair  $(x_i, y) \in X$ , if so the challenger responds with the value  $y$ . If not, the challenger uniformly selects  $y \in \{0, 1\}^m$ , and sets  $X = X \cup (x_i, y)$ , and returns  $y$  to the adversary.

The adversary is said to win the game if the adversary’s guess  $b'$  is equal to the challenger’s bit  $b$ . A PRF is said to be secure if any probabilistic polynomial-time adversary has a negligible (as a function of  $\lambda$ ) probability of winning the above game.

### 3.3 OT

Oblivious Transfer (OT) [Rab81, EGL85] is a two party protocol that allows one party (Bob) to privately select one of two input strings held by the other party (Alice).

**Definition 2** (OT). *Oblivious Transfer (OT) is a two party protocol that securely realizes the following functionality.*

**inputs:** Alice inputs strings  $x_0, x_1$ . Bob inputs a choice bit  $b$ .

**outputs:** Alice receives nothing. Bob receives  $x_b$ .

Sender Random-OT (ROT) is similar to oblivious transfer, except the strings  $x_0, x_1$  are not provided by Alice, but instead are randomly generated by the protocol itself. Although ROT seems to be a weaker primitive than OT, they are known to be equivalent [Cré87].

It is known that OT is symmetric (i.e., reversible) [WW06] and that OT is sufficient (“complete”) for general secure multiparty computation [Kil88, IPS08].

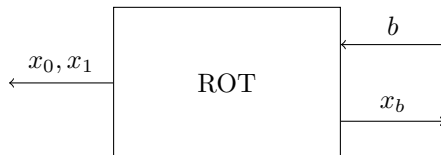


Figure 2: Random Oblivious Transfer. In the ROT functionality, the *sender* provides no input to the protocol, and instead the values  $x_0, x_1$  are generated uniformly at random by the protocol itself.

OT can be constructed generically from many different cryptographic primitives, including PIR [CMO00], projective hash proofs [Kal05, HK07], dual-mode encryption [PVW08] and noisy channels [IKO<sup>+</sup>11]. On the other hand, a black-box construction of OT from one-way permutations would imply  $P \neq NP$  [IR89], (perfect) OT cannot be constructed from quantum mechanical processes [Lo97], and quantum mechanics doesn’t even allow OT extension [SSS09, WW10].

One of the key features of OT is that a small number of “base” or “seed” OTs can be extended into a huge number of OTs with low overhead in terms of computation and communication [IKNP03]. Thus, although OT is inherently a public-key primitive (OT implies public-key encryption, but not vice-versa [GKM<sup>+</sup>00]), protocols that require a large number of OTs (e.g., [PSZ14, PSSZ15, PSZ16]) do not require a large number of public-key operations.

### 3.4 OPRFs

An Oblivious pseudorandom function (OPRF) is a two-party protocol for securely computing a PRF. The OPRF protocol securely realizes the functionality where one party (Alice) provides a PRF key,  $\kappa$ , and the other party (Bob) provides an input value,  $x$ . In the ideal functionality, Alice learns nothing, while Bob learns  $F_\kappa(x)$ . See Figure 3.

OPRFs were introduced in [FIPR05] as a means of achieving private keyword search, and since that time, many OPRF protocols have been introduced. A generic method for realizing the OPRF functionality is to use a general-purpose MPC protocol (e.g., garbled circuits) to implement a PRF.

**Definition 3** (OPRF). *An Oblivious Pseudo-Random Function (OPRF) is a two party protocol that securely realizes the following functionality, where  $F$  is a cryptographically secure pseudorandom function (PRF).*

**inputs:** Alice inputs a key  $\kappa$ . Bob inputs a value  $x$ .

**outputs:** Alice receives nothing. Bob receives  $F_\kappa(x)$ .

### 3.5 One-time OPRFs

A one-time OPRF is essentially an OPRF where the PRF key is randomly generated by the protocol (instead of specified by one of the players). In this

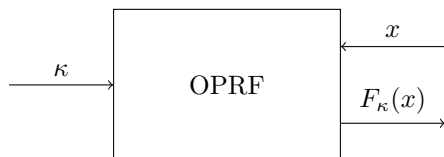


Figure 3: The Oblivious PRF (OPRF) functionality. The sender provides a PRF key,  $\kappa$ , and the receiver provides an input,  $x$ . The receiver learns  $F_\kappa(x)$ , and the sender learns nothing.



Figure 4: The one-time OPRF functionality. In the one-time OPRF functionality, the PRF key is not specified by the sender, but instead is generated by the protocol itself.

sense, the relationship between a one-time OPRF and an OPRF is similar to the relationship between Random OT and OT. The qualifier one-time indicates that the one-time OPRF protocol can only be used to securely evaluate  $F_\kappa(\cdot)$  once, since each additional run of the one-time OPRF will evaluate the PRF at a different key.

**Definition 4** (one-time OPRF). *A one-time Oblivious Pseudo-Random Function (one-time OPRF) is a two party protocol that securely realizes the following functionality, where  $F$  is a cryptographically secure pseudorandom function (PRF).*

**inputs:** *Alice inputs nothing. Bob inputs a value  $x$ .*

**outputs:** *Alice receives a uniformly chosen key,  $\kappa$ . Bob receives  $F_\kappa(x)$ .*

### 3.6 Multiple-choice hashing and Cuckoo hashing

Multiple-choice hashing is a technique for representing a set  $S$  using an array  $A$ , that has efficient space utilization and allows constant look-up time. First, a set of hash functions  $h_1, \dots, h_k$  is chosen, where  $h_i : \mathcal{U} \rightarrow [n]$ . An element  $x_j \in S$ , is stored in an index  $h_i(x_j)$  of  $A$ , for some  $i \in [k]$ . Look-ups of an element,  $x$ , simply require checking whether  $A[h_i(x)] = x$  for any  $i \in [k]$ . Multiple choice hashing assumes that the entire set is known in advance, so items in the set can be allocated indices in  $A$  using a 2D matching protocol.

A related technique, Cuckoo hashing, handles situations where the entire set is not known in advance, or may be modified dynamically. Like in Multiple-

choice hashing, an element  $x \in S$  will be stored in  $h_i(x)$  for some  $i \in [k]$ . The name Cuckoo hashing refers to the method used for insertion: a new element  $x$  will be placed in one of its allowed indices  $h_i(x)$  and if the spot is already occupied it will, like a Cuckoo bird taking over a nest, evict the existing element from that location. The evicted element will then seek a new home from its allowed indices, possibly evicting yet another element in the process. This process continues until an element can move to a space that was previously empty, or a pre-determined threshold in the number of evictions is reached, in which case the insertion fails.

Many previous PSI works refer to Cuckoo hashing. However, as outlined in more detail in Section 4.2 in the PSI case the items *are* known in advance, so static multiple-choice hashing can be used instead of dynamic Cuckoo hashing. This means that an allocation can always be found if one exists.

However, there is the chance that no allocation exists. In the case where  $k = 2$ , if the table is of size  $2(1 + \epsilon)n$ , for any  $\epsilon > 0$ , the probability that multiple-choice hashing fails is  $\Theta(1/n)$ , where the constant in the Theta-notation depends on  $\epsilon$  [Kut06]. In many situations, a failure probability of  $\Theta(1/n)$  is not acceptable. A solution to this is to have a “stash”, of maximum size  $s$ , in which to store any elements that were unable to be placed in the table. The probability of failure becomes  $\Theta(1/n^{s+1})$  [Mit09]. To achieve negligible failure probability,  $s = \omega(1)$  is needed. To achieve a fixed statistical failure probability  $2^{-\sigma}$ , it is required that  $s \geq \frac{\sigma}{\log(n)} + \frac{\log(c)}{\log(n)} - 1$  where  $c$  is the constant from the Theta-notation above.

## 4 Two generic optimizations for hashing-based PSI protocols

Here we present two improvements which can be used for hashing-based PSI protocols in general. Both are orthogonal to the primary contributions of this paper and the rest of the paper should still make sense if this section is skipped.

### 4.1 Permutation-based hashing [ANS10, PSSZ15]

When hashing elements into buckets, the size of the representation of each element can be reduced using the notion of permutation-based hashing [ANS10, PSSZ15]. In general, it takes  $\log |\mathcal{U}|$  bits to store an element  $x \in \mathcal{U}$ . In a traditional hash table, a hash function,  $h : \mathcal{U} \rightarrow [n]$  is chosen, and the element  $x$  is stored in the location indexed by  $h(x)$ . Notice, however, that the bucket index,  $h(x)$ , carries  $\log n$  bits of information, so it should be possible to reduce the information stored in the bucket from  $\log |\mathcal{U}|$  bits to  $\log |\mathcal{U}| - \log n$  bits, while still retaining the ability to uniquely recover an element  $x$ . This reduction in storage is possible, if we choose our hash function carefully. Permutation-based hashing provides a method for doing this, using a Feistel-style trick. First, we choose a public random permutation  $\pi : \mathcal{U} \rightarrow \mathcal{U}$  and, for each element  $x$  in a set, compute

$\hat{x} = \pi(x)$ . Suppose  $\hat{x}$  has bit-representation  $\hat{x} = \hat{x}_1 || \hat{x}_2$ , where  $\hat{x}_1$  has length  $\log n$ , and  $\hat{x}_2$  has length  $\log |\mathcal{U}| - \log n$ . Let  $f : \{0, 1\}^{\log |\mathcal{U}| - \log n} \rightarrow \{0, 1\}^{\log n}$  be a uniform hash function. Then we define  $h(\hat{x}) = \hat{x}_1 \oplus f(\hat{x}_2)$ , and we store  $\hat{x}_2$  in the bin defined by  $h(\hat{x})$ . The crucial observation here is that if  $\hat{x}$  and  $\hat{y}$  are in the same bin, and the stored values are the same (*i.e.*,  $\hat{x}_2 = \hat{y}_2$ ), then that means  $f(\hat{x}_2) = f(\hat{y}_2)$ , and since  $h(\hat{x}) = h(\hat{y})$ , we conclude that  $\hat{x}_1 = \hat{y}_1$ , which means  $\hat{x} = \hat{y}$  and therefore  $x = y$ .

