

# Efficient Circuit-based PSI via Cuckoo Hashing

## (Full Version)\*

Benny Pinkas<sup>1</sup>, Thomas Schneider<sup>2</sup>, Christian Weinert<sup>2</sup>, and Udi Wieder<sup>3</sup>

<sup>1</sup> Bar-Ilan University

benny@pinkas.net

<sup>2</sup> TU Darmstadt

{thomas.schneider,christian.weinert}@crisp-da.de

<sup>3</sup> VMware Research

udi.wieder@gmail.com

**Abstract.** While there has been a lot of progress in designing efficient custom protocols for computing Private Set Intersection (PSI), there has been less research on using generic Multi-Party Computation (MPC) protocols for this task. However, there are many variants of the set intersection functionality that are not addressed by the existing custom PSI solutions and are easy to compute with generic MPC protocols (e.g., comparing the cardinality of the intersection with a threshold or measuring ad conversion rates).

Generic PSI protocols work over circuits that compute the intersection. For sets of size  $n$ , the best known circuit constructions conduct  $O(n \log n)$  or  $O(n \log n / \log \log n)$  comparisons (Huang et al., NDSS'12 and Pinkas et al., USENIX Security'15). In this work, we propose new circuit-based protocols for computing *variants of the intersection* with an almost linear number of comparisons. Our constructions are based on new variants of Cuckoo hashing in two dimensions.

We present an asymptotically efficient protocol as well as a protocol with better concrete efficiency. For the latter protocol, we determine the required sizes of tables and circuits experimentally, and show that the run-time is concretely better than that of existing constructions.

The protocol can be extended to a larger number of parties. The proof technique for analyzing Cuckoo hashing in two dimensions is new and can be generalized to analyzing standard Cuckoo hashing as well as other new variants of it.

**Keywords:** Private Set Intersection, Secure Computation

## 1 Introduction

Private Set Intersection (PSI) refers to a protocol which enables two parties, holding respective input sets  $X$  and  $Y$ , to compute the intersection  $X \cap Y$  without revealing any information about the items which are not in the intersection. The PSI functionality is useful for applications where parties need to apply a JOIN operation to private datasets. There are multiple constructions of secure protocols for computing PSI, but there is an advantage for computing PSI by applying a generic Multi-Party Computation (MPC) protocol to a circuit computing the intersection (see §1.1). The problem is that a naive circuit computes  $O(n^2)$  comparisons, and even the most recent circuit-based constructions require  $O(n \log n)$  or  $O(n \log n / \log \log n)$  comparisons (see §1.4).

In this work, we present a new circuit-based protocol for computing PSI variants. In our protocol, each party first inserts its input elements into bins according to a new hashing algorithm, and then the intersection is computed by securely computing a Boolean comparison circuit over the bins. The insertion of the items is based on new Cuckoo hashing variants which guarantee that if the two

---

\* Please cite the conference version of this work published at EUROCRYPT'18 [PSWW18].

parties have the same input value, then there is exactly one bin to which both parties map this value. Furthermore, the total number of bins is  $O(n)$  and there are  $O(1)$  items mapped to each bin, plus  $\omega(1)$  items which are mapped to a special stash. Hence, the circuit that compares (1) for each bin, the items that the two parties mapped to it, and (2) all stash items to all items of the other party, computes only  $\omega(n)$  comparisons.

## 1.1 Motivation for Circuit-based PSI

PSI has many applications, as is detailed for example in [PSZ18]. Consequently, there has been a lot of research on efficient secure computation of PSI, as we describe in §1.4. However, most research was focused on computing the intersection itself, while there are interesting applications for the ability to securely compute arbitrary functions of the intersection. We demonstrate the need for efficient computation of PSI using generic protocols through the following arguments:

**Adaptability.** Assume that you are a cryptographer and were asked to propose and implement a protocol for computing PSI. One approach is to use a specialized protocol for computing PSI. Another possible approach is to use a protocol for generic secure computation, and apply it to a circuit that computes PSI. A trivial circuit performs  $O(n^2)$  comparisons, while more efficient circuits, described in [HEK12] and [PSSZ15], perform only  $O(n \log n)$  or  $O(n \log n / \log \log n)$  comparisons, respectively. The most efficient specialized PSI protocols are faster by about two orders of magnitude than circuit-based constructions (see [PSSZ15]), and therefore you will probably choose to use a specialized PSI protocol. However, what happens if you are later asked to change the protocol to compute another function of the intersection? For example, output only the size of the intersection, or output 1 iff the size is greater than some threshold, or output the most “representative” item that occurs in the intersection (according to some metric). Any change to a specialized protocol will require considerable cryptographic know-how, and might not even be possible. On the other hand, the task of writing a new circuit component that computes a different function of the intersection is rather trivial, and can even be performed by undergrad students.

Consider the following function as an example of a variant of the PSI functionality for which we do not know a specialized protocol: Suppose that you want to compute the size of the intersection, but you also wish to preserve the privacy of users by ensuring differential privacy. This is done by adding some noise to the exact count before releasing it. This functionality can easily be computed by a circuit, but it is unclear how to compute it using other PSI protocols. (See [PL15] for constructions that add noise to the results of MPC computation in order to ensure differential privacy.)

**Existing code base.** Circuit-based protocols benefit from all the work that was invested in recent years in designing, implementing, and optimizing very efficient systems for generic secure computation. Users can download existing secure computation software, e.g., [HEKM11, EFL12, DSZ15], and only need to design the circuit to be computed and implement the appropriate hashing technique.

**Existing applications.** There are existing applications that need to compute functions over the results of the set intersection. For example, Google reported [Yun15, Kre17] a PSI-based application for measuring ad conversion rates, namely the revenues from ad viewers who later perform a related transaction. This computation can be done by comparing the list of people who have seen an ad with those who have completed a transaction. These lists are held by the advertiser (say, Google or Facebook), and by merchants, respectively. A simple (non-private) solution is for one side to disclose its list of customers to the other side, which computes the necessary statistics. Another option is to run a secure computation over the results of the set intersection. For example, the merchant inputs pairs of the customer-identity and the value of the transactions made by this

customer, and the computation calculates the total revenue from customers who have seen an ad, namely customers in the intersection of the sets known to the advertiser and the merchant. Google reported implementing this computation using a Diffie-Hellman-based PSI cardinality protocol (for computing the cardinality of the intersection) and Paillier encryption (for computing the total revenues) [IKN<sup>+</sup>17]. This protocol additionally reveals the number of items in the intersection, and seems less efficient than our protocol as it uses public key operations, rather than efficient symmetric cryptographic operations.<sup>4</sup>

## 1.2 Our Contributions

This work provides the following contributions:

**Circuit-based PSI protocols with almost linear overhead.** We show a new circuit-based construction for computing any symmetric function on top of PSI, with an asymptotic overhead of only  $\omega(n)$  comparisons. (More accurately, for any function  $f \in \omega(n)$ , the overhead of the construction is  $o(f(n))$ .) This construction is based on standard Cuckoo hashing.

**Small constants.** Standard measures of asymptotic security are not always a good reflection of the actual performance on reasonable parameters. Therefore, in addition to the asymptotic improvement, we also show a concrete circuit-based PSI construction. This construction is based on a new variant of Cuckoo hashing, *two-dimensional Cuckoo hashing*, that we introduce in this work. We carefully handle implementation issues to improve the actual overhead of our protocols, and make sure that all constants are small. In particular, we ran extensive experiments to analyze the failure probabilities of the hashing scheme, and find the exact parameters that reduce this statistical failure probability to an acceptable level (e.g.,  $2^{-40}$ ). Our analysis of the concrete complexities is backed by extensive experiments, which consumed about 5.5 million core hours on the Lichtenberg high performance computer of the TU Darmstadt and were used to set the parameters of the hashing scheme. Given these parameters we implemented the circuit-based PSI protocol and tested it.

**Implementation and experiments.** We implemented our protocols using the ABY framework for secure two-party computation [DSZ15]. Our experiments show that our protocols are considerably faster than the previously best circuit-based constructions. For example, for input sets of  $n = 2^{20}$  elements of arbitrary bitlength, we improve the circuit size over the best previous construction by up to a factor of 3.8x.

**New Cuckoo hashing analysis.** Our two-dimensional Cuckoo hashing is based on a new Cuckoo hashing scheme that employs two tables and each item is mapped to either *two* locations in the first table, or *two* locations in the second table. This is a new Cuckoo hashing variant that has not been analyzed before. In addition to measuring its performance using simulations, we provide a probabilistic analysis of its performance. Interestingly, this analysis can also be used as a new proof technique for the success probability of standard Cuckoo hashing.

## 1.3 Computing Symmetric Functions

A trivial circuit for PSI that performs  $O(n^2)$  comparisons between all pairs of the input items of the two parties allows the parties to set their inputs in any arbitrary order. On the other hand,

<sup>4</sup> Facebook is running a computation of this type with companies that have transaction records for a large part of loyalty card holders in the US. According to the report in <https://www.eff.org/deeplinks/2012/09/deep-div-e-facebook-and-datalogix-whats-actually-getting-shared-and-how-you-can-opt>, the computation is done using an insecure PSI variant based on creating pseudonyms using naive hashing of the items.

there exist more efficient circuit-based PSI constructions where each party first independently orders its inputs according to some predefined algorithm: the sorting network-based construction of [HEK12] requires each party to sort its input to the circuit, while the hashing-based construction of [PSSZ15] requires the parties to map their inputs to bins using some public hash functions. (These constructions are described in §1.4.) The location of each input item thus depends on the identity of the other inputs of the input owner, and must therefore be kept hidden from the other party.

In this work, we focus on constructing a circuit that computes the intersection. The outputs of this circuit can be the items in the intersection, or some functions of the items in the intersection: say, a “1” for each intersecting item, or an arbitrary function of some data associated with the item (for example, if the items are transactions, we might want to output a financial value associated with each transaction that appears in the intersection). On top of that circuit it is possible to add circuits for computing any function that is based on the intersection. In order to preserve privacy, the output of that function must be a *symmetric* function of the items in the intersection. Namely, the output of the function must not depend on the *order* of its inputs. There are many examples of interesting symmetric functions of the intersection. (In fact, it is hard to come up with examples for interesting non-symmetric functions of the intersection, except for the intersection itself.) Examples of symmetric functions include:

- Computing the size of the intersection, i.e., PSI cardinality (PSI-CA).
- Computing a threshold function that is based on the size of the intersection. For example, outputting “1” if the size of the intersection is greater than some threshold (PSI-CAT), or outputting a rounded value of the percentage of items that are in the intersection. An extension of PSI-CAT, where the intersection is revealed only if the size of the intersection is greater than a threshold, can be used for privacy-preserving ridesharing [HOS17].
- Computing the size of the intersection while preserving the privacy of users by ensuring differential privacy [Dwo06]. This can be done by adding some noise to the exact count.
- Computing the sum of values associated with the items in the intersection. This is used for measuring ad-generated revenue (cf. §1.1). Similarly, there could be settings where each party associates a value with each transaction, and the output is the sum of the differences between these assigned values in the intersection, or the sum of the squares of the differences, etc.