This trick allows us to reduce the size of the representation of elements within the bins, while still providing the property that if two elements have the same representation and are in the same bin, they must be equal. In the context of PSI, this means that secure computations are done on elements with smaller representations, and this can result in noticeable efficiency improvements (quantified in [PSSZ15]). The random permutation,  $\pi(\cdot)$ , is needed because  $h(x)$  may no longer be a uniform hash function, even though  $f(x)$  is. If the input to  $h$  is correlated, this can increase the probability of a high number of hash-table collisions. The permutation guarantees that the inputs to  $h$  is uncorrelated. The permutation-based hashing trick can be used in all hashing-based PSI protocols, including ours. In the situation where there are multiple hash functions, the ID of the hash function must also be stored in the bin to allow equality tests [Lam16].

## 4.2 Static (offline) hashing

In many hashing-based PSI functions, Bob will allocate each element one of  $k$  possible hash functions, and Bob wishes to choose an optimal allocation, namely one that minimizes the number of elements that need to be stored in his stash.

Existing PSI protocols [PSZ14, PSSZ15, KKRT16] have used *online* Cuckoo hashing to compute Bob's allocation. Cuckoo hashing is designed to allow *dynamic* insertion of elements, but we observe that in the PSI setting, both parties know their sets in advance. In particular, Bob can allocate hash functions to elements based on knowledge of his entire set.

Instead of using dynamic cuckoo hashing, Bob can statically calculate an optimal allocation. He does this by treating the problem as a 2D matching problem (matching elements to potential bins). For dynamic cuckoo hashing, however, it remains an open problem of whether an optimal allocation can be found [Mit09]. As such this simple modification provides a guarantee that hash functions are allocated optimally. It also eases reasoning about the probability of failure for a given stash size.

## 5 Linear communication PSI via hashing and OPRFs

Here we present a PSI protocol with linear communication cost. In short, it does this by combining a hashing-based PSI protocol with an unbalanced PSI

scheme for handling the stash. Previous PSI protocols with linear communication tended to rely on oblivious polynomial evaluation, and instantiated the idea with an additively homomorphic cryptosystem [FNP04, HN10, KS05, DSMRY09] or require concrete number-theoretic assumptions [HL10, DCT10, DCT12, JL09, JL10]. By contrast, our protocol uses a hashing-based PSI protocol (like [PSZ14, PSSZ15, KKRT16]) which can be instantiated from OT, combined with an OPRF (which can also be instantiated with OT).

## 5.1 Construction

At a high level, our construction is as follows: Alice and Bob choose  $k$  hash functions,  $h_i : \{0, 1\}^* \rightarrow [n]$ . Then Alice hashes each of her elements into  $k$  bins, and Bob uses multiple-choice hashing to hash each of his elements into 1 (out of  $k$  possible) bins. Then they perform pairwise comparisons on the bins. For constant  $k$  we can set  $n$  to be  $O(m)$  and the hashing will succeed with probability  $1 - o(n)$ . To make this failure probability negligible, we allow Bob to keep a super-constant sized “stash” of elements that could not be allocated to a single bin.

Then, for each of Bob’s bins we use secure comparison protocol (like those described in [PSZ14, PSSZ15, KKRT16]) to compare the element in Bob’s bin to the elements in Alice’s corresponding bin. Finally, we need to compare Bob’s stash (of size  $\omega(1)$ ) to Alice’s elements (of size  $O(m)$ ). To do this, we use an unbalanced PSI protocol like those described in [CLR17, CHLR18, RA18, KLS<sup>+</sup>17]. For concreteness, we use the OPRF-based unbalanced PSI protocol described in [KLS<sup>+</sup>17] because its use of OPRFs is conceptually closest to the one-time OPRFs used in [KKRT16], and because many implementations of OPRFs were available for our benchmarks.

**Remark 1.** *Our construction is also compatible with silent-OT extension [BCG<sup>+</sup>19] that could be used to reduce the communication complexity of [PSZ14, PSSZ15, KKRT16]. This reduction in communication cost would come at an increase in computation cost, and currently silent-OT extension is only known to be secure under a variant of the LPN assumption.*

1. Set  $k \geq 2$ , and  $n = O(m)$ , and the players choose  $k$  hash functions  $h_i : \mathcal{U} \rightarrow [n]$ .
2. Bob will hash his elements into  $n$  bins using a 1-out-of- $k$  hashing scheme, such that each bin obtains at most one element. Elements that cannot be allocated to a single bin will be put in a “stash”,  $\text{Stash}_B$ , of size  $s = \omega(1)$ . Bob can compute this allocation efficiently. Let  $B[i]$  denote the element stored in Bob’s  $i$ th bin. Let  $\text{Stash}_B[i]$  for  $i \in [s]$  denote the  $i$ th element of the stash.

3. Alice will hash her elements into  $n$  bins, using a  $k$ -out-of- $k$  hashing scheme. Let  $C[i]$  denote the number of elements in Alice's  $i$ th bin, and let  $A[i][j]$  denote the element in the bin for  $1 \leq i \leq n$ ,  $1 \leq j \leq C[i]$ .
4. Alice and Bob will engage in  $n$  parallel executions of a one-time OPRF, where for  $i \in [n]$ , Alice learns a key  $\kappa_i$ , and Bob learns  $S_B^i \stackrel{\text{def}}{=} \text{PRF}_{\kappa_i}(B[i])$ . For each  $i \in [n]$ , Alice will locally compute  $S_{A,j}^i = \text{PRF}_{\kappa_i}(A[i][j])$  for  $j = 1, \dots, C[i]$ .
5. Alice will shuffle  $\{S_{A,j}^i\}_{i \in [n], j \in [C[i]]}$  and send the shuffled set to Bob. Note that this set will have exactly  $km$  elements.
6. Bob will locally compute the intersection of his non-stash set with Alice's set by finding which of the  $S_B^i$  are in the set received from Alice.
7. To handle the stash, Alice will generate a key,  $\kappa$ , for an OPRF, and Alice and Bob will engage in  $s$  executions of an OPRF protocol, where Bob learns  $R_B^i \stackrel{\text{def}}{=} \text{PRF}_{\kappa}(\text{Stash}_B[i])$  for  $i \in [s]$ .
8. Alice will compute  $R_A^i \stackrel{\text{def}}{=} \text{PRF}_{\kappa}(A[i][j])$  for  $i \in [n]$ ,  $j \in [C[i]]$  shuffle the set, and send it Bob who will locally compare these values to  $\{R_B^i\}_{i \in [s]}$  to find the intersection of Alice's set with the stash. Note that the set Alice sends will have exactly  $m$  elements.

Figure 5: The high-level outline of our algorithm

The communication cost of the protocol is the sum of the costs of the following:

- $n$  parallel executions of a one-time OPRF. The cost of this will depend on how the one-time OPRF is instantiated.
- Alice will send Bob  $km$  evaluations of her inputs under the one-time OPRF keys.
- $s = \omega(1)$  secure computations of the Stash OPRF.<sup>3</sup>
- Alice will send Bob  $m$  evaluations of her inputs under the Stash OPRF key.

The exact communication cost will depend on how the one-time OPRF and OPRF are instantiated, but standard constructions allow the entire protocol to run with communication linear in  $m$ .

<sup>3</sup> If the one-time OPRF and the OPRF compute the same PRF, a small optimization is possible. Rather than Alice generating the OPRF key, it could be generated by a one-time OPRF on the first element of the stash. The protocol would then use  $(n+1)$  parallel one-time OPRFs and  $(s-1)$  OPRFs. The OPRF and one-time OPRF we used represent different PRFs, so we could not use this.



In Sections 5.2 and 5.3, we review existing methods for implementing the one-time OPRF using methods from [PSSZ15] and [KKRT16], and give the communication complexity of each approach. Then we describe how to implement the OPRF using methods from [KLS<sup>+</sup>17]. Finally, in Section 5.5 we compare the concrete communication complexity of our modified protocol with the best existing protocol [KKRT16].

## 5.2 An OT-masking-based protocol

In this section, we review the OT-masking-based OPRF protocol from [PSZ16]. When instantiated with OT, this protocol implements a standard (multi-time) OPRF, when instantiated with ROT, this protocol implements a one-time OPRF, which is sufficient for the PSI applications.

1. Bob will represent his element,  $B[i]$ , (the contents of the  $i$ th bucket) as a bit vector of length  $t$ , including a tag for which of the  $k$  hash functions was used. Using permutation-based hashing (Section 4.1), this can be done with  $t = \log_2 |\mathcal{U}| - \log_2 n + \log_2 k$ .
2. Alice constructs a vector of length  $2^t$  of random masks,  $M$ . This vector is the OPRF key (in fact, actually a truth-table for a truly random function). For input  $j \in [2^t]$ ,  $M_j$  is the evaluation of the OPRF on input  $j$ . Since Alice can choose the key, this is a normal OPRF rather than a one-time OPRF. If Alice and Bob are using ROT (instantiating a one-time OPRF) they can skip this step.
3. To evaluate the (one-time) OPRF on his input,  $B[i]$ , Alice and Bob engage in a  $\binom{2^t}{1}$  string-OT (respectively  $\binom{2^t}{1}$  string-ROT). Bob uses his input  $B[i] \in [2^t]$  as the input to the OT. From this, Bob learns  $M_{B[i]}$ , which is the evaluation of the PRF on  $B[i]$ .

Figure 6: Implementing OPRF using Oblivious Transfer

This protocol requires a single  $\binom{2^t}{1}$ -string OTs (for random strings) for each bucket, and thus a total of  $n \binom{2^t}{1}$ -string OTs. Using OT extension [IKNP03], these can be generated using  $\lambda$  base OTs.

The outline above follows the construction of [PSZ16], which fixed an error in earlier OT-masking approaches [PSZ14, PSSZ15].

## 5.3 Batching OPRFs

In this section, we describe how to use our hashing scheme with the Batched, Related-Key OPRFs (BaRK-OPRFs) introduced in [KKRT16]. At a high level, this protocol is very similar to the OT-based protocol described in Section 5.2.

As this is the most efficient instantiation of the one-time OPRF used in our scheme we provide a description in more detail below.

Before outlining the actual construction of [KKRT16], we review some of the necessary terminology. First, a *relaxed*-PRF is a pair  $(F, \tilde{F})$  where  $F$  is a PRF such that 1)  $F_\kappa(x)$  can be computed from  $\tilde{F}_\kappa(x)$  and 2)  $\tilde{F}_\kappa(x)$  does not improve the adversary's distinguishing advantage in the PRF security experiment. In other words, given query access to  $\tilde{F}_\kappa(\cdot)$ ,  $F_\kappa(\cdot)$  appears pseudo-random on all unqueried points.

The original OT-extension protocol of [IKNP03] relied on a *correlation-robust* hash function. In [KKRT16] the notion of correlation-robustness is extended as follows.

**Definition 5** (*k*-Hamming Correlation Robustness). *Let  $H : \{0, 1\}^* \rightarrow \{0, 1\}^v$  be a hash function such that for all  $z_i \in \{0, 1\}^*$ , and  $a_i, b_i \in \{0, 1\}^n$  for  $i = 1, \dots, m$ , with  $\text{wt}(b_i) \geq k$ , we have*

$$\begin{aligned} & \left\{ \{H(z_i \| a_i \oplus [b_i \cdot s])\}_{i \in [m]} \mid s \xleftarrow{\$} \{0, 1\}^n \right\} \\ & \approx \\ & \left\{ \{r_i\}_{i \in m} \mid r_i \xleftarrow{\$} \{0, 1\}^v \right\} \end{aligned}$$

The definition of correlation robustness in [IKNP03] required  $k = n$ , i.e.,  $b_i \cdot s = s$ .

A pseudo-random code is a relaxation of an error-correcting code such that for any two distinct messages, their encodings have a large minimum distance with high probability over the choice of a specific code from the family.

**Definition 6** (Pseudo-Random Code [KKRT16]). *A family of functions,  $\mathcal{C}$ , is called a  $(d, \epsilon)$  pseudorandom code (PRC) if for all strings  $x \neq x'$ ,*

$$\Pr_{C \leftarrow \mathcal{C}} [\text{wt}(C(x) \oplus C(x')) < d] \leq 2^{-\epsilon}$$

The key security definition of [KKRT16] is the notion of related-key, relaxed PRFs. These are PRFs that remain secure when the challenger chooses  $n$  keys,  $\kappa_1, \dots, \kappa_n$ , and the adversary sees the *relaxed* output for each key, then any  $m$  additional outputs (of the PRF,  $F$ ) corresponding to any of the keys are indistinguishable from random.

**Definition 7** (*m*-related key PRFs [KKRT16]). *An  $m$ -related key PRF is a pair of functions  $F, \tilde{F}$ , and security is defined relative to the following game:*

1. *The adversary chooses input strings  $\{x_j\}_{j \in [n]}$  and pairs  $\{(j_i, y_i)\}_{i \in [m]}$ , with  $j_i \in [n]$  and  $y_i \neq x_{j_i}$  for  $i \in [m]$ .*
2. *The challenger chooses PRF keys  $\kappa^*, \kappa_1, \dots, \kappa_n$  and a challenge bit  $b \in \{0, 1\}$ .*
  - *If  $b = 0$ , the challenger sends  $\left\{ \tilde{F}_{(\kappa^*, \kappa_j)}(x_j) \right\}_j$  and  $\left\{ F_{(\kappa^*, \kappa_{j_i})}(y_i) \right\}_i$  to the adversary.*

- If  $b = 1$ , the challenger generates  $m$  random strings  $z_i \leftarrow \{0, 1\}^v$ , and sends  $\left\{ \tilde{F}_{(\kappa^*, \kappa_j)}(x_j) \right\}_j$  and  $\{z_i\}_i$  to the adversary.

3. The adversary outputs a guess  $b'$ , and the adversary wins if  $b' = b$ . We say the adversary's advantage is  $\Pr[b' = b] - 1/2$ .

We say the pair  $F, \tilde{F}$  is secure if the adversary's advantage is negligible. Intuitively, the pair  $F, \tilde{F}$  is an  $m$ -related key PRF if the relaxed output,  $\tilde{F}$ , on  $n$  distinct keys, does not reveal information that would allow an adversary to distinguish  $F$  from a true random function, even if the inputs are arbitrarily correlated across keys.

Related-key PRFs can be efficiently instantiated using a pseudo random code, and a correlation robust hash function as follows.

**Lemma 1** (Lemma 5 in [KKRT16]). *If  $C$  is a  $(d, \epsilon + \log m)$ -pseudo random code, and  $2^{-\epsilon}$  is negligible and  $H$  is a  $d$ -Hamming correlation robust hash function, then*

$$\begin{aligned} F_{((C,s),(q,j))}(r) &= H(j || q \oplus [C(r) \cdot s]) \\ \tilde{F}_{((C,s),(q,j))}(r) &= (j, C, q \oplus [C(r) \cdot s]) \end{aligned}$$

is an  $m$ -related key PRF as in Definition 7.

1. Alice chooses a random PRC  $C \xleftarrow{\$} \mathcal{C}$  and sends the code to Bob.
2. Alice chooses a key  $s \xleftarrow{\$} \{0, 1\}^k$ .
3. Bob generates two matrices  $T_0, T_1 \in \{0, 1\}^{m \times k}$  as follows. For  $j = 1, \dots, m$ ,
  - (a) The  $j$ th row of  $T_0$  is generated uniformly at random  $t_{0,j} \xleftarrow{\$} \{0, 1\}^k$ .
  - (b) The  $j$ th row of  $T_1$  is defined as  $t_{1,j} = C(r_j) \oplus t_{0,j}$ .
 The  $i$ th columns of  $T_0$  and  $T_1$  are denoted  $t_0^i, t_1^i$  respectively.
4. Alice and Bob engage in  $k$  parallel instances of  $\binom{2}{1}$ -OT for strings of length  $m$  as follows:
  - Bob acts as sender and his inputs are the  $k$  columns of  $T_0$  and  $T_1$ , i.e., Bob's inputs to the OT are  $\{t_0^i, t_1^i\}_{i \in [k]}$
  - Alice acts as receiver and her inputs are the  $k$  bits of  $s$ , i.e.,  $\{s_i\}_{i \in [k]}$

- Alice receives  $k$  outputs (each of length  $m$ ) denoted  $\{q^i\}_{i \in [k]}$

Alice creates the  $m \times k$  matrix  $Q$  whose columns are the received vectors  $\{q^i\}$ . Thus the  $i$ th column of  $Q$  is  $t_{s_i}^i$ . Let  $q_j$  denote the  $j$ th row of  $Q$ . Then

$$q_j = ((t_{0,j} \oplus t_{1,j}) \cdot s) \oplus t_{0,j} = t_{0,j} \oplus (C(r_j) \cdot s)$$

5. For  $j \in [m]$ , Alice keeps the PRF seed  $((C, s), (j, q_j))$
6. For  $j \in [m]$ , Bob keeps the relaxed PRF output  $(j, C, t_{0,j})$ .

Figure 7: Instantiating a BaRK-OPRF using OT [KKRT16].

At the end of the protocol in Figure 7, Alice has the keys  $\kappa^* = (C, s)$ , and  $\kappa_j = (j, q_j)$ , which allows her to evaluate the BaRK-OPRF

$$F_{(\kappa^*, \kappa_j)}(r) = F_{((C, s), (j, q_j))}(r) = H(j || q_j \oplus [C(r) \cdot s])$$

for any  $r$ , and Bob has the relaxed PRF outputs

$$(j, C, t_{0,j}) = (j, C, q_j \oplus [C(r_j) \cdot s]) = \tilde{F}_{(\kappa^*, \kappa_j)}(r_j)$$

which allows him to compute  $F_{(\kappa^*, \kappa_j)}(r_j)$ .

The BaRK-OPRF protocol allows Alice and Bob to securely compute  $n$  parallel, one-time OPRFs in a very efficient manner. Since the protocol is inherently “batched,” it does not fit exactly into the one-time OPRF paradigm described in Figure 5. Conceptually, the idea is essentially the same. Alice and Bob will engage in a BaRK-OPRF protocol, where Bob learns  $\tilde{F}_{(\kappa^*, \kappa_j)}(B[i])$  for  $i = 1, \dots, n$ , which allows him to learn the PRF outputs  $F_{(\kappa^*, \kappa_j)}(B[i])$ . Alice will learn the PRF keys,  $\kappa^*, \kappa_1, \dots, \kappa_n$ .