The circuits for computing all these functions are of size  $O(n)$ . Therefore, with our new construction, the total size of the circuits for computing these functions is  $\omega(n)$ , whereas circuit-based PSI protocols [HEK12, PSSZ15] had size  $O(n \log n)$ .

If one wishes to compute a function that is not symmetric, or wishes to output the intersection itself, then the circuit must first shuffle the values in the intersection (in order to assign a random location to each item in the intersection) and then compute the function over the shuffled values, or output the shuffled intersection. A circuit for this “shuffle” step has size  $O(n \log n)$ , as described in [HEK12]. (It is unclear, though, why a circuit-based protocol should be used for computing the intersection, since this job can be done much more efficiently by specialized protocols, e.g., [KKRT16, PSZ18].)

## 1.4 Related Work

**PSI.** Work on protocols for private set intersection was presented as early as [Sha80, Mea86], which introduced public key-based protocols using commutative cryptography, namely the Diffie-Hellman function. A survey of PSI protocols appears in [PSZ14]. The goal of these protocols is to let one party

learn the intersection itself, rather than to enable the secure computation of arbitrary functions of the intersection. Other PSI protocols are based on oblivious polynomial evaluation [FNP04], blind RSA [DT10], and Bloom filters [DCW13]. Today’s most efficient PSI protocols are based on hashing the items to bins and then evaluating an oblivious pseudo-random function per bin, which is implemented using oblivious transfer (OT) extension. These protocols have linear complexity and were all implemented and evaluated, see, e.g., [PSZ14, PSSZ15, KKRT16, PSZ18]. In cases where communication cost is a crucial and computation cost is a minor factor, recent solutions based on fully homomorphic encryption represent an interesting alternative [CLR17]. PSI protocols have also been adapted to the special requirements of mobile devices [HCE11, ADN<sup>+</sup>13, KLS<sup>+</sup>17].

**Circuit-based PSI.** Circuit-based PSI protocols compute the set intersection functionality by running a secure evaluation of a Boolean circuit. These protocols can easily be adapted to compute different variants of the PSI functionality. The straightforward solution to the PSI problem requires  $O(n^2)$  comparisons – one comparison for each pair of items belonging to the two parties. Huang et al. [HEK12] designed a circuit for computing PSI based on sorting networks, which computes  $O(n \log n)$  comparisons and is of size  $O(\sigma n \log n)$ , where  $\sigma$  is the bitlength of the inputs. A different circuit, based on the usage of Cuckoo hashing by one party and simple hashing by the other party, was proposed in [PSSZ15]. The size of that circuit is  $O(\sigma n \log n / \log \log n)$ . In our work we propose efficient circuits for PSI variants with an asymptotic size of  $\omega(\sigma n)$  and better concrete efficiency. We give more details and a comparison of the concrete complexities of circuit-based PSI protocols in §6.2.

**PSI Cardinality (PSI-CA).** A specific interesting function of the intersection is its cardinality, namely  $|X \cap Y|$ , and is referred to as PSI-CA. There are several protocols for computing PSI-CA with linear complexity based on public key cryptography, e.g., [DGT12] which is based on Diffie-Hellman and is essentially a variant of the DH-based PSI protocol of [Sha80, Mea86] (see also references given therein for other less efficient public key-based protocols); or [DD15] which is based on Bloom filters and the public key cryptosystem of Goldwasser-Micali; [EFG<sup>+</sup>15] additionally allows one to compute the private set union cardinality using also Bloom filters together with ElGamal encryption. In these protocols, one of the parties learns the cardinality. As we show in our experiments in §6.3, these protocols are slower than our constructions already for relatively small set sizes ( $n = 2^{12}$ ) in the LAN setting and for large set sizes ( $n = 2^{20}$ ) in the WAN setting, since they are based on public key cryptography. An advantage of these protocols is that they achieve the lowest amount of communication, but it seems hard to extend them to compute arbitrary functions of the intersection. Protocols for private set intersection and union and their cardinalities with linear complexity are given in [DC17]. They use Bloom filters and computationally expensive additively homomorphic encryption, whereas our protocols can flexibly be adapted to different variants and are based on efficient symmetric key cryptography.

## 2 Preliminaries

**Setting.** We consider two parties, which we denote as Alice and Bob. They have input sets,  $X$  and  $Y$ , respectively, which are each of size  $n$  and each item has bitlength  $\sigma$ . We assume that both parties agree on a symmetric function  $f$  and would like to securely compute  $f(X \cap Y)$ . They also agree on a circuit that receives the items in the intersection as input and computes  $f$ .

**Security Model.** The secure computation literature considers *semi-honest* adversaries, which try to learn as much information as possible from a given protocol execution, but are not able to deviate from the protocol steps, and *malicious* adversaries, which are able to deviate arbitrarily from the protocol. The semi-honest adversary model is appropriate for scenarios where execution of the intended software is guaranteed via software attestation or business restrictions, and yet an untrusted third party is able to obtain the transcript of the protocol after its execution, either by stealing it or by legally enforcing its disclosure. Most protocols for private set intersection, as well as this work, focus on solutions that are secure against semi-honest adversaries. PSI protocols for the malicious setting exist, but they are less efficient than protocols for the semi-honest setting (see, e.g., [FNP04, HL08, DSMRY09, HN10, DKT10, FHNP16, RR17a, RR17b]).

**Secure Computation.** There are two main approaches for generic secure two-party computation with security against semi-honest adversaries that allow to securely evaluate a function that is represented as a Boolean circuit: (1) Yao’s garbled circuit protocol [Yao86] has a constant round complexity and with today’s most efficient optimizations provides free XOR gates [KS08], whereas securely evaluating an AND gate requires a constant number of fixed-key AES evaluations [BHKR13] and sending two ciphertexts [ZRE15]. (2) The GMW protocol [GMW87] also provides free XOR gates and requires two ciphertexts of communication per AND gate using OT extension [ALSZ13]. The main advantage of the GMW protocol is that *all* symmetric cryptographic operations can be pre-computed in a constant number of rounds in a setup phase, whereas the online phase is very efficient, but requires interaction for each layer of AND gates. In more detail, the setup phase is independent of the actual inputs and pre-computes multiplication triples for each AND gate using OT extension in a constant number of rounds (cf. [ALSZ13]). The online phase runs from the time the inputs are provided until the result is obtained and involves sending one message for each layer of AND gates. A detailed description and a comparison between Yao and GMW is given in [SZ13].

**Cuckoo Hashing.** In its simplest form, Cuckoo hashing [PR01] uses two hash functions  $h_0, h_1$  to map  $n$  elements to two tables  $T_0, T_1$ , each containing  $(1 + \varepsilon)n$  bins. Each bin accommodates at most a single element. The scheme avoids collisions by relocating elements when a collision is found using the following procedure: Let  $b \in \{0, 1\}$ . An element  $x$  is inserted into a bin  $h_b(x)$  in table  $T_b$ . If a prior item  $y$  exists in that bin, it is evicted to bin  $h_{1-b}(y)$  in  $T_{1-b}$ . The pointer  $b$  is then assigned the value  $1 - b$ . The procedure is repeated until no more evictions are necessary, or until a threshold number of relocations has been performed. In the latter case, the last element is mapped to a special stash. It was shown in [KMW09] that, for any constant  $s$ , the probability that the size of the stash is greater than  $s$  is at most  $O(n^{-(s+1)})$ . After inserting all items, each item can be found in one of two locations or in the stash. A lookup therefore requires checking only  $O(1)$  locations.

Many variants of Cuckoo hashing were suggested and analyzed. See [Wie16] for a thorough discussion and analysis of different Cuckoo hashing schemes. A variant of Cuckoo hashing that is similar to our constructions was given in [AP11], although in a different application domain. It considers a setting with three tables, where an item must be placed in two out of three tables. The analysis of this construction uses a different proof technique than the one we present, and we have not attempted to generalize their proof to a general number of item insertions (as we do for our construction in §A). Furthermore, there is no tight analysis of the stash size in [AP11]. The work in [EGMT17] builds on the construction of [AP11] and proves that the failure probability when using a stash of size  $s$  behaves as  $\tilde{O}(n^{-s})$ . However, the experiments of [EGMT17, Fig. 6] reveal that

the size of the stash is rather large and actually *increasing* in  $n$  within the range of 1 000 to 100 000 elements. For example, for table size  $7.1n$ , a stash of at least size 4 is required for inserting 10 000 elements, whereas a stash of at least size 11 is required for inserting 100 000 elements. Since each item in the stash must be compared to all items of the other party, and since these comparisons cannot use a shorter representation based on permutation-based hashing, the effect of the stash is substantial, and in the context of circuit-based PSI it is therefore preferable to use constructions that place very few or no items in the stash.

**PSI based on Hashing.** Some existing constructions of circuits for PSI require the parties to reorder their inputs before inputting them to the circuit: The sorting-network based construction of [HEK12] requires the parties to sort their inputs. The hashing based construction of [PSSZ15] requires that each party maps its items to bins using a hash function. It was observed as early as [FNP04] that if the two parties agree on the same hash function and use it to map their respective input to bins, then the items that one party maps to a specific bin need to be compared only to the items that the other party maps to the same bin. However, the parties must be careful not to reveal to each other the number of items they mapped to each bin, since this data leaks information about their other items. Therefore, they agree beforehand on an upper bound  $m$  for the maximum number of items that can be mapped to a bin (such upper bounds are well known for common hashing algorithms, and can also be substantiated using simulation), and pad each bin with random dummy values until it has exactly  $m$  items in it. If both parties use the same hash algorithm, then this approach considerably reduces the overhead of the computation from  $O(n^2)$  to  $O(\beta \cdot m^2)$ , where  $m$  is the maximum number of items mapped to any of the  $\beta$  bins.