1. Bob has inputs  $B[i]$ , for  $i = 1, \dots, n$
2. Alice has no inputs
3. Alice and Bob will use the BaRK-OPRF protocol (Figure 7) with Bob providing the contents of his  $n$  buckets as his inputs. At the end of the protocol, Alice has keys:

$$\kappa^*, \kappa_1, \dots, \kappa_n$$

and Bob has relaxed PRF outputs  $\tilde{F}_{(\kappa^*, \kappa_i)}(B[i])$

4. From the relaxed outputs, Bob can compute  $S_B^i \stackrel{\text{def}}{=} F_{(\kappa^*, \kappa_i)}(B[i])$ .

5. From the PRF keys  $\kappa^*, \kappa_1, \dots, \kappa_n$ , Alice can compute  $F_{(\kappa^*, \kappa_i)}(x)$  for any  $x$  of her choosing.

Figure 8: Using BaRK-OPRFs to instantiate  $n$  parallel one-time OPRFs. Concretely, our implementation will use the BaRK-OPRF protocol to compute the  $n$  parallel one-time OPRFs, and then use a different (multi-use) PRF for the stash comparison. See section 5.5.

## 5.4 Security

The PSI protocols outlined in Section 5.1 are formed by taking existing one-time OPRF-based PSI protocols (e.g. [PSZ14, PSSZ15, PSZ16, KKRT16]) and combining them (in parallel) with an OPRF-based PSI protocol (e.g. [KLS<sup>+</sup>17]) to handle the unbalanced stash. The proof of security (against honest-but-curious adversaries) follows in a straightforward manner from the security of the two underlying protocols.

For completeness, we outline the security of our protocol here. We define security in the standard, simulation paradigm, i.e., a PSI protocol is secure if there exists an efficient (probabilistic polynomial-time) simulator that can simulate the view of each player given the output of the protocol alone.

In our basic protocol (Figure 5), Alice’s view consists of  $n$  parallel executions of a one-time-OPRF (from which she receives  $n$  random keys), and  $s$  applications of an OPRF (from which she receives nothing). Thus her view is completely independent of Bob’s input and can be trivially simulated since her view consists solely of random strings. Bob’s view consists of  $m$  PRF-outputs provided by the one-time OPRF protocol,  $km$  PRF outputs provided by Alice,  $s$  PRF outputs provided by the OPRF protocol for the stash, and  $m$  PRF outputs provided by Alice to compare with the stash. Given the intersection of Alice and Bob’s sets, Bob’s view can be easily simulated by providing  $km + m$  random strings (corresponding to the one-time OPRF round) with the correct intersection pattern, and  $s + m$  random strings (corresponding to the OPRF round) with the correct intersection pattern.

## 5.5 Concrete communication benchmarks

In [KMW09] it was shown that under the assumption that the dynamic (online) cuckoo hashing algorithm achieves the optimal load, allocating  $m$  elements to  $n = O(m)$  buckets with a stash of size  $s$  will succeed with probability at least  $1 - O(n^{-(s+1)})$ . The protocols of [PSZ14, PSSZ15, PSZ16, KKRT16] set  $n = 1.2m$ . Thus, under the assumption that online cuckoo hashing achieves the optimal offline load, to achieve a negligible probability of failure, the stash size,  $s$ , must be  $s = \omega(1)$  (e.g.,  $s = O(\log n)$ ).

In the context of PSI, a cuckoo-hashing failure reveals information about the

players’ sets, so the failure probability should be set below the security threshold. In the work of [PSZ14, PSSZ15], they set security threshold to be a constant  $2^{-\sigma}$ , with  $\sigma = 40$ , independent of the set sizes. The asymptotic failure rate is known to be  $\Theta(n^{-(s+1)})$  [KMW09], which implies  $s \geq \frac{\sigma}{\log(n)} + \frac{\log(c)}{\log n} - 1$  provides failure probability below  $2^{-\sigma}$  for some constant  $c$ . In [PSSZ15], they empirically tested the failure probability with 2-way cuckoo hashing, and different stash sizes for  $m$  up to about  $2^{14}$ . Then they extrapolated these failure probabilities up to larger set sizes, and used these values to choose the stash sizes to be  $s \sim 40/\log(n)$  (see [PSSZ15] Figure 2). This is consistent with the analysis above, assuming the constant  $c$  is not very large. The scheme of [KKRT16] uses similar stash parameters even though they are hashing with three hash functions instead of two.

In practice, fixing a concrete security parameter that does not increase with  $m$  means that the necessary stash size *decreases* as the set sizes increase, whereas an asymptotic analysis (which assumes that the failure probability should be negligible in  $m$ ) requires that the stash size *increases* as the set sizes increase. Thus this choice of a concrete security parameter means that the concrete and asymptotic *performance* metrics diverge as the set sizes increase.

For a hash table of size  $n$ , and a stash of size  $s$ , the hashing protocol of [KKRT16] requires  $n$  applications of a one-time OPRF, and the communication of  $n$  (truncated) PRF outputs in the main phase, and  $s$  one-time OPRF evaluations and the communication of  $ms$  PRF outputs in the stash phase.

We replace the stash computation with an unbalanced PSI protocol based on a standard (reusable) OPRF, reducing the communication in the stash phase to  $s$  evaluations of an OPRF followed by sending  $m$  PRF outputs.

In practice, because equality testing is done by comparing the pseudorandom masks (PRF outputs), there is no need to transmit the entire mask. Instead, to achieve error probability less than  $2^{-\sigma}$ , it is sufficient to transmit only about  $v = \sigma + 2 \log m$  bits of the mask, and in this case, by the birthday bound, the probability of a spurious collision between masks is negligible in  $m$ . This is what is done in practice by [PSSZ15, KKRT16].

Let  $d$  denote the number of bits required for an OPRF application, and  $d' < d$  be the number of bits required for a one-time OPRF evaluation. Then our protocol obtains a concrete performance improvement whenever  $sd + mv < sd' + msv$ .

We therefore obtain a concrete improvement whenever

$$\frac{s(d - d')}{v(s - 1)} < m$$

In [ARS<sup>+</sup>15], they report that a single (GMW-based) AES evaluation using ABY can be computed using only 170 kb of communication. Using their custom, “MPC-friendly” PRF, the garbled circuit can be computed using only 23 kb of communication ([ARS<sup>+</sup>15] Table 6). Similarly, 1024 garbled AES representations were computed using 185 Mb of communication ([KLS<sup>+</sup>17] Table 5), which reduces to about 177 kb per AES circuit (including regenerating the OTs).

Table 1: This table shows the communication cost of a single OPRF evaluation, and the minimum number of elements for which replacing the OPRF-based comparison of the stash in [KKRT16] would be improved by switching to our protocol. The AES Obliv-C and AES EMP benchmarks were obtained via our internal tests. The ABY benchmarks were taken from [ARS<sup>+</sup>15], the Yao benchmarks were taken from [KLS<sup>+</sup>17], and the LowMC benchmarks were taken from [ARS<sup>+</sup>15].

| OPRF                                     | Comm. Cost ( $d$ ) | $m$ at which our protocol outperforms [KKRT16] |
|--|--------------------|--|
| AES (Yao - Obliv-C)                      | 8 Mb               | $2^{21}$                                       |
| AES (Yao - EMP)                          | 212 Kb             | $2^{16}$                                       |
| AES (GMW - ABY) [ARS <sup>+</sup> 15]    | 170 Kb             | $2^{15}$                                       |
| AES (Yao - OblivM) [KLS <sup>+</sup> 17] | 177 Kb             | $2^{15}$                                       |
| LowMC [ARS <sup>+</sup> 15]              | 23 Kb              | $2^{12}$                                       |

Although our protocol outperforms [KKRT16] asymptotically, the exact point at which it starts to outperform [KKRT16] concretely depends on the exact implementation of the OPRF. In Table 1 we show the communication costs of four different standard OPRF implementations, and give the number of elements ( $m$ ) at which our protocol starts to outperform that of [KKRT16].

Thus, using an off-the-shelf AES128 implementation, we start to see concrete improvements in communication costs by replacing the one-time OPRF protocol of [KKRT16] with a true OPRF, whenever the sets being compared have more than  $2^{15} = 32768$  elements. Using the “MPC-friendly” PRF, LowMC [ARS<sup>+</sup>15], we start to see concrete efficiency improvements for sets of size  $2^{12} = 4096$ .

Figure 9 shows the overall reduction in communication when the BaRK protocol [KKRT16] is replaced with our protocol instantiated with LowMC. Figure 10 shows the overall reduction in communication (over [KKRT16]) when our protocol is instantiated using different OPRF instantiations. Since the OPRF cost does not increase with  $m$ , all three instantiations of our protocols approach the same asymptotic improvement (about 10%) over the protocol of [KKRT16]. Note that these results are conservative, in that we use the same stash sizes as [KKRT16] which assumes that the probability of failure is constant and therefore the stash sizes decrease as  $m$  increases. The step-down points of non-differentiability in Figures 9 and 10 correspond to these decreases in the stash size. From a theoretical perspective to maintain negligible failure rates, the stash must be *increasing* in  $m$ , and in that setting we would see corresponding step-up points in the graphs at each point where the stash size increased.