When a random hash function  $h$  is used to map  $n$  items to  $n$  bins, where  $x$  is mapped to bin  $h(x)$ , the most occupied bin has w.h.p.  $m = \frac{\ln n}{\ln \ln n}(1 + o(1))$  items [Gon81] (a careful analysis shows, e.g., that, for  $n = 2^{20}$  and an error probability of  $2^{-40}$ , one needs to set  $m = 20$ ). Cuckoo hashing is much more promising, since it maps  $n$  items to  $2(1 + \varepsilon)n$  bins, where each bin stores at most  $m = 1$  items. Cuckoo hashing typically uses two hash functions  $h_0, h_1$ , where an item  $x$  is mapped to one of the two locations  $h_0(x), h_1(x)$ , or to a stash of a small size. It is tempting to let both parties, Alice and Bob, map their items to bins using Cuckoo hashing, and then only compare the item that one party maps to a bin with the item that the other party maps to the same bin. The problem is that Alice might map  $x$  to  $h_0(x)$  whereas Bob might map it to  $h_1(x)$ . They cannot use a protocol where Alice's value in bin  $h_0(x)$  is compared to the two bins  $h_0(x), h_1(x)$  in Bob's input, since this reveals that Alice has an item that is mapped to these two locations. The solution used in [FHNP16, PSZ14, PSSZ15] is to let Alice map her items to bins using Cuckoo hashing, and Bob map his items using simple hashing. Namely, each item of Bob is mapped to both bins  $h_0(x), h_1(x)$ . Therefore, Bob needs to pad his bins to have  $m = O(\log n / \log \log n)$  items in each bin, and the total number of comparisons is  $O(n \log n / \log \log n)$ .

### 3 Analyzing the Failure Probability

Efficient cryptographic protocols that are based on probabilistic constructions are typically secure as long as the underlying probabilistic constructions do not fail. Our work is based on variants of Cuckoo hashing, and the protocols are secure as long as the relevant tables and stashes do not overflow. (Specifically, hashing is computed using random hash functions which are chosen independently of the data. If a party observes that these functions cannot successfully hash its data,

it can indeed ask to replace the hash functions, or remove some items from its input. However, the hash functions are then no longer independent of this party’s input and might therefore leak some information about the input.)

There are two approaches for arguing about the failure probability of cryptographic protocols:

1. For an **asymptotic analysis**, the failure probability must be negligible in  $n$ .
2. For a **concrete analysis**, the failure probability is set to be smaller than some threshold, say  $2^{-\lambda}$ , where  $\lambda$  is a statistical security parameter.

In typical experiments, the statistical security parameter is set to  $\lambda = 40$ . This means that “unfortunate” events that leak information happen with a probability of at most  $2^{-40}$ . In particular,  $\lambda = 40$  was used in all PSI constructions which are based on hashing (e.g., [DCW13, PSZ14, PSSZ15, FHNP16, KKRT16]).

With regards to the probabilistic constructions, there are different levels of analysis of the failure probability:

1. For simple constructions, it is sometimes possible to compute the **exact failure probability**. (For example, suppose that items are hashed to a table using a random hash function, and a failure happens when two items are mapped to the same location. In this case it is trivial to compute the exact failure probability.)
2. For some constructions there are known **asymptotic bounds** for the failure probability, but no concrete expressions. (For example, for Cuckoo hashing with a stash of size  $s$ , it was shown in [KMW09] that the overflow probability is  $O(n^{-(s+1)})$ , but the exact constants are unknown.)<sup>5</sup>
3. For other constructions there is no analysis for the failure probability, even though they **perform very well in practice**. For example, Cuckoo hashing variants where items can be mapped to  $d > 2$  locations, or where each bin can hold  $k > 1$  items, were known to have better space utilization than standard Cuckoo hashing, but it took several years to theoretically analyze their performance [Wie16]. There are also insertion algorithms for these Cuckoo hashing variants which are known to perform well but which have not yet been fully analyzed.

### 3.1 Using Probabilistic Constructions for Cryptography

Suppose that one is using a probabilistic construction (e.g., a hash table) in the design of a cryptographic protocol. An asymptotic analysis of the cryptographic protocol can be done if the hash table has either an exact analysis or an asymptotic analysis of its failure probability (items 1 and 2 in the previous list).

If the aim is a concrete analysis of the cryptographic protocol, then exact values for the parameters of the hash construction must be identified. If an exact analysis is known (item 1), then it is easy to plug in the desired failure probability ( $2^{-\lambda}$ ) and compute the values for the different parameters. However, if only an asymptotic analysis or experimental evidence is known (items 2 and 3), then experiments must be run in order to find the parameters that set the failure probability to be smaller than  $2^{-\lambda}$ .

We stress that a concrete analysis is needed whenever a cryptographic protocol is to be used in practice. In that case, even an asymptotic analysis is insufficient since it does not specify any constants, which are crucial for deriving the exact parameter values.

---

<sup>5</sup> We note though that many probabilistic constructions are analyzed in the algorithms research literature to have a failure probability of  $o(1)$ , which is fine for many applications, but is typically insufficient for cryptographic applications.



### 3.2 Experimental Parameter Analysis

Verifying that the failure probability is smaller than  $2^{-\lambda}$  for  $\lambda = 40$  requires running many repetitions of the experiments. Furthermore, for large input sizes (large values of  $n$ ), each single run of the experiment can be rather lengthy. (And one could justifiably argue that the more interesting results are for the larger values of  $n$ , since for smaller  $n$  we can use less optimal constructions and still get reasonable performance.)

**Examining the failure probability for a specific choice of parameters.** For a specific choice of parameters, running  $2^\lambda$  repetitions of an experiment is insufficient to argue about a  $2^{-\lambda}$  failure probability, since it might happen that the experiments were very unlucky and resulted in no failure even though the failure probability is somewhat larger than  $2^{-\lambda}$ . Instead, we can argue about a confidence interval: namely, a confidence interval of  $1 - \alpha$  (say, 95%, or 99.9%) states that if the failure probability is greater than  $2^{-\lambda}$ , then we would have *not* seen the results of the experiment, except with a probability that is smaller than  $\alpha$ . Therefore, either the experiment was very unlucky, or the failure probability is sufficiently small. For example, an easy to remember confidence level used in statistics is the “rule of three”, which states that if an event has not occurred in  $3 \cdot s$  experiments, then the 95% confidence interval for its rate of occurrence in the population is  $[0, 1/s]$ . For our purposes this means that running  $3 \cdot 2^\lambda$  experiments with no failure suffices to state that the failure probability is smaller than  $2^{-\lambda}$  with 95% confidence. (We will report experiments in §6.1 which result in a 99.9% confidence interval for the failure probability.)

**Examining the failure probability as a function of  $n$ .** For large values of  $n$  (e.g.,  $n = 2^{20}$ ), it might be too costly to run sufficiently many (more than  $2^{40}$ ) experiments. Suppose that the experiments spend just 10 cycles on each item. This is an extremely small lower bound, which is probably optimistic by orders of magnitude compared to the actual run-time. Then the experiments take at least  $10 \cdot 2^{60}$  cycles. This translates to about a million core hours on 3 GHz machines.

In order to be able to argue about the failure probability for large values of  $n$ , we can run experiments for progressively increasing values of  $n$  and identify how the failure probability behaves as a function of  $n$ . If we observe that the failure probability is decreasing, or, better still, identify the dependence on  $n$ , we can argue, given experimental results for medium-sized  $n$  values, about the failure probabilities for larger values of  $n$ .

### 3.3 Our Constructions

**Asymptotic overhead.** We present in §4 a construction of circuit-based PSI that we denote as the “mirror” construction. This construction uses four instances of standard Cuckoo hashing and therefore we know that a stash of size  $s$  guarantees a failure probability of  $O(n^{-(s+1)})$  [KMW09]. (Actually, the previously known analysis was only stated for  $s = O(1)$ . We show in §C that this failure probability also holds for  $s$  that is not constant.)

The bound on the failure probability implies that for any constant security parameter  $\lambda$ , a stash of constant size is sufficient to ensure that the failure probability is smaller than  $2^{-\lambda}$  for sufficiently large  $n$ . In order to achieve a failure probability that is negligible in  $n$ , we can set the stash size  $s$  to be slightly larger than  $O(1)$ , e.g.,  $s = \log \log n$ ,  $s = \log^* n$ , or any  $s = \omega(1)$ . The result is a construction with an overhead of  $\omega(n)$ . (More accurately, the overhead is as close as desired to being linear: for any  $f(n) \in \omega(n)$ , the overhead is  $o(f(n))$ .)

**Concrete overhead.** In §5 we present a new variant of Cuckoo hashing that we denote as two-dimensional (or 2D) Cuckoo hashing. We analyze this construction in §A and show that when no stash is used, then the failure probability (with tables of size  $O(n)$ ) is  $O(1/n)$ , as in standard Cuckoo hashing.

We only have a sketch of an analysis for the size of the stash of the construction in §5, but we observed that this construction performed much better than the asymptotic construction. Also, performance was improved with the heuristic of using half as many bins but letting each bin store two items instead of one. (This variant is known to perform much better also in the case of standard Cuckoo hashing, see [Wie16].)

Since we do not have a theoretical analysis of this construction, we ran extensive experiments in order to examine its performance. These experiments follow the analysis paradigm given in §3.2, and are described in §6.1. For a specific ratio between the table size and  $n$ , we ran  $2^{40}$  experiments for  $n = 2^{12}$  and found that the failure probability is at most  $2^{-37}$  with 99.9% confidence. We also ran experiments for increasing values of  $n$ , up to  $n = 2^{12}$ , and found that the failure probability has linear dependence on  $n^{-3}$  (an explanation of this behavior appears in §B). Therefore, we can argue that for  $n \geq 2^{13} = 2 \cdot 2^{12}$  the failure probability is at most  $2^{-37} \cdot 2^{-3} = 2^{-40}$ .

## 4 An Asymptotic Construction through Mirror Cuckoo Hashing

We show here a construction for circuit-based PSI that has an  $\omega(n)$  asymptotic overhead. The analysis in this section is not intended to be tight, but rather shows the asymptotic behavior of the overhead.

The analysis is based on a construction which we denote as *mirror Cuckoo hashing* (as the placement of the hash functions that are used in one side is a mirror image of the hash functions of the other side). Hashing is computed in a single iteration. The main advantage of this construction is that it is based on four copies of standard Cuckoo hashing. Therefore, we can apply known bounds on the failure probability of Cuckoo hashing. Namely, applying the result of [KMW09] that the failure probability when using a stash of size  $s$  is  $O(n^{-(s+1)})$ . Given this result, a stash of size  $\omega(1)$  guarantees that the failure probability is negligible in  $n$  (while a constant stash size guarantees that for sufficiently large  $n$  the failure probability is smaller than any constant, and in particular smaller than  $2^{-40}$ ). We note that while the known results about the size of the stash are only stated for  $s = O(1)$ , we show in §C that the  $O(n^{-(s+1)})$  bound on the failure probability also applies to a non-constant stash size.

### 4.1 Mirror Cuckoo Hashing

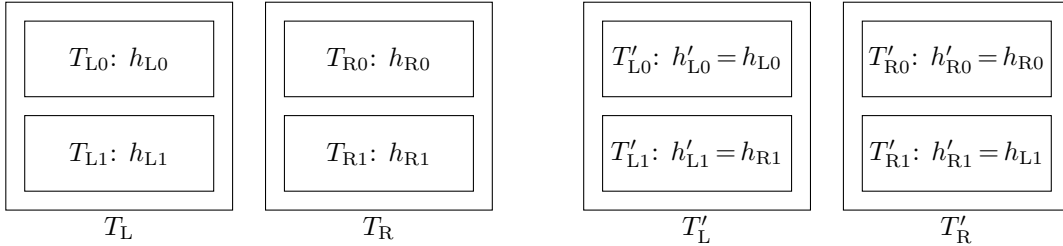
We describe a hashing scheme that uses two sets of tables. A left set including tables  $T_L, T_R$ , and a right set including tables  $T'_L, T'_R$ . Each table is also denoted as a “column”. Each table has two subtables, or “rows”. So overall there are four tables (columns), each containing two subtables (rows).

Bob maps each of his items to one subtable in each table, namely to one row in each column. Alice maps each of her items to the two subtables in one of the tables, namely to both rows in just one of the columns. These mappings ensure that for any item  $x$  that is owned by both Alice and Bob, there is exactly one subtable to which it is mapped by both parties.

*The tables.* The construction uses two sets of tables,  $T_L, T_R$  and  $T'_L, T'_R$ . Each table is of size  $2(1 + \varepsilon)n$  and is composed of two subtables of size  $(1 + \varepsilon)n$  ( $T_L$  contains the subtables  $T_{L0}, T_{L1}$ , etc.). Each subtable is associated with a hash function that will be used by both parties. E.g., function  $h_{L0}$  will be used for subtable  $T_{L0}$ , etc. The tables and the hash functions are depicted in Figure 1.

*The hash functions.* The hash functions associated with the tables are defined as follows:

- The functions for the left two tables (columns)  $T_L, T_R$ , i.e.,  $h_{L0}, h_{L1}, h_{R0}, h_{R1}$ , are chosen at random. Each function maps items to the range  $[0, (1 + \varepsilon)n - 1]$ , which corresponds to the number of bins in each of  $T_{L0}, T_{L1}, T_{R0}, T_{R1}$ .
- The functions for the two right tables  $T'_L, T'_R$  are defined as follows:
  - The two functions of the upper subtables are equal to the functions of the upper subtables on the left. Namely,  $h'_{L0} = h_{L0}$  and  $h'_{R0} = h_{R0}$ .
  - The two functions of the lower subtables are the *mirror image* of the functions of the lower subtables on the left. Namely,  $h'_{L1}, h'_{R1}$  are defined such that  $h'_{L1} = h_{R1}$ , and  $h'_{R1} = h_{L1}$ .



**Fig. 1.** The tables  $T_L, T_R$  and  $T'_L, T'_R$ . The hash functions in the upper subtables of  $T'_L, T'_R$  are the same as in  $T_L, T_R$ , and those in the lower subtables are in reverse order.

*Bob's insertion algorithm.* Bob needs to insert each of his items to one subtable in each of the tables  $T_L, T_R, T'_L, T'_R$ . He can do so by simply using Cuckoo hashing for each of these tables. For example, for the table  $T_L$  and its subtables  $T_{L0}, T_{L1}$ , Bob uses the functions  $h_{L0}, h_{L1}$  to insert each input  $x$  to either  $T_{L0}$  or  $T_{L1}$ . The same is applied to  $T_R, T'_L$ , and  $T'_R$ . In addition, Bob keeps a small stash of size  $\omega(1)$  for each of the four tables. Overall, based on known properties of Cuckoo hashing, we can claim that the construction guarantees the following property:

*Claim.* With all but negligible probability, it holds that for every input  $x$  of Bob, and for each of the four tables  $T_L, T_R, T'_L, T'_R$ , Bob inserts  $x$  to exactly one of the two subtables or to the stash.

*Alice's insertion algorithm.* Alice's operation is a little more complex and is described in Algorithm 1. Alice considers the two upper subtables on the left,  $T_{L0}, T_{R0}$ , as two subtables for standard Cuckoo hashing. Similarly, she considers the two lower subtables on the left,  $T_{L1}, T_{R1}$ , as two subtables for standard Cuckoo hashing. In other words, she considers the left top row and the left bottom row as standard Cuckoo hashing tables.

Alice then inserts each input item of hers to each of these two tables using standard Cuckoo hashing. (She also uses stashes of size  $\omega(1)$  to store items which cannot be placed in the Cuckoo tables.) For some input items  $x$  it happens that  $x$  is inserted in the top row to  $T_{L0}$  and in the

**Algorithm 1 (Mirror Cuckoo hashing)**

1. Alice uses Cuckoo hashing to insert each item  $x$  to one of the subtables  $T_{L0}, T_{R0}$ , using the hash functions  $h_{L0}, h_{R0}$ .
2. Similarly, Alice uses Cuckoo hashing to insert each item  $x$  to one of the subtables  $T_{L1}, T_{R1}$ , using the hash functions  $h_{L1}, h_{R1}$ .
3. At this point, Alice observes the result of the first two steps. For some inputs  $x$  it happened that they were mapped to the same “column” in both of these steps. Namely,  $x$  was mapped to both  $T_{L0}$  and  $T_{L1}$ , or to both  $T_{R0}$  and  $T_{R1}$ . These are the “good” items, since they were mapped to the same column, as is required for all of Alice’s inputs.
4. The other inputs of Alice, the “bad” items, were mapped to one column in Step 1 and to the other column in Step 2. Alice applies the following procedure to these items:
  - (a) Each “bad” item  $x$  is removed from both locations to which it was mapped in Steps 1 and 2.
  - (b)  $x$  is now inserted in either of  $T'_{L0}, T'_{R0}$  using the hash functions  $h'_{L0} := h_{L0}, h'_{R0} := h_{R0}$  with the same mapping as in Step 1.
  - (c)  $x$  is also inserted in either of  $T'_{L1}, T'_{R1}$  using the hash functions  $h'_{L1} := h_{R1}, h'_{R1} := h_{L1}$  with the same mapping as in Step 2.

bottom row to  $T_{L1}$ ; or  $x$  is inserted in the top row to  $T_{R0}$  and in the bottom row to  $T_{R1}$ . Therefore,  $x$  is inserted in two subtables in the same column. ( $x$  is denoted as “good” since this is the outcome that we want.)

Let  $x'$  be one of the other, “bad”, items. Thus,  $x'$  is inserted in the top row to  $T_{L0}$  and in the bottom row to  $T_{R1}$ , or vice versa. In this case, Alice removes  $x'$  from the tables on the left and inserts it to the tables  $T'_L, T'_R$  on the right. Since the hash functions that are used in  $T'_L, T'_R$  are equal to the functions used on the left side (where in the bottom row the functions are in reverse order), Alice does not need to run a Cuckoo hash insertion algorithm on the right side: Assume that  $x'$  was stored in locations  $T_{L0}[h_{L0}(x')]$  and  $T_{R1}[h_{R1}(x')]$  on the left. Then Alice inserts it to locations  $T'_{L0}[h'_{L0}(x')] = T'_{L0}[h_{L0}(x')]$  and  $T'_{L1}[h'_{L1}(x')] = T'_{L1}[h_{R1}(x')]$  on the right.

In other words, in a global view, one can see the algorithm as composed of the following steps: (1) First, all items are placed in the left tables. (2) Each subtable is divided in two copies, where one copy contains the good items and the other copy contains the bad items. (3) The subtable copies with the good items are kept on the left, whereas the copies with the bad items are moved to the right, where in the bottom row on the right we replace the order of the subtables.

This algorithm has two important properties: First, all items that were successfully inserted in the first step to the left tables will be placed in tables on either the left or the right hand sides. Moreover, each item will be placed in two subtables in the same column — the good items happened to initially be placed in this way in the left tables; whereas the bad items were in different columns on the left side but were moved to the same column on the right side. Hence, we can state the following claim:

*Claim.* With all but negligible probability, Alice inserts each of her inputs either to two locations in exactly one of  $T_L, T_R, T'_L, T'_R$  and to no locations in other tables, or to a stash.

*Tables size.* The total size of the tables is  $8(1 + \varepsilon)n$ .

*Stash size.* With regards to stashes, each party needs to keep a stash for each of the Cuckoo hashing tables that it uses. Since Alice runs the Cuckoo hashing insertion algorithm only for the left tables and re-uses the mapping for the right tables, she needs only two stashes. Bob on the other hand runs the Cuckoo hashing insertion algorithm four times and hence needs four stashes. (In order to preserve

simplicity, we omitted the stashes in Figure 1 and Algorithm 1.) Given the result of [KMW09], and our observation in §C about its applicability to non-constant stash sizes, it holds that a total stash of size  $\omega(1)$  elements suffices to successfully map all items, except with negligible probability. We note that the size of the stash can be arbitrarily close to constant, e.g., it can be set to be  $O(\log \log n)$  or  $O(\log^* n)$ . Essentially, for any function  $f(n) \in \omega(n)$ , the size of the stash can be  $o(f(n))$ .

## 4.2 Circuit-based PSI from Mirror Cuckoo Hashing

Mirror Cuckoo hashing lets the parties map their inputs to tables of size  $O(n)$  and stashes of size  $\omega(1)$ , with negligible failure probability. It is therefore straightforward to construct a PSI protocol based on this hashing scheme:

1. The parties agree on the parameters that define the size of the tables and the stash for mirror Cuckoo hashing. They also agree on the hash functions that will be used in each table.
2. Each party maps its items to the tables using the hash functions that were agreed upon.
3. The parties evaluate a circuit that performs the following operations:
  - (a) For each bin in the tables, the circuit compares the item that Alice mapped to the bin to the item that Bob mapped to the same bin.
  - (b) Each item that Bob mapped to his stashes is compared with all items of Alice. Similarly, each item that Alice mapped to her stashes is compared with all items of Bob.

The properties of mirror Cuckoo hashing ensure: (1) If an item  $x$  is in the intersection, then there is exactly one comparison in which  $x$  is input by both Alice and Bob. (2) The number of comparisons in Step 3 is  $\omega(n)$ .

## 5 A Concretely Efficient Construction through 2D Cuckoo Hashing

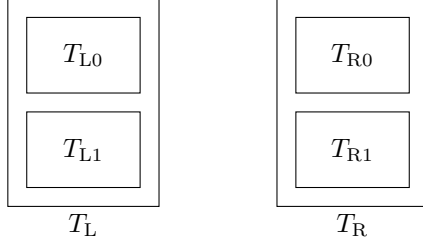
Two-dimensional Cuckoo hashing (a.k.a. 2D Cuckoo hashing) is a new construction with the following properties:

- It uses overall  $O(n)$  memory (specifically,  $8(1 + \varepsilon)n$  in our construction, where we set  $\varepsilon = 0.2$  in our experiments).
- Both, Alice and Bob, map each of their items to  $O(1)$  memory locations (specifically, to two or four memory locations in our construction).
- If  $x$  appears in the input of both parties, then there is exactly one location to which both Alice and Bob map  $x$ .