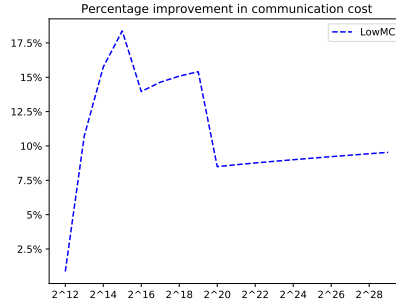


Figure 9: The percentage improvement in overall communication cost when modifying the BaRK protocol to use the LowMC-based PRF to compare the stash. For most values of  $m$ , our modification reduces the overall communication cost of the protocol by 10 to 15%. The points of non-differentiability in the graph correspond to the places where the stash sizes drop (we use the same stash sizes as [PSSZ15, KKRT16])

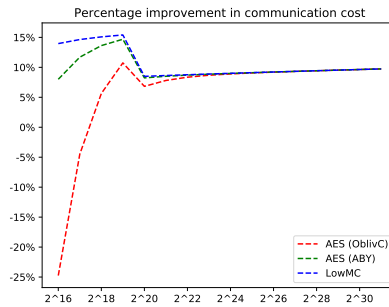


Figure 10: As  $m$  increases, and the stash size stays at 2, all three OPRF protocols tend towards a 10% improvement in communication cost over the BaRK protocol of [KKRT16]



## 6 Conclusion

PSI is one of the most basic and fundamental types of secure computation, and numerous diverse types of PSI protocols have been proposed and implemented. Existing PSI protocols are highly optimized and extremely efficient.

In this work, we show how simple, modular composition of existing PSI protocols leads to both asymptotic and concrete improvements in efficiency. Specifically we show that the one-time OPRF protocols of [PSZ14, PSSZ15, KKRT16] can be improved both asymptotically (from super-linear to linear) and, for  $n \geq 2^{12}$ , concretely by replacing the stash-comparison step with an unbalanced PSI protocol [KLS<sup>+</sup>17].

## References

- [ANS10] Yuriy Arbitman, Moni Naor, and Gil Segev. Backyard cuckoo hashing: Constant worst-case operations with a succinct representation. In *FOCS*, pages 787–796. IEEE, 2010.
- [AP11] Rasmus Resen Amossen and Rasmus Pagh. A new data layout for set intersection on gpus. In *Parallel & Distributed Processing Symposium (IPDPS), 2011 IEEE International*, pages 698–708. IEEE, 2011.
- [ARS<sup>+</sup>15] Martin R Albrecht, Christian Rechberger, Thomas Schneider, Tyge Tiessen, and Michael Zohner. Ciphers for MPC and FHE. In *EUROCRYPT*, pages 430–454. Springer, 2015.
- [BCG<sup>+</sup>19] Elette Boyle, Geoffroy Couteau, Niv Gilboa, Yuval Ishai, Lisa Kohl, and Peter Scholl. Efficient pseudorandom correlation generators: Silent OT extension and more. IACR ePrint 2019/448, 2019.
- [CHLR18] Hao Chen, Zhicong Huang, Kim Laine, and Peter Rindal. Labeled PSI from fully homomorphic encryption with malicious security. In *CCS*, pages 1223–1237. ACM, 2018.
- [CLR17] Hao Chen, Kim Laine, and Peter Rindal. Fast private set intersection from homomorphic encryption. In *CCS*, pages 1243–1255, 2017.
- [CMO00] Giovanni D. Crescenzo, Tal Malkin, and Rafail Ostrovsky. Single Database Private Information Retrieval Implies Oblivious Transfer. In *Eurocrypt*, volume 1807, pages 122–138, 2000.
- [CO18] Michele Ciampi and Claudio Orlandi. Combining private set-intersection with secure two-party computation. In *SCN*, 2018.
- [Cré87] Claude Crépeau. Equivalence between two flavours of oblivious transfers. In *CRYPTO*, pages 350–354, 1987.

- [DCT10] Emiliano De Cristofaro and Gene Tsudik. Practical private set intersection protocols with linear complexity. In *FC*, volume 10, pages 143–159. Springer, 2010.
- [DCT12] Emiliano De Cristofaro and Gene Tsudik. Experimenting with fast private set intersection. *Trust*, 7344:55–73, 2012.
- [DCW13] Changyu Dong, Liqun Chen, and Zikai Wen. When private set intersection meets big data: an efficient and scalable protocol. In *CCS*, pages 789–800, 2013.
- [DSMRY09] Dana Dachman-Soled, Tal Malkin, Mariana Raykova, and Moti Yung. Efficient robust private set intersection. In *Applied Cryptography and Network Security*, pages 125–142. Springer, 2009.
- [DSZ15] Daniel Demmler, Thomas Schneider, and Michael Zohner. Aby-a framework for efficient mixed-protocol secure two-party computation. In *NDSS*, 2015.
- [EGL85] Shimon Even, Oded Goldreich, and Abraham Lempel. A randomized protocol for signing contracts. *Communications of the ACM*, 28(6):637–647, 1985.
- [FBI19] FBI. National crime information center (ncic). <https://www.fbi.gov/services/cjis/ncic>, July 2019.
- [FIPR05] Michael J Freedman, Yuval Ishai, Benny Pinkas, and Omer Reingold. Keyword search and oblivious pseudorandom functions. In *TCC*, volume 3378, pages 303–324, 2005.
- [FNP04] Michael J Freedman, Kobbi Nissim, and Benny Pinkas. Efficient private matching and set intersection. In *EUROCRYPT*, pages 1–19, 2004.
- [GKM<sup>+</sup>00] Yael Gertner, Sampath Kannan, Tal Malkin, Omer Reingold, and Mahesh Viswanathan. The relationship between public key encryption and oblivious transfer. In *FOCS*, page 325, 2000.
- [GMW87] Oded Goldreich, Silvio Micali, and Avi Wigderson. How to play any mental game. In *STOC*, pages 218–229, 1987.
- [Gol01] Oded Goldreich. *Foundations of cryptography: volume 1*. Cambridge university press, 2001.
- [Gol04] Oded Goldreich. *Foundations of cryptography: volume 2*. Cambridge university press, 2004.
- [GW10] Pu Gao and Nicholas C Wormald. Load balancing and orientability thresholds for random hypergraphs. In *Proceedings of the forty-second ACM symposium on Theory of computing*, pages 97–104. ACM, 2010.

- [HEK12] Yan Huang, David Evans, and Jonathan Katz. Private set intersection: Are garbled circuits better than custom protocols? In *NDSS*, 2012.
- [HK07] Dennis Hofheinz and Eike Kiltz. Secure hybrid encryption from weakened key encapsulation. In *CRYPTO*, Berlin, Heidelberg, 2007.
- [HL10] Carmit Hazay and Yehuda Lindell. Efficient protocols for set intersection and pattern matching with security against malicious and covert adversaries. *Journal of cryptology*, 23(3):422–456, 2010.
- [HN10] Carmit Hazay and Kobbi Nissim. Efficient set operations in the presence of malicious adversaries. In *PKC*, volume 6056, pages 312–331, 2010.
- [IKNP03] Yuval Ishai, Joe Kilian, Kobbi Nissim, and Erez Petrank. Extending oblivious transfers efficiently. In *CRYPTO*, volume 2729, pages 145–161. Springer, 2003.
- [IKO<sup>+</sup>11] Yuval Ishai, Eyal Kushilevitz, Rafail Ostrovsky, Manoj Prabhakaran, Amit Sahai, and Jürg Wullschleger. Constant-Rate Oblivious Transfer from Noisy Channels. In *CRYPTO*, volume 6841, chapter 38, pages 667–684. 2011.
- [IKOS08] Yuval Ishai, Eyal Kushilevitz, Rafail Ostrovsky, and Amit Sahai. Cryptography with constant computational overhead. In *STOC*, pages 433–442, 2008.
- [IPS08] Yuval Ishai, Manoj Prabhakaran, and Amit Sahai. Founding Cryptography on Oblivious Transfer - Efficiently. In *CRYPTO*, pages 572–591, 2008.
- [IR89] Russell Impagliazzo and Steven Rudich. Limits on the Provable Consequences of One-Way Permutations. In *STOC*, pages 44–61, 1989.
- [JL09] Stanisław Jarecki and Xiaomin Liu. Efficient oblivious pseudorandom function with applications to adaptive ot and secure computation of set intersection. In *TCC*, volume 5444, pages 577–594. Springer, 2009.
- [JL10] Stanisław Jarecki and Xiaomin Liu. Fast secure computation of set intersection. *Security and Cryptography for Networks*, pages 418–435, 2010.
- [Kal05] Yael T. Kalai. Smooth Projective Hashing and Two-Message Oblivious Transfer. In *EUROCRYPT*, pages 78–95, 2005.

- [Kil88] Joe Kilian. Founding cryptography on oblivious transfer. In *STOC*, STOC '88, pages 20–31, 1988.
- [KKRT16] Vladimir Kolesnikov, Ranjit Kumaresan, Mike Rosulek, and Ni Trieu. Efficient batched oblivious PRF with applications to private set intersection. In *CCS*, pages 818–829, 2016.
- [KL02] Michał Karoński and Tomasz Łuczak. The phase transition in a random hypergraph. *Journal of Computational and Applied Mathematics*, 142(1):125–135, 2002.
- [KLS<sup>+</sup>17] Ágnes Kiss, Jian Liu, Thomas Schneider, N Asokan, and Benny Pinkas. Private set intersection for unequal set sizes with mobile applications. *PoPETs*, 4:97–117, 2017.
- [KMW09] Adam Kirsch, Michael Mitzenmacher, and Udi Wieder. More robust hashing: Cuckoo hashing with a stash. *SIAM Journal on Computing*, 39(4):1543–1561, 2009.
- [KS05] Lea Kissner and Dawn Song. Privacy-preserving set operations. In *CRYPTO*, volume 3621, pages 241–257, 2005.
- [Kut06] Reinhard Kutzelnigg. Bipartite Random Graphs and Cuckoo Hashing. In Philippe Chassaing et al., editors, *Fourth Colloquium on Mathematics and Computer Science Algorithms, Trees, Combinatorics and Probabilities*, volume DMTCS Proceedings vol. AG, Fourth Colloquium on Mathematics and Computer Science Algorithms, Trees, Combinatorics and Probabilities of *DMTCS Proceedings*, pages 403–406, Nancy, France, 2006. Discrete Mathematics and Theoretical Computer Science.
- [Lam16] Mikkel Lambæk. Breaking and fixing private set intersection protocols. IACR ePrint 2016/665, 2016.
- [Lel12] Marc Lelarge. A new approach to the orientation of random hypergraphs. In *Proceedings of the twenty-third annual ACM-SIAM symposium on Discrete Algorithms*, pages 251–264. SIAM, 2012.
- [Lo97] Hoi K. Lo. Insecurity of quantum secure computations. *Physical Review A*, 56:1154–1162, 1997.
- [LP14] Po-Shen Loh and Rasmus Pagh. Thresholds for extreme orientability. *Algorithmica*, 69(3):522–539, 2014.
- [MBD12] Dilip Many, Martin Burkhart, and Xenofontas Dimitropoulos. Fast private set operations with SEPIA. Technical Report 345, ETH Zurich, march 2012.
- [Mit09] Michael Mitzenmacher. Some open questions related to cuckoo hashing. In *ESA*, pages 1–10, 2009.