The construction uses two tables,  $T_L, T_R$ , located on the left and the right side, respectively. Each of these tables is of size  $4(1 + \varepsilon)n$  and is composed of two smaller subtables:  $T_L$  is composed of the two smaller subtables  $T_{L0}, T_{L1}$ , while  $T_R$  is composed of the two smaller tables  $T_{R0}, T_{R1}$ . The hash functions  $h_{L0}, h_{L1}, h_{R0}, h_{R1}$  are used to map items to  $T_{L0}, T_{L1}, T_{R0}, T_{R1}$ , respectively. The tables are depicted in Figure 2.

Hashing is performed in the following way:

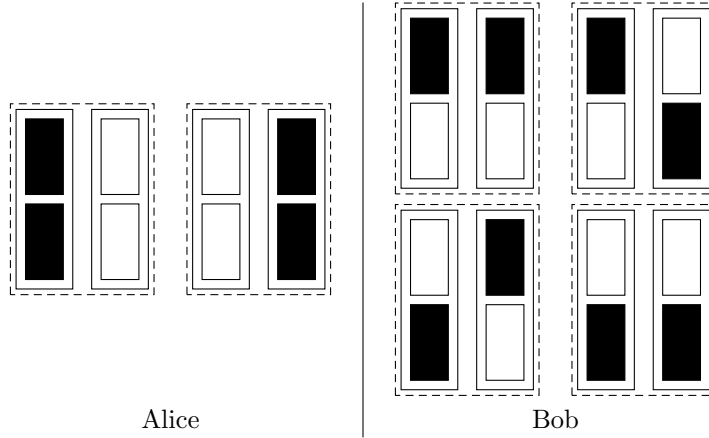
- Alice maps each of her items to all subtables on one of the two sides. Namely, each item  $x$  of Alice is either mapped to both bins  $T_{L0}[h_{L0}(x)]$  and  $T_{L1}[h_{L1}(x)]$  on the left side, or to bins  $T_{R0}[h_{R0}(x)]$  and  $T_{R1}[h_{R1}(x)]$  on the right side. In other words, ALICE maps each item to ALL subtables on one side.



**Fig. 2.** The tables  $T_L$  and  $T_R$ , consisting of  $T_{L0}, T_{L1}$  and  $T_{R0}, T_{R1}$ , respectively.

- Bob maps each of his items to one subtable on each side. This is done using standard Cuckoo hashing. Namely, each input  $x$  of Bob is mapped to one of the locations  $T_{L0}[h_{L0}(x)]$  or  $T_{L1}[h_{L1}(x)]$  on the left side, as well as mapped to one of the locations  $T_{R0}[h_{R0}(x)]$  or  $T_{R1}[h_{R1}(x)]$  on the right side. In other words, BOB maps each item to one subtable on BOTH sides.

The possible options for hashing an item  $x$  by both parties are depicted in Figure 3. It is straightforward to see that if both parties have the same item  $x$ , there is exactly one table out of  $T_{L0}, T_{L1}, T_{R0}, T_{R1}$  that is used by both Alice and Bob to store  $x$ .



**Fig. 3.** The possible combinations of locations to which Alice and Bob map their inputs.

We next describe a construction of 2D Cuckoo hashing, followed by a variant based on a heuristic optimization that stores two items in each table entry. The asymptotic behavior of the basic construction is analyzed in §A. In §6.1 we describe simulations for setting the parameters of the heuristic construction in order to reduce the hashing failure probability to below  $2^{-40}$ .

### 5.1 Iterative 2D Cuckoo Hashing

This construction uses two tables,  $T_L, T_R$ , each of  $4(1 + \varepsilon)n$  entries. (In this construction, there is no need to assume that each table is composed of two subtables. We note that even in the mirror Cuckoo hashing construction of §4.1 it was possible to consider each of  $T_L, T_R$  as a single table that is not divided into subtables, but we think that a construction using two subtables per table

**Algorithm 2 (Iterative 2D Cuckoo hashing)**

1. Alice maps all of her items to table  $T_L$ , using simple hashing. That is, each item  $x$  is inserted in locations  $h_{L0}(x), h_{L1}(x)$ . Obviously, there will be entries in  $T_L$  that will have more than a single item mapped to them.  
Denote  $T_L$  as the active table.
2. For each entry in the active table with more than one item in it: remove all items – except for the item that was mapped to this entry most recently – and move them to the “relocation pool”. For each of the removed items, remove the item also from its other appearance in the active table. (At the end of this step, all entries in the active table have at most one entry. However, there might be items in the relocation pool.)
3. If the relocation pool is empty, then stop (found a successful mapping).
4. Change the designation of the active table to point to the other table.
5. Move each item  $x$  from the relocation pool to locations  $h_{R0}(x), h_{R1}(x)$  in the active table. (For example, if  $T_R$  is the active table, move  $x$  to  $h_{R0}(x), h_{R1}(x)$ .)
6. Go to Step 2.

is conceptually simpler to understand.) The parties associate two hash functions with each table, namely  $h_{L0}, h_{L1}$  for  $T_L$ , and  $h_{R0}, h_{R1}$  for  $T_R$ .

As in the mirror Cuckoo hashing construction described in §4.1, Bob uses Cuckoo hashing to insert each of his items into one location in each of the tables.

Alice inserts each item  $x$  either into the two locations  $h_{L0}(x)$  and  $h_{L1}(x)$  in  $T_L$ , or into the two locations  $h_{R0}(x)$  and  $h_{R1}(x)$  in  $T_R$ . This is achieved by Alice running a modified Cuckoo insertion algorithm that maps an item to two locations in one table, “kicks out” any item that is currently present in these locations and also removes the other occurrence of this item from the table, and then tries to insert this item into its two locations in the other table, and so on.

This is a new variant of Cuckoo hashing, where inserting an item into a table might result in four elements that need to be stored in the other table: storing  $x$  in  $h_{L0}(x), h_{L1}(x)$  might remove two items,  $y_0, y_1$ , one from each location. These items are also removed from their other occurrences in  $T_L$ . They must now be stored in locations  $h_{R0}(y_0), h_{R1}(y_0), h_{R0}(y_1), h_{R1}(y_1)$  in  $T_R$ .

It is not initially clear whether such a mapping is possible (with high probability, given random choices of the hash functions). We analyze the construction in §A and show that it only fails with probability  $O(1/n)$ . We ran extensive simulations, showing that the algorithm (when using a stash and a certain choice of parameters) fails with very small probability, smaller than  $2^{-40}$ .

The insertion algorithm of Alice is described in Algorithm 2. The choice made in Step 2 of the algorithm, to first remove the oldest items that were mapped to the entry, is motivated by the intuition that it is more likely that the locations to which these items are mapped in the other table are free.

**Storing two items per bin.** It is known that the space utilization of Cuckoo hashing can be improved by storing more than one item per bin (cf. [Pan05, DW07] or the review of multiple choice hashing in [Wie16]). We take a similar approach and use two tables of size  $2(1 + \epsilon)n$  where each entry can store *two* items. (These tables have half as many entries as before, but each entry can store two items rather than one. The total size of the tables is therefore unchanged.) The change to the insertion algorithm is minimal and affects only Step 2. The new algorithm is defined in Algorithm 3.

Our experiments in §6.1 show that when using the same amount of space, then this variant of iterative 2D Cuckoo hashing performs better than the basic protocol with bins of size one. That is,

**Algorithm 3 (Iterative 2D Cuckoo hashing with bins of size 2)**

The algorithm is identical to Algorithm 2, except for the following change in Step 2:

2. For each entry in the active table with more than *two* items in it: remove all items – except for the *two* items that were mapped to this entry most recently – and move them to the “relocation pool”. For each of the removed items, remove the item also from its other appearance in the active table.

it achieves a lower probability of hashing failure, namely of the need to use the stash, and requires less iterations to finish.

## 5.2 Circuit-based PSI from 2D Cuckoo Hashing

This section describes how 2D Cuckoo hashing can be used for computing PSI. In addition, we describe two optimizations which substantially improve the efficiency of the protocol. The first optimization has the parties use permutation-based hashing [ANS10] (as was done in [PSSZ15]) in order to reduce the size of the items that are stored in each bin, and hence reduce the number of gates in the circuit. The second optimization is based on having each party use a single stash instead of using a separate stash for each Cuckoo hashing instance.

The PSI protocol is pretty straightforward given 2D Cuckoo hashing:

First, the parties agree on the hash functions to be used in each table. (These functions must be chosen at random, independently of the inputs, in order not to disclose any information about the inputs. Therefore, a participant cannot change the hash functions if some items cannot be mapped, and thus we seek parameter values that make the hashing failure probability negligible, e.g., smaller than  $2^{-40}$ .)

Then, each party maps its items to bins using 2D Cuckoo hashing and the chosen hash functions. The important property is that if Alice and Bob have the same input item then there exists exactly one bin into which both parties map this item (or, alternatively, at least one of them places this item in a stash). Empty bins are padded with dummy elements. This ensures that no information is leaked by how empty the tables and stashes are.

Afterwards, the parties construct a circuit that compares, for each bin, the items that both parties stored in it. In addition, this circuit compares each item that Alice mapped to the stash with all of Bob’s items, and vice versa. Since the number of bins is  $O(n)$ , the number of items in each bin is  $O(1)$ , and the number of items in the stash is  $\omega(1)$ , the total size of this circuit is  $\omega(n)$ . The parties can define another circuit that takes the output of this circuit and computes a desired function of it, e.g., the number of items in the intersection.

Finally, the parties run a generic MPC protocol that securely evaluates this circuit (cf. §6.3 for a concrete implementation and benchmarks).

**Permutation-based Hashing.** The protocol uses permutation-based hashing to reduce the bitlength of the elements that are stored in the bins and thus reduces the size of the circuit comparing them. This idea was introduced in [ANS10] and used for PSI in [PSSZ15]. It is implemented in the following way. The hash function  $h$  that is used to map an item  $x$  to one of the  $\beta$  bins is constructed as follows: Let  $x = x_L|x_R$  where  $|x_L| = \log \beta$ . We first assume that  $\beta$  is a power of 2 and then describe the general case. Let  $f$  be a random function with range  $[0, \beta - 1]$ . Then  $h$  maps an



element  $x$  to bin  $x_L \oplus f(x_R)$  and the value stored in the bin is  $x_R$ . The important property is that the stored value has a reduced bitlength of only  $|x| - \log \beta$ , yet there are no collisions (since if  $x, y$  are mapped to the same bin and store the same value, then  $x_R = y_R$  and  $x_L \oplus f(x_R) = y_L \oplus f(y_R)$  and therefore  $x = y$ ).

In the general case, where  $\beta$  is not a power of two, the output of  $h$  is reduced modulo  $\beta$  and a stored extra bit indicates if the output was reduced or not.

For Cuckoo hashing the protocol uses two hash functions to map the elements to the bins in one table. To avoid collisions among the two hash functions, a stored extra bit indicates which hash function was used.

**Using a Combined Stash.** Recall that Alice uses 2D Cuckoo hashing, for which we show experimentally in §6.1 that no stash is needed. Bob, on the other hand, uses two invocations of standard Cuckoo hashing, and therefore when he does not succeed in mapping an item to a table, he must store it in a stash and compare it with all items of Alice. In this case, the parties cannot encode their items using permutation-based hashing, and therefore these comparisons must be of the full-length original values and not of the shorter values computed using permutation-based hashing as described before. Therefore, the size of the circuits that handle the stash values have a considerable effect on the total overhead of the protocol.

We observe that, instead of keeping several stashes, Bob can collect all the values that he did not manage to map to any of the tables in a *combined* stash. Suppose that he maps items to  $c$  tables and that we have an upper bound  $s$  which holds w.h.p. on the size of each stash. A naive approach would use  $c$  stashes of that size, resulting in a total stash size of  $c \cdot s$ . A better approach would be to use a single stash for all these items, since it is very unlikely that all stashes will be of maximal size, and therefore we can show that with the same probability, the size  $s'$  of the combined stash is much smaller than  $c \cdot s$ . To do so, we determine the upper bounds for the combined stash for  $c = 2$ : The probability of having a combined stash of size  $s'$  is  $\sum_{i=0}^{s'} P(i) \cdot P(s' - i)$ , where  $P(i)$  denotes the probability of having a single stash of size  $i$ . The value of  $P(i)$  is  $O(n^{-i}) - O(n^{-(i+1)}) \approx O(n^{-i})$  [KMW09]. We can estimate the exact values of these probabilities based on the experiments conducted by [PSSZ15]: they performed  $2^{30}$  Cuckoo hashing experiments for each  $n \in \{2^{11}, 2^{12}, 2^{13}, 2^{14}\}$  and counted the required stash sizes. Using linear regression, we extrapolated the results for larger sets of  $2^{16}$  and  $2^{20}$  elements. Table 1 shows the required stash sizes when binding the probability to be below  $2^{-40}$ : it turns out that for  $2^{12}$  and  $2^{16}$  elements the combined stash should include only one more element compared to the upper bound for a single stash, whereas for  $2^{20}$  even the same stash size is sufficient. All in all, when comparing to the naive solution with two separate stashes, the combined stash size is reduced by almost a factor of 2x.

**Table 1.** Stash sizes required for binding the error probability to be below  $2^{-40}$  when inserting  $n \in \{2^{12}, 2^{16}, 2^{20}\}$  elements into  $2.4n$  bins using Cuckoo hashing.

Number of elements $n$	$2^{12}$	$2^{16}$	$2^{20}$
Single stash size $s$ (from [PSSZ15, Table 4])	6	4	3
Stash size for two separate stashes $s' = 2s$	12	8	6
Combined stash size $s'$	7	5	3

### 5.3 Extension to a Larger Number of Parties

Computing PSI between the inputs of more than two parties has received relatively little interest. (The challenge is to compute the intersection of the inputs of all parties, without disclosing information about the intersection of the inputs of any subset of the parties.) Specific protocols for this task were given, e.g., in [FNP04, HV17, KMP<sup>+</sup>17]. We note that our 2D Cuckoo hashing can be generalized to  $m$  dimensions in order to obtain a circuit-based protocol for computing the intersection of the inputs of  $m$  parties. The caveat is that the number of tables grows to  $2^m$  and therefore the solution is only relevant for a small number of parties.

We describe the case of three parties: The hashing will be to a set of eight tables  $T_{x,y,z}$ , where  $x, y, z \in \{0, 1\}$ . Any input item of  $P_1$  is mapped to either all tables  $T_{0,0,0}, T_{0,0,1}, T_{0,1,0}, T_{0,1,1}$ , or to all tables  $T_{1,0,0}, T_{1,0,1}, T_{1,1,0}, T_{1,1,1}$ . Namely, the index  $x$  is set to either 0 or 1, and the input item is mapped to all tables with that value of  $x$ . Every input of  $P_2$  is mapped either to all tables whose  $y$  index is 0, or to all tables where  $y = 1$ . Every input of  $P_3$  is mapped either to all tables whose  $z$  index is 0, or to all tables where  $z = 1$ .

It is easy to see that regardless of the choices of the values of  $x, y, z$ , the sets of tables to which all parties map an item intersect in exactly one table. Therefore, the parties can evaluate a simple circuit that checks every bin for equality of the values that were mapped to it by the three parties. It is guaranteed that if the same value is in the input sets of all parties, then there is exactly one bin to which this value is mapped by all three parties. If some items are mapped to a stash by one of the parties, they must be compared with all items of the other parties, but the overhead of this comparison is  $\omega(n)$  if the stash is of size  $\omega(1)$ .

The remaining issue is the required size of the tables. In §A.4 we show that inserting an item into one of two (big) tables, such that the item is mapped to  $k$  locations in that table, requires tables of size greater than  $k^2(1 + \varepsilon)n$ . When computing PSI between three parties using the method described above, we have eight (small) tables, where each party must insert its items to four tables in one plane or to four tables in the other plane. Each such set of four small tables corresponds to a big table in the analysis of §A.4 and is therefore of size  $16(1 + \varepsilon)n$ . The total size of the tables is therefore  $32(1 + \varepsilon)n$ .

### 5.4 No Extension to Security against Malicious Adversaries

We currently do not see how to extend our hashing-based protocols to achieve security against malicious adversaries. As pointed out by [RR17b], it is inherently hard to extend protocols based on Cuckoo hashing to obtain security against malicious adversaries. The reason is that the placement of items depends on the exact composition of the input set, and therefore a malicious party might learn the placement used by the other party.

Coming up with a similar argument as in [RR17b], assume that in our construction in Figure 3, Bob maps an item  $x$  to the two upper subtables and Alice maps  $x$  to the two left subtables. Now assume Alice maliciously deviates from the protocol and places  $x$  only in the upper left subtable, but not in the lower left one. This deviation may allow Alice to learn whether Bob placed  $x$  in the upper or lower subtables: For example, in a PSI-CA protocol Alice could use only dummy elements and  $x$  as an input set and if the cardinality turns out to be 1, then she knows that Bob placed  $x$  in the upper left subtable. However, the locations in which Bob places an item cannot be simulated in the ideal world as they depend on other items in his input set. Therefore, we see no trivial way to provide security against malicious adversaries based on 2D Cuckoo hashing.

## 6 Evaluation

This section describes extensive experiments that set the parameters for the hashing schemes, the resulting circuit sizes, and the results of experiments evaluating PSI using these circuits.

### 6.1 Simulations for Setting the Parameters of 2D Cuckoo Hashing

We experimented with the iterative 2D Cuckoo hashing scheme described in §5.1, set concrete sizes for the tables, and examined the failure probabilities of hashing to the tables.

Our implementation is written in C and available online at <http://encrypto.de/code/2DCuckooHashing>. It repeatedly inserts a set of random elements into two tables using random hash functions. The insertion algorithm is very simple: All elements are first inserted into the two locations to which they are mapped (by the hash functions) in the first table. Obviously, many table entries will contain multiple items. Afterwards, the implementation iteratively moves items between the tables, in order to reduce the maximum bin occupancy below a certain threshold (cf. Algorithm 2 and Algorithm 3 in §5.1).

**Run-time.** We report in §6.3 the results of experiments analyzing the run-time of the 2D Cuckoo hashing insertion algorithm. Overall, the insertion time (a few milliseconds) is negligible compared to the run-time of the entire PSI protocol.

**Hashing to bins of size 1.** First, we checked if it is possible to use a maximum bin occupation of 1. For this, we set the sizes of each of the two tables to be  $4.8n$  (corresponding to the threshold size of  $4(1 + \epsilon)n$  in the analysis in §A, as well as twice the recommended size for Cuckoo hashing, since all elements are inserted twice). We ran the experiment 100 000 times with input size  $n = 2^{12}$  and bitlength 32. For all except 828 executions it was possible to reduce the maximum bin occupation to 1 after at least 7 and at most 129 iterations of the insertion algorithm. On average, 20 iterations of the insertion algorithm were necessary to achieve the desired result. In said 828 cases there remained at least one bin with more than one item even after 500 iterations of the insertion algorithm. This implies that iterative 2D Cuckoo hashing works in principle, but, as standard Cuckoo hashing, requires a stash for storing the elements of overfull bins.

**Hashing to bins of size 2.** For PSI protocols it would be desirable to avoid having an additional stash on Alice’s side. In standard Cuckoo hashing it is possible to achieve better memory utilization and less usage of the stash by using fewer bins, where each bin can store two items [Wie16]. Therefore, we changed the parameters as follows: the table size is halved and reduced to  $2.4n$ , but each bin is allowed to contain two elements. This way, while consuming the same amount of memory as before, we try to achieve better utilization. We followed the paradigm that was described in §3.2 for the experimental analysis of the failure probability. Namely, we ran massive sets of experiments to measure the number of failures for several values of  $n$  and several table sizes, and given this data we (1) found confidence intervals for the failure probability for specific values of the parameters, and (2) found how the failure probability behaves as a function of  $n$ .

Our first experiment ran  $2^{40}$  tests within  $\sim 2$  million core hours on the Lichtenberg<sup>6</sup> high performance computer of the TU Darmstadt for input size  $n = 2^{12}$ . We chose input size  $2^{12}$  (instead

<sup>6</sup> See <http://www.hh1r.tu-darmstadt.de/hh1r/index.en.jsp> for details on the hardware configuration.

of larger sizes like  $2^{16}$  or  $2^{20}$ ) since running experiments with larger values of  $n$  would have taken even more time and would have simply been impractical. It turned out that the insertion algorithm was successful in reducing the maximum bin size to 2 (after at most 18 iterations) in all but one test.

Given this data, we calculated the confidence interval of the failure probability  $p$ . The probability of observing one failure in  $N$  experiments is  $N \cdot p \cdot (1 - p)^{N-1}$ , where in our experiments  $N = 2^{40}$ . We checked the values of  $p$  for which the probability of this observation is greater than 0.001 and concluded that with 99.9% confidence, the failure probability for iterative 2D Cuckoo hashing with set size  $n = 2^{12}$  and table size  $2.4n$  lies within  $[2^{-50}, 2^{-37}]$ . (Namely, there is at most a 0.001 probability that we would have seen one failure in  $2^{40}$  runs if  $p$  was greater than  $2^{-37}$  or smaller than  $2^{-50}$ .)