- [PRTY19] Benny Pinkas, Mike Rosulek, Ni Trieu, and Avishay Yanai. SpOT-light: Lightweight private set intersection from sparse OT extension. 2019.
- [PSSW09] Benny Pinkas, Thomas Schneider, Nigel P Smart, and Stephen C Williams. Secure two-party computation is practical. In *ASIACRYPT*, volume 9, pages 250–267, 2009.
- [PSSZ15] Benny Pinkas, Thomas Schneider, Gil Segev, and Michael Zohner. Phasing: Private set intersection using permutation-based hashing. In *USENIX Security Symposium*, pages 515–530, 2015.
- [PSTY19] Benny Pinkas, Thomas Schneider, Oleksandr Tkachenko, and Avishay Yanai. Efficient circuit-based psi with linear communication. In *EUROCRYPTO*, pages 122–153, 2019.
- [PSWW18] Benny Pinkas, Thomas Schneider, Christian Weinert, and Udi Wieder. Efficient circuit-based PSI via cuckoo hashing. In *EUROCRYPT*, 2018.
- [PSZ14] Benny Pinkas, Thomas Schneider, and Michael Zohner. Faster private set intersection based on ot extension. In *USENIX*, pages 797–812, 2014.
- [PSZ16] Benny Pinkas, Thomas Schneider, and Michael Zohner. Scalable private set intersection based on OT extension. IACR Cryptology ePrint Archive, 2016.
- [PVW08] Chris Peikert, Vinod Vaikuntanathan, and Brent Waters. A Framework for Efficient and Composable Oblivious Transfer. In David Wagner, editor, *CRYPTO*, volume 5157, pages 554–571, 2008.
- [RA18] Amanda C Davi Resende and Diego F Aranha. Faster unbalanced private set intersection. *FC 2018*, 2018.
- [Rab81] Michael O. Rabin. How to exchange secrets with oblivious transfer. Technical Report TR-81, Harvard University, 1981.
- [RR17] Peter Rindal and Mike Rosulek. Improved private set intersection against malicious adversaries. In *EUROCRYPT*, pages 235–259, 2017.
- [SFF14] Aaron Segal, Bryan Ford, and Joan Feigenbaum. Catching bandits and only bandits: Privacy-preserving intersection warrants for lawful surveillance. In *4th USENIX Workshop on Free and Open Communications on the Internet (FOCI 14)*, San Diego, CA, 2014. USENIX Association.
- [SSS09] Louis Salvail, Christian Schaffner, and Miroslava Sotáková. On the power of two-party quantum cryptography. In *ASIACRYPT*, pages 70–87, 2009.

- [Wak68] Abraham Waksman. A permutation network. *Journal of the ACM (JACM)*, 15(1):159–163, 1968.
- [WW06] Stefan Wolf and Jürg Wullschleger. Oblivious transfer is symmetric. In *EUROCRYPT*, pages 222–232, 2006.
- [WW10] Severin Winkler and Jürg Wullschleger. On the Efficiency of Classical and Quantum Oblivious Transfer Reductions. In *CRYPTO*, volume 6223, pages 707–723, Berlin, Heidelberg, 2010.

## A Bounds for static (offline) multiple-choice hashing

In our protocol Bob uses a multiple-choice hashing scheme to hash his  $m$  elements into  $n$  buckets. In general, each element will be hashed into  $k$  out of  $d$  possible locations. In previous work, the [PSZ14, PSSZ15, KKRT16], Bob uses 1-out-of- $k$  cuckoo hashing with a stash. Cuckoo hashing is designed for dynamic (online) hashing, whereas Bob is assumed to know his entire set at the time he makes the hash allocation.

Since the communication complexity of the resulting protocols depends on  $n$  (the number of hash buckets), we seek the minimum value of  $n$  that results in an acceptably low failure probability. Error bounds for multiple-choice hashing schemes have been well studied, and most results are obtained by analyzing the orientability of certain hypergraphs. In the following sections, we review some of the pertinent results.

### A.1 Orienting hypergraphs

Given  $k$  hash functions,  $h_i : \mathcal{U}[n]$ , the problem of hashing  $m$  elements into  $n$  bins such that each bin contains at most  $b$  elements, and each element is hashed using  $d$  functions can be analyzed by translating the problem into one of orienting hypergraphs.

Formally, we have,

**Definition 8** (Orientability of hypergraphs). *A  $k$ -uniform hypergraph  $H = (V, E)$  is called  $(d, b)$ -orientable if there exists an assignment of each hyperedge  $e \in E$  to exactly  $d$  of its vertices  $v \in e$  such that no vertex is assigned more than  $b$  hyperedges.*

Translating from our original context, hash-bucket will correspond to a vertex in the hypergraph (thus there are  $n$  vertices), and each element will correspond to an hyper-edge (thus there are  $m$  hyper-edges). The hyper-edge corresponding to an element  $x \in \mathcal{U}$  will contain the  $k$  vertices corresponding to  $h_1(x), \dots, h_k(x)$ . An  $(d, b)$  orientation of this hypergraph then corresponds to choosing a set of  $d$  hash buckets (vertices) for each element (hyper-edge) such that each vertex (bucket) is chosen at most  $b$  times.

The problem of orienting random hypergraphs is well-studied, and many deep results are known [KL02, GW10, Lel12, LP14].

[GW10] give asymptotic thresholds for  $(d, b)$ -orientability in terms of  $n$  and  $m$ , but their results only hold when  $b$  is “sufficiently large.”

[AP11] show how a similar 2-out-of-3 hash structure can be used to compute set intersections in GPUs. Their work focuses on the online setting, and they show the naive insertion algorithm has expected constant running time when  $n = O(m)$ .

[Lel12] gives asymptotic thresholds for  $(d, b)$ -orientability in terms of  $n$  and  $m$ , for almost all values of  $d, b$ . The main result of [Lel12] is the following: for any positive integers,  $k, d, b$  there is an explicit (although complicated)  $c^* = c_{k,d,b}$  such that if  $m > c^*n$  then the probability a random  $n, m, k$  hypergraph is  $(d, b)$ -orientable tends to 0 as  $n \rightarrow \infty$ , and if  $m < c^*n$ , this probability tends to 1.

**Theorem 1** (Theorem 1 [Lel12]). *For integers  $k > d \geq 1$ , and  $b \geq 1$ , then if  $\xi$  is the unique solution to*

$$kd = \xi \frac{E[\max(d - \text{Bin}(k, 1 - Q(\xi, b)), 0)]}{Q(\xi, b + 1) \Pr[\text{Bin}(k - 1, 1 - Q(\xi, b)) < d]}$$

and

$$c_{k,d,b} = \frac{\xi}{k \Pr[\text{Bin}(k - 1, 1 - Q(\xi, b)) < \ell]}$$

where

$$Q(x, y) \stackrel{\text{def}}{=} e^{-x} \sum_{j=y}^{\infty} \frac{x^j}{j!}$$

and  $\text{Bin}(n, p)$  is the binomial distribution with  $\Pr[\text{Bin}(n, p) = k] = \binom{n}{k} p^k (1 - p)^{n-k}$  then

$$\lim_{n \rightarrow \infty} \Pr[H_{n,m,k} \text{ is } (d, b)\text{-orientable}] = \begin{cases} 0 & \text{if } m > c_{k,d,b}n \\ 1 & \text{if } m < c_{k,d,b}n \end{cases}$$

In the case  $d = k - 1$  and  $b = 1$ , we have  $c_{k,k-1,1} = \frac{1}{k(k-1)}$ .

## A.2 Finding an optimal allocation

Given  $2k - 1$  hash functions,  $h_i : \mathcal{U} \rightarrow [n]$ , and  $m$ , we would like to find an allocation such that each of the  $m$  values is hashed into  $k$  of its possible locations, and each bucket contains only a single value. As described above, this corresponds to finding a  $(k, 1)$  orientation on a  $(2k - 1)$ -regular hypergraph  $(V, E)$ .