**Measuring the dependence on the parameters.** To get a better understanding on how the failure probability behaves for different input and table sizes, we performed a set of experiments that required another  $\sim 3.5$  million core hours. Concretely, we ran  $2^{40}$  tests for each set size  $n \in \{2^6, 2^8, 2^{10}\}$  and each table size in the range  $2.2n$ ,  $2.4n$ , and  $2.6n$ . We also tested the table size  $3.6n$  for  $n \in \{2^6, 2^8\}$  as well as table sizes  $3.0n$  and  $3.2n$  for  $n = 2^{10}$ . The results for all experiments are given in Table 2 and are depicted in Figure 4.

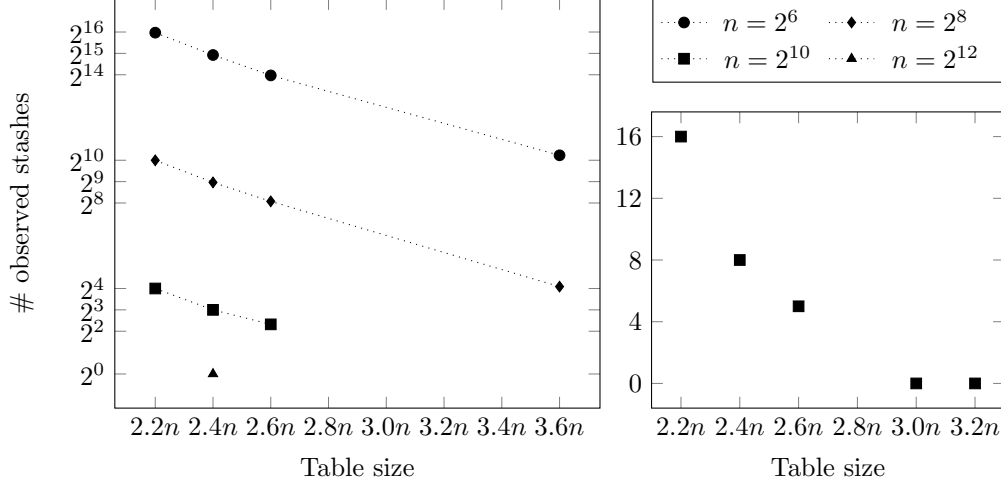
The results demonstrate that, w.r.t. the dependence on  $n$ , for set sizes  $n \in \{2^6, 2^8, 2^{10}\}$  it can be observed that increasing the set size by factor 4x reduces the failure probability by factor 64x. (For larger set sizes, the number of failures is too small to be meaningful.) These experiments also demonstrate that the dependence of the failure probability on  $n$  is  $O(n^{-3})$ . An intuitive theoretical explanation why the probability behaves this way is given in §B. As for the dependence on the table size, the failure probability decreases by a factor of 2x when increasing the table size in steps of  $0.2n$  within the tested range  $2.2n$  to  $3.6n$ .

From these results (a failure probability of at most  $2^{-37}$  for  $n = 2^{12}$  with table size  $2.4n$  and a dependence of  $O(n^{-3})$  of the failure probability on  $n$ ) we conclude that the failure probability for  $n \geq 2^{13}$  and table size  $2.4n$  is at most  $2^{-40}$ .

In total we spent about 5.5 million core hours on our experiments on the Lichtenberg high performance computer of the TU Darmstadt.

**Table 2.** Number of observed stashes for different table sizes and set sizes  $n$  when performing  $2^{40}$  tests of iterative 2D Cuckoo hashing.

Table size	Stash size	$n = 2^6$	$n = 2^8$	$n = 2^{10}$	$n = 2^{12}$
$2.2n$	1	64 020	1 021	16	—
	2	154	1	0	—
	3	4	0	0	—
$2.4n$	1	31 033	499	8	1
	2	65	0	0	0
$2.6n$	1	16 014	270	5	—
	2	33	0	0	—
$3.0n$	1	—	—	0	—
$3.2n$	1	—	—	0	—
$3.6n$	1	1 202	17	—	—



**Fig. 4.** Number of observed stashes for different table and set sizes when performing  $2^{40}$  tests of iterative 2D Cuckoo hashing.

## 6.2 Circuit Complexities

We compare the complexities of the different circuit-based PSI constructions for two sets, each with  $n$  elements that have bitlength  $\sigma$ . We consider two possible bitlengths:

1. **Fixed bitlength:** Here, the elements have fixed bitlength  $\sigma = 32$  bits (e.g., for IPv4 addresses).
2. **Arbitrary bitlength:** Here, the elements have arbitrary bitlength and are hashed to values of length  $\sigma = 40 + 2 \log_2(n) - 1$  bits, with a collision probability that is bounded by  $2^{-40}$ . (See Appendix A of the full version of [PSZ14] for an analysis.) Therefore, we set the bitlength to  $\sigma = 40 + 2 \log_2(n) - 1$  bits.

For all protocols we report the circuit size where we count only the number of AND gates, since many secure computation protocols provide free computation of XOR gates. We compute the size of the circuits up to the step where single-bit wires indicate if a match was found for the respective element. We note that for many circuits computing functions of the intersection, this part of the circuit consumes the bulk of the total size. For example, computing the Hamming weight of these bits is equal to computing the cardinality of the intersection (PSI-CA). The size-optimal Hamming weight circuit of [BP06] has size  $x - w_H(x)$  and depth  $\log_2 x$ , where  $x$  is the number of inputs and  $w_H(\cdot)$  is the Hamming weight. The size of the Hamming weight circuit is negligible compared to the rest of the circuit. As another example, if the cardinality is compared with a threshold (yielding a PSI-CAT protocol), this only adds  $3 \log_2 n$  AND gates and depth  $\log_2 \log_2 n$  using the depth-optimized construction described in [SZ13], which is also negligible.

*The size of the Sort-Compare-Shuffle circuit.* The Sort-Compare-Shuffle circuit [HEK12] has three phases. In the SORT phase, the two sorted lists of inputs are merged into one sorted list, which takes  $2\sigma n \log_2(2n)$  AND gates. In the COMPARE phase, neighboring elements are compared to find the elements in the intersection, which takes  $\sigma(3n - 1) - n$  AND gates. The SHUFFLE phase randomly permutes these values and takes  $\sigma(n \log_2(n) - n + 1)$  AND gates. To have a fair comparison with our protocols, we remove the SHUFFLE phase and let the COMPARE phase output only a

single bit that indicates if a match was found for the respective element or not; this removes  $n$  multiplexers of  $\sigma$ -bit values from the COMPARE phase, i.e.,  $\sigma n$  AND gates. Hence, the total size is  $2\sigma n \log_2(n) + 2\sigma n - n - \sigma + 2$  AND gates.

*The size of the Circuit-Phasing circuit.* The Circuit-Phasing circuit [PSSZ15] has  $2.4nm(\sigma - \log_2(2.4n) + 1) + sn(\sigma - 1)$  AND gates where  $m$  is the maximum occupancy of a bin for simple hashing and  $s$  is the size of the stash.

*The size of our mirror Cuckoo hashing construction of §4.2.* In this protocol, Bob uses a total of four Cuckoo hash tables, so he needs four stashes. Alice uses Cuckoo hashing for the first set of tables, and then she essentially applies the same Cuckoo hashing in reverse order for the second set of tables, so she needs only two stashes. The total number of stashes is hence 6. The protocol compares 4 times the shortened representations (using permutation-based hashing) for Cuckoo tables with  $2.4n$  bins and 6 stashes of  $s$  elements each with  $n$  elements using the full representation of  $\sigma$  bits. Hence, the total complexity is  $4 \cdot 2.4n(\sigma - \log_2(2.4n) + 1) + 3s'n(\sigma - 1)$ , where  $s' \leq 2s$  is the size of the combined stash. (The concept of the combined stash, cf. §5.2, can also be applied to the mirror Cuckoo hashing construction).

*The size of our iterative 2D Cuckoo hashing construction of §5.2.* Each of the following operations is performed twice for the left and right side: (1) For each of the  $2.4n$  bins the shortened representation (cf. §5.2) of the single item in Bob’s bin is compared with the two elements in the corresponding bin of Alice. (2) Bob has a stash of size  $s'$ . Each item in the stash is compared to all of Alice’s items (using the full bitlength representation). Hence, the overall complexity is  $4 \cdot 2.4n(\sigma - \log_2(2.4n) + 1) + s'n(\sigma - 1)$  AND gates, where  $s'$  is the size of the combined stash. By comparing this formula with the one for the mirror construction above we directly see that the iterative construction is always better as it requires only one combined stash.

**Concrete Circuit Sizes.** The Sort-Compare-Shuffle construction [HEK12] has a circuit of size  $O(\sigma n \log n)$ . The Circuit-Phasing construction [PSSZ15] has circuit size  $O(\sigma n \log n / \log \log n)$ , while the asymptotic construction we present in this paper has a size of  $\omega(\sigma n)$  and the iterative 2D Cuckoo hashing construction has an even smaller size.

For a comparison of the concrete circuit sizes, we use the parameters from the analysis in [PSSZ15]: For  $n = 2^{12}$  elements the maximum bin size for simple hashing is  $m = 18$ , for  $n = 2^{16}$  we set  $m = 19$ , and for  $n = 2^{20}$  we set  $m = 20$ . We set the stash size  $s$  and the combined stash size  $s'$  according to Table 1 (on page 17).

On the left side of Table 3 we compare the concrete circuit sizes for *fixed* bitlength  $\sigma = 32$  bit. Our best protocol (“Ours Iterative Combined”) improves over the best previous protocol by factor 2.0x for  $n = 2^{12}$  (over [HEK12]), by factor 2.7x for  $n = 2^{16}$  (over [PSSZ15]), and by factor 3.2x for  $n = 2^{20}$  (over [PSSZ15]).

On the right side of Table 3 we compare the concrete circuit sizes for *arbitrary* bitlength  $\sigma$ . Our best protocol (Ours Iterative Combined) improves over the best previous protocol by factor 1.8x for  $n = 2^{12}$  (over [HEK12]), by factor 2.8x for  $n = 2^{16}$  (over [HEK12]), and by factor 3.8x for  $n = 2^{20}$  (over [PSSZ15]).

Two obvious conclusions from the tables are that the iterative construction always results in a smaller circuit than the mirror construction, and that a combined stash always results in a smaller circuit than using separate stashes. Our constructions always have smaller circuits than both former

constructions, and, due to our better asymptotic size, the savings become greater as  $n$  increases. (The Sort-Compare-Shuffle construction of [HEK12] performs better than the Phasing construction of [PSSZ15] for smaller values of  $n$ , but for larger values it is the opposite. This is expected, as the asymptotic sizes are  $O(\sigma n \log n)$  and  $O(\sigma n \log n / \log \log n)$ , respectively.)