If such an allocation exists, it can be found efficiently using a max-flow algorithm.

| $k$ | $d$ | $b$ | $c_{k,d,b}$ |
|-----|-----|-----|-------------|
| 3   | 2   | 1   | .1666       |
| 5   | 3   | 1   | .2119       |
| 7   | 4   | 1   | .1747       |
| 9   | 5   | 1   | .1453       |
| 11  | 6   | 1   | .1236       |
| 13  | 7   | 1   | .1074       |
| 15  | 8   | 1   | .0948       |

Table 2: The threshold for orientability in random hypergraphs. If  $m$  is the number of elements, and  $n$  is the number of buckets, a random hashing scheme that hashes each element into  $d$  out of  $k$  buckets, where each bucket has size  $b$  will succeed with probability approaching one if and only if  $m/n < c_{k,d,b}$ . In particular, if we hash each element into 3 out of 5 buckets, then (with high probability) there will be no collisions as long as  $n > 5m$ .

| $k$ | $d$ | $b$ | $c_{k,d,b}$ |
|-----|-----|-----|-------------|
| 2   | 1   | 1   | .4999       |
| 3   | 1   | 1   | .9179       |
| 4   | 1   | 1   | .9768       |
| 5   | 1   | 1   | .9924       |
| 6   | 1   | 1   | .9973       |

Table 3: The threshold for 1-out-of- $k$  orientability in random hypergraphs. If  $m$  is the number of elements, and  $n$  is the number of buckets, a random hashing scheme that hashes each element into  $d$  out of  $k$  buckets, where each bucket has size  $b$  will succeed with probability approaching one if and only if  $m/n < c_{k,d,b}$ . In particular, if we hash each element into 1 out of 3 buckets, then (with high probability) there will be no collisions as long as  $n > 1.09 \cdot m$ .



**Theorem 2.** *Given  $2k - 1$  hash functions,  $h_i : \mathcal{U} \rightarrow [n]$ , and  $\{x_i\}_{i=1}^m \subset \mathcal{U}$ , it is possible to find an allocation  $A : [m] \rightarrow \mathcal{P}([2k - 1])$  that allocates each element  $k$  hash functions, satisfying*

1. *Each element is allocated exactly  $k$  hash functions, i.e.,  $|A(i)| = k$  for all  $i \in [m]$ .*
2. *Each bucket has at most one element, i.e., for all  $i \neq i' \in [m]$ ,  $\{h_j(x_i)\}_{j \in A(i)}$  and  $\{h_j(x_{i'})\}_{j \in A(i')}$  are disjoint.*

*and the allocation can be found in time  $O(nk)$ .*

*Proof.* We begin by translating the problem to finding a  $(k, 1)$ -orientation in an  $(2k - 1)$ -regular hypergraph  $G_{\text{hyper}} = (V_{\text{hyper}}, H_{\text{hyper}})$ , with  $n$  vertices and  $m$  edges. As described above, each bucket corresponds to a vertex in the hypergraph, and each element  $x \in \mathcal{U}$  corresponds to a hyperedge, containing the  $2k - 1$  vertices  $\{h_1(x), \dots, h_{2k-1}(x)\}$ .

Finding a  $(k, 1)$ -orientation of this hypergraph corresponds to finding an orientation, where each hyper-edge “points” to  $k$  of its vertices, and each vertex has in-degree at most 1.

This problem can be translated to a problem on a flow network as follows.

Consider a network with node set with nodes for each vertex and each edge in the hypergraph, thus the set of nodes is  $H_{\text{hyper}} \cup V_{\text{hyper}} \cup \{s, t\}$ . For each edge in the hypergraph, put a capacity 1 edge connected the node corresponding to that edge with the  $2k - 1$  nodes corresponding to the vertices in  $V_{\text{hyper}}$  that it contains. Add capacity  $k$  edges from the source  $s$  to each vertex corresponding to an edge in  $H_{\text{hyper}}$ , and capacity 1 edges from each vertex corresponding to a vertex in  $V_{\text{hyper}}$  to the sink,  $t$ .

Now, note that any  $|H_{\text{hyper}}|(k)$  flow from  $s$  to  $t$  in the flow network induces a  $(k, 1)$  orientation in the hypergraph  $G_{\text{hyper}}$ . Conversely, any  $(k, 1)$  orientation of the hypergraph corresponds exactly to an  $|H_{\text{hyper}}|(k)$  flow from  $s$  to  $t$  in the induced flow network.

Using the Ford-Fulkerson algorithm, the max flow can be computed in time  $O((|V_{\text{hyper}}| + |H_{\text{hyper}}|)(k)) = O((m + n)k) = O(nk)$ .  $\square$

## B OPRFs based on OT-masking

The works of [PSZ14, PSSZ15] outlined a method for constructing a Private Set Membership (PSM) protocol using an OPRF, and then instantiate the OPRF using OT-based masking.

If Bob has an input  $x \in [N]^t$ , and Alice has set of inputs  $Y \subset [N]^t$ , OT-masking protocol works as follows:

- Alice and Bob engage in  $t$  parallel invocations of  $\binom{N}{1}$ -ROT, with the digits of Bob’s input acting as the selection string.
- Thus Alice receives  $2t$  random “masks”,  $\{M_b^i\}$  for  $i = 1, \dots, t$ , and  $b \in \{0, 1\}$ .

- Bob receives  $t$  random masks  $\{M_{x_i}^i\}$  for  $i = 1, \dots, t$ .
- Bob computes his “PRF output”

$$S_B \stackrel{\text{def}}{=} \bigoplus_{i=1}^t M_{x_i}^i$$

as the XOR of his  $t$  masks.

- For each  $y$  in Alice’s set, Alice computes the “PRF”

$$f(y) \stackrel{\text{def}}{=} \bigoplus_{i=1}^t M_{y_i}^i$$

i.e., Alice uses the  $t$  digits of  $y$  to select  $t$  of the masks, and XORs them together.

- Alice sends these  $m$  composite masks to Bob.
- Bob checks if his mask  $S_B$  is in the set received from Alice.

Although this “OT-masking” is secure as a “private equality test” (where Alice has only 1 input), it does not securely realize the one-time OPRF functionality (or the PSM functionality). In particular, the masks on Alice’s elements are not independent, and thus given masks for  $y$  and  $y'$ , Bob can check whether  $y \oplus y'$  is in Alice’s set.

This error appears to have been fixed in the full-version [PSZ16], where the outline a true OT-based one-time OPRF by having Alice and Bob engage in a single  $\binom{N^t}{1}$ -ROT, where Bob’s input is his value  $x \in [N]^t$ , and Alice essentially receives the “truth-table” of a truly random function (with domain  $[N]^t$ ).

## C Log-log-linear MPC-friendly PSI via Hashing and SCS

In this section, we outline a second PSI protocol which allows participants to compute a *secret-sharing* of the intersection, rather than revealing the intersection in the clear. Our protocol uses generic MPC techniques and achieves  $O(mt \log \log m)$  communication for  $t$ -bit values. This protocol combines a basic hashing scheme with the Sort-Compare-Shuffle (SCS) protocol [HEK12] to compare elements within each bucket. This provides asymptotic communication improvements over existing schemes, using only generic MPC techniques, and is naturally composable with larger secure computations.

The generic one-time OPRF protocol (Section 5.1) and its instantiation with BaRKs (Section 5.3) are very efficient, but they reveal the intersection to the players. In many situations, however, PSI is used as a sub-protocol in a larger secure computation, and the intersection itself should never be revealed (e.g. in secure database-joins). In these situations, the protocols outlined in Section 5.1 are insufficient and an “MPC-friendly” PSI protocol, like the one outlined in this section, is necessary.

All existing “MPC-friendly” PSI protocols are significantly less efficient than protocols that reveal the intersection in the clear (like the protocols of [PSZ14, PSSZ15, KKRT16] and our protocol in Section 5.1). It is an interesting open problem whether “MPC-friendly” PSI protocols can be made as efficient as PSI protocols that are not “MPC-friendly.”

Two concurrent, independent works also addressed the question of building PSI protocols that don’t output the intersection in the clear, and thus are suitable for combination within larger MPC protocols. The works of [CO18, PSWW18] both build on the hashing-based PSI protocols of [PSZ14, PSSZ15, PSZ16] to obtain an “MPC-friendly” PSI protocol with  $\omega(m\lambda)$  communication. Their solutions differ from each other, and from our solution which achieves  $O(mt \log \log m)$  communication (where  $t$  is the bit-length of the strings).

## C.1 Construction

Our construction is based on a simple hashing-based scheme, but we use a small number of hash buckets. In particular, instead of setting  $n > m$  (the number of buckets is larger than the number of elements being hashed) we set  $n = m/b$ , for some  $b > 1$ , and allow each bucket to have  $(1 + \delta)b$  elements.

The players will then use the generic, circuit-based, Sort-Compare-Shuffle (SCS) paradigm of [HEK12] to compare the elements within each of the buckets. Somewhat surprisingly, for appropriate choices of  $b$  and  $\delta$  this yields both asymptotic and concrete improvements in communication cost over the basic SCS paradigm of [HEK12].

1. Set  $n = m/b$ , and let  $h : \mathcal{U} \rightarrow [n]$  be a hash function
2. Alice and Bob will each hash all of their elements to  $n$  buckets, using the hash function,  $h$ . If any bucket contains more than  $(1 + \delta)b$  elements, Alice and Bob abort the protocol. Since the abort leaks information, we will show that this happens with probability that is negligible (in  $m$ ).
3. For each bucket, Alice and Bob will engage in a 2-party secure computation, implementing the sort-compare-shuffle (SCS) protocol of [HEK12].

Figure 11: Combining Sort-Compare-Shuffle and hashing to reduce the asymptotic communication cost.

## C.2 Security

The Sort-Compare-Shuffle protocol is implemented in a generic secure multi-party setting such as ABY [DSZ15]. As such, no information is leaked apart

from the function result. If there is no abort, the only information revealed is the intersect.

Now, let us examine the probability that the protocol is aborted, *i.e.*, the probability that a bucket contains more than  $(1 + \delta)b$  elements. Fix a player, and a bucket, and let  $X_i$  denote the random variable that is 1 if the player's  $i$ th element lands in the chosen bucket, and 0 otherwise. Let  $X = \sum_{i=1}^m X_i$  denote the number of elements that land in the given bucket. If we model  $h$  as a truly random hash function, each  $X_i$  is an independent random variable with  $\Pr[X_i = 1] = E[X_i] = \frac{1}{n}$ . Then,  $E[X] = b$ , and by a Chernoff Bound,

$$\Pr[X > (1 + \delta)b] \leq e^{-\frac{\delta^2 b}{3}}$$

Taking a union bound over the 2 players and the  $n$  buckets, we find that the probability of abort (and hence information leakage) is upper bounded by

$$\frac{2me^{-\frac{\delta^2 b}{3}}}{b} < 2^{-\frac{\delta^2 b}{3} + \log m}$$

Asymptotically, setting  $\delta = 1$ , and  $b = \log^2 m$ , we have that the failure probability is negligible in  $m$ . Since the SCS protocol requires  $O(tb \log b)$  AND gates to compare each bucket of size  $O(b)$ , the total number of AND gates is  $O(mt \log \log m)$ . Using the GMW protocol [GMW87], each AND gate can be implemented using 2 OTs. Using OT-extensions, the entire secure two-party computation can be implemented using computation and communication that is linear in the size of the circuit [IKOS08].

Concretely, to keep the overall failure probability below  $2^{-\sigma}$ , we set

$$\begin{aligned} \frac{2me^{-\frac{\delta^2 b}{3}}}{b} &< 2^{-\frac{\delta^2 b}{3} + \log m} < 2^{-\sigma} \\ &\Rightarrow \log m - \frac{\delta^2 b}{3} < -\sigma \\ &\Rightarrow \frac{3(\log m + \sigma)}{\delta^2} < b \end{aligned}$$

### C.3 Communication Complexity

As previously mentioned, the SCS protocol is implemented in a generic secure multiparty computation framework. Since SCS involves 2 parties and consists of a large number of comparisons, it is natural to implement this protocol in a framework that uses Garbled Circuits as the backend such as ABY [DSZ15]. In this setting, the communication cost primarily depends on the number of AND gates in the circuit.

We assume that Alice and Bob's sets can be represented using  $t$  bits. If we use permutation-based hashing (as in [PSSZ15]), the elements can be represented in buckets using  $t - \log n$  bits. Using the SCS protocol of [HEK12], to compare

$\ell$  elements of length  $\tau$ , requires

$$\begin{aligned}
& \frac{5}{3}\tau\ell \log \ell + 2\tau\ell + ((3\ell - 1)\tau - \ell) - \frac{\tau\ell + \tau}{3} \\
&= \frac{5}{3}\tau\ell \log \ell + \frac{14}{3}\tau\ell - \frac{4}{3}\tau - \ell \\
&< \frac{5}{3}\tau\ell \log \ell + \frac{14}{3}\tau\ell \\
&= \frac{\tau\ell}{3} (5 \log \ell + 14)
\end{aligned}$$

AND gates. In our situation, each bucket has  $\ell = (1 + \delta)b$  elements and each element is of length  $\tau = t - \log n$  bits. Thus the calculation requires

$$\begin{aligned}
& \frac{\tau\ell}{3} (5 \log \ell + 14) \\
&= \frac{(t - \log n)(1 + \delta)b}{3} (5 \log((1 + \delta)b) + 14) \\
&= \frac{(t - \log m + \log b)(1 + \delta)b}{3} (5 \log(1 + \delta) + 5 \log b + 14)
\end{aligned}$$

AND gates. For given values of number of elements,  $m$ , bit-length,  $t$ , and a given error threshold (e.g.  $\sigma = 40$ ), we can find the value of  $\delta$  to minimize the total number of AND gates required. Table 4 summarizes the number of AND gates per element required to implement our hashing-SCS scheme (with  $\sigma = 40$ , and gives optimal choices of the parameters  $\delta$  and  $b$ ).

Thus for real-world parameter choices, the entire PSI protocol can be computed by a circuit using approximately  $500m$  AND gates, and this circuit can be evaluated with constant overhead (independent of the security parameter) using the generic techniques of [IKOS08].

#### C.4 Comparison with concurrent MPC-friendly PSI protocols

Two concurrent, independent works, [PSWW18, CO18], developed MPC-friendly PSI protocols. We implemented our protocol of Figure 11 using the ABY [DSZ15] framework, and obtained concrete communication metrics. In this section, we compare these performance numbers with those found in [PSWW18] and [CO18].

We find that our protocol has somewhat higher concrete communication complexity than that of [PSWW18].

Figure 5 shows the number secure AND gates required to compute a PSI using our protocol compared to that of [PSWW18]. The number of AND gates required by the protocol is a good metric for the complexity of the protocol because the number of AND gates is the primary driver of performance of all circuit-based MPC protocols, but unlike communication or runtime benchmarks, it is

| $\log(m)$ | $t$ | $\delta$ | $b$ | Number of AND gates per element |
|-----------|-----|----------|-----|---------------------------------|
| 8         | 12  | 0.88     | 130 | 370                             |
| 8         | 16  | 0.62     | 263 | 499                             |
| 8         | 20  | 0.52     | 378 | 621                             |
| 12        | 16  | 0.85     | 150 | 378                             |
| 12        | 20  | 0.61     | 295 | 507                             |
| 12        | 24  | 0.51     | 421 | 631                             |
| 16        | 20  | 0.83     | 169 | 386                             |
| 16        | 24  | 0.60     | 327 | 515                             |
| 16        | 28  | 0.50     | 465 | 640                             |
| 20        | 24  | 0.81     | 190 | 392                             |
| 20        | 28  | 0.59     | 361 | 523                             |
| 24        | 28  | 0.80     | 211 | 399                             |
| 24        | 32  | 0.58     | 395 | 530                             |
| 28        | 32  | 0.78     | 233 | 405                             |
| 28        | 36  | 0.58     | 430 | 537                             |
| 32        | 36  | 0.77     | 255 | 411                             |

Table 4: Number of AND gates required per element in the hashing-SCS protocol.  $m$  is the total number of elements,  $t$  is the bit length of each element,  $b$  is the expected bucket size, and  $(1 + \delta)b$  is the maximum allotted bucket size. Thus the overall computation requires  $m/b$  SCS sub-computations, each on sets of size  $(1 + \delta)b$ .

| $\log(n)$ | $t$      | Number of AND gates per element |          |                |
|-----------|----------|---------------------------------|----------|----------------|
|           |          | Ours                            | [PSWW18] | BWA<br>[HEK12] |
| 12        | 32       | 871                             | 406      | $2^{20}$       |
| 16        | 32       | 761                             | 306      | $2^{16}$       |
| 20        | 32       | 648                             | 205      | $2^{12}$       |
| 12        | $\infty$ | 1780                            | 921      | $2^{51}$       |
| 16        | $\infty$ | 1915                            | 875      | $2^{55}$       |
| 20        | $\infty$ | 2050                            | 797      | $2^{59}$       |

Table 5: Comparing the number of AND gates per element using the best (“Iterative Combined”) scheme of [PSWW18].  $n$  is the number of elements, and  $t$  is the bit length of each element. When  $t = \infty$ , we use the same method of [PSWW18] to first hash each element to a string of length  $40 + 2 \log_2(n) - 1$ , which ensures that the collision probability remains below  $2^{-40}$ .

agnostic to the exact MPC implementation. We find that our protocol requires 2-3 times as many AND gates as that of [PSWW18]. The protocol of [CO18] is not compatible with a generic, circuit-based comparison protocol, and thus it is not included in Table 5.

In Table 6 we compare the communication cost of our protocol with those of [PSWW18] and [CO18]. The work of [CO18] does not include an implementation, but the concrete communication cost of their protocol can be calculated analytically. Since the protocol of [CO18] primarily requires sending AES ciphertexts, we can calculate the exact communication cost of the protocol by calculating the maximum number of elements in each hash bin, and using the communication costs provided in [CO18]. This provides a slight under-estimate of the communication cost of [CO18] since any actual implementation would require sending some additional meta-information (e.g. packet headers, error-correction, etc).

We find that the concrete communication cost of our protocol is 10-20 times higher than that of [PSWW18] and on par with that of [CO18]. All three protocols require 100-1000 times more communication than a protocol like [KKRT16] (or our protocol from Section 5.1) that reveals the intersection in the clear to both parties.

| $\log(n)$ | Communication Cost (mb) |          |        |                   |
|-----------|-------------------------|----------|--------|-------------------|
|           | MPC-friendly            |          |        | Not MPC- friendly |
|           | Ours                    | [PSWW18] | [CO18] | [KKRT16]          |
| 12        | 611                     | 52       | 507    | 0.53              |
| 16        | 7730                    | 638      | 6864   | 8                 |
| 20        | 121411                  | 6950     | 109824 | 127               |

Table 6: The total communication cost in mb of our MPC-friendly protocol compared to the concurrent MPC-friendly works of [PSWW18] and [CO18] as well as the not MPC-friendly work of [KKRT16]. For the [CO18] communication, we use AES128 as the cipher, and maximum bins of size 7, 7, 8 (which yields a failure probability bounded by  $2^{-40}$ ). Notice that the best MPC-friendly PSI protocols are orders of magnitude more communication intensive than those that reveal the intersection in the clear (like [KKRT16]). All costs are for elements represented by 32 bits (*i.e.*,  $t = 32$ ).