**Table 3.** Concrete circuit sizes in #ANDs for PSI variants on  $n$  elements of fixed bitlength  $\sigma = 32$  (left) and arbitrary bitlength hashed to  $\sigma = 40 + 2 \log_2(n) - 1$  bits (right).

Protocol	Fixed Bitlength $\sigma = 32$			Arbitrary Bitlength		
	$n = 2^{12}$	$n = 2^{16}$	$n = 2^{20}$	$n = 2^{12}$	$n = 2^{16}$	$n = 2^{20}$
Sort-Compare-Shuffle [HEK12]	3 403 746	71 237 602	1 408 237 538	6 705 091	158 138 299	3 478 126 515
Circuit-Phasing [PSSZ15]	4 254 256	55 155 466	688 258 388	10 501 475	181 928 305	3 201 695 060
Separate stashes $s' = 2s$						
Ours Mirror Separate	5 347 225	58 659 626	703 253 572	11 137 330	144 538 001	2 063 466 359
Ours Iterative Separate	2 299 801	26 153 770	313 183 300	5 042 482	71 137 681	1 081 999 223
Combined stash $s'$ (cf. Table 1)						
Ours Mirror Combined	3 442 585	40 375 082	410 700 868	7 328 050	103 250 321	1 327 366 007
Ours Iterative Combined	1 664 921	20 058 922	215 665 732	3 772 722	57 375 121	836 632 439

**Circuit Depths.** For some protocols, the circuit depth is a relevant metric (e.g., for the GMW protocol the depth determines the round complexity of the online phase). Our constructions have the same depth as the Circuit-Phasing protocol of [PSSZ15], i.e.,  $\log_2 \sigma$ . This is much more efficient than the depth of the Sort-Compare-Shuffle circuit of [HEK12] which is  $O(\log \sigma \cdot \log n)$  when using depth-optimized comparison circuits.

**Further Optimizations.** So far, we computed the comparisons with a Boolean circuit consisting of 2-input gates: For elements of bitlength  $\ell$ , the circuit XORs the elements and afterwards computes a tree of  $\ell - 1$  non-XOR gates s.t. the final output is 1 if the elements are equal or 0 otherwise. This circuit allows to use an arbitrary secure computation protocol based on Boolean gates, e.g., Yao or GMW. However, the construction of evaluating 2-input gates in GMW using 1-out-of-4 OTs described in [Gol04, Sect. 7.3.3] can easily be generalized to evaluating  $d$ -input gates with 1-out-of- $2^d$  OTs which can be instantiated very efficiently with the 1-out-of- $N$  OT extension protocol of [KK13].<sup>7</sup> Further optimizations and an implementation are described in [DKS<sup>+</sup>17] who call these  $d$ -input gates lookup tables (LUTs). Their best LUT has 7 inputs and requires only 372 bits of total communication (cf. [DKS<sup>+</sup>17, Tab. IV]). For computing equality, 6 of the non-XOR gates in the tree can be combined into one 7-input LUT. This improves communication of the Circuit-Phasing protocol of [PSSZ15] and our protocols by factor  $6 \cdot 256/372 = 4.1x$ .

### 6.3 Performance

We empirically compare the performance of our iterative 2D Cuckoo hashing PSI-CAT protocol with a combined stash described in §5.2 with the Circuit-Phasing PSI-CAT protocol of [PSSZ15]. As

<sup>7</sup> This observation was made concurrently and independently of each other in [DKS<sup>+</sup>17, KKW17].

a baseline, we also compare with the public key-based PSI-CA protocol of [Sha80, Mea86, DGT12] that leaks the cardinality to one party, and the currently best specialized PSI protocol of [KKRT16] that cannot be easily modified to compute variants of the set intersection functionality. We do not compare with our mirror hashing scheme from §4.1, as the analysis in §6.2 already showed that this protocol is less efficient than the iterative one.

**Implementation.** Pinkas et al. [PSSZ15] provide the implementation of their Circuit-Phasing PSI protocol as part of the ABY framework [DSZ15]. This framework allows to securely evaluate the PSI circuit using either Yao’s garbled circuit or the GMW protocol, both implemented with most recent optimizations (cf. §2). However, since the evaluation in [PSSZ15] showed that using the GMW protocol yields much better run-times, we focus only on GMW. We also compare with the LUT-based evaluation protocols that are described, optimized, and implemented in [DKS<sup>+</sup>17] (cf. §6.2) and integrated in ABY. For the Circuit-Phasing PSI-CAT protocol, we extended the existing codebase with the Hamming weight circuit of [BP06] and the depth-optimized comparison circuit of [SZ13] to compare the Hamming weight with a threshold. Based on this, we implemented our iterative 2D Cuckoo hashing PSI-CAT protocol by duplicating the code for simple hashing and Cuckoo hashing, combining the stashes, and implementing the iterative insertion algorithm. Our implementation is available online as part of the ABY framework at <http://crypto.de/code/ABY>. For the DH/ECC-based protocol of Shamir/Meadows/De Cristofaro et al. [Sha80, Mea86, DGT12], we use the ECC-based implementation of [PSSZ15] available online at <http://crypto.de/code/PSI> that already supports computing the cardinality (PSI-CA). The implementation of the special purpose BaRK-OPRF PSI protocol of [KKRT16] is taken from <https://github.com/osu-crypto/BaRK-OPRF>.

*Benchmarking Environment.* For our benchmarks we use two machines, each equipped with an Intel Core i7-4790 CPU @ 3.6 GHz and 16 GB of RAM. The CPUs support the AES-NI instruction set for fast AES evaluations. We distinguish two network settings: a LAN setting and a WAN setting. For the LAN setting, we restrict the bandwidth of the network interfaces to 1 Gbit/s and enforce a round-trip time of 1 ms. For the WAN setting, we limit the bandwidth to 100 Mbit/s and set a round-trip time of 100 ms. We instantiate all protocols corresponding to a computational security parameter of 128 bit and a statistical security parameter of 40 bit. All reported run-times are the average of 10 executions with less than 10% variance.

**Benchmarking Results.** In Table 4, we give the run-times for  $n \in \{2^{12}, 2^{16}, 2^{20}\}$  elements<sup>8</sup> of bitlength  $\sigma = 32$  (suitable, e.g., for IPv4 addresses). The corresponding communication is given in Table 6. We do not use the LUT-based evaluation in the LAN setting since there is little need for better communication while the run-times are not competitive. However, to demonstrate the advantages of the LUT-based evaluation in the WAN setting, we compare the protocols when running with a single thread and four threads.<sup>9</sup>

*Run-times (Table 4 and Table 5).* In comparison with the Circuit-Phasing PSI-CAT protocol of [PSSZ15] in Table 4, our iterative combined PSI-CAT protocol is faster by factor 1.4x for  $n = 2^{12}$  and up to factor 2.8x for  $n = 2^{20}$ . This holds when the circuit is evaluated with GMW in both

<sup>8</sup> Unfortunately, the LUT-based implementation of [DKS<sup>+</sup>17] was not capable of evaluating the PSI circuits for  $n = 2^{20}$  elements.

<sup>9</sup> We do not provide benchmarks with multiple threads for the DH/ECC PSI-CA protocol since the implementation of [PSSZ15] does not support multi-threading.



**Table 4.** Total run-times in ms for PSI variants on  $n$  elements of bitlength  $\sigma = 32$  bit.

Network setting		LAN			WAN				
Circuit evaluation protocol		GMW [GMW87]			GMW [GMW87]			LUT [Gol04, DKS <sup>+</sup> 17]	
Protocol	Set size $n$	$2^{12}$	$2^{16}$	$2^{20}$	$2^{12}$	$2^{16}$	$2^{20}$	$2^{12}$	$2^{16}$
DH/ECC PSI-CA [Sha80, Mea86, DGT12]		3 296	49 010	7 904 054	4 082	51 866	8 008 771	4 082	51 866
BaRK-OPRF PSI [KKRT16]		113	295	3 882	540	1 247	14 604	540	1 247
<i>1 Thread</i>									
Circuit-Phasing PSI-CAT [PSSZ15]		3 170	20 401	242 235	15 143	99 433	1 042 712	19 951	117 438
Ours Iterative Separate PSI-CAT		2 433	11 251	122 008	11 210	57 474	547 950	15 656	70 545
Ours Iterative Combined PSI-CAT		2 220	9 076	86 648	10 060	45 252	389 891	12 999	56 179
<i>4 Threads</i>									
Circuit-Phasing PSI-CAT [PSSZ15]		2 333	10 600	123 765	12 492	97 480	987 459	15 471	76 184
Ours Iterative Separate PSI-CAT		1 903	6 273	64 324	9 361	56 141	541 677	11 946	46 797
Ours Iterative Combined PSI-CAT		1 694	5 177	49 417	8 793	44 596	376 591	9 413	39 272

network settings and for both 1 and 4 threads. With the implementation of the LUT-based protocols of [DKS<sup>+</sup>17], we observe a further improvement for the circuit-based protocols by about 13% in the WAN setting, but only for medium set sizes of  $n = 2^{16}$  and 4 threads due to the higher computation complexity.

The circuit-based protocols have two steps: mapping the input items to the tables, and securely evaluating the circuit. The run-times of the hashing step are shown in Table 5. The times for Cuckoo hashing into two tables in our PSI-CAT protocol are exactly twice of those for Cuckoo hashing into one table in [PSSZ15]. Compared to simple hashing, our 2D Cuckoo hashing is slower by factor 1.6x up to factor 2.1x due to the additional iterations. However, all in all, the hashing procedures are by 2-3 orders of magnitude faster than the times for securely evaluating the circuit, and therefore negligible w.r.t. the overall run-time.

In comparison with the DH-based PSI-CA protocol of [Sha80, Mea86, DGT12], our iterative combined PSI-CAT protocol is faster by factor 1.5x for  $n = 2^{12}$  up to factor 91x for  $n = 2^{20}$  in the LAN setting with a single thread. Also in the WAN setting with a single thread, our protocol is faster (except for small sets with  $n = 2^{12}$ ), despite the substantially lower communication of the DH-based protocol described below. In both network settings even the best measured run-times of our PSI-CAT protocol are between 19x to 36x slower than the BaRK-OPRF specialized PSI protocol of [KKRT16], but our protocols are generic.

**Table 5.** Run-times in ms for hashing  $n$  elements of bitlength  $\sigma = 32$  bit.

Hashing Procedure	Set size $n$	$2^{12}$	$2^{16}$	$2^{20}$
<i>Circuit-Phasing PSI-CAT [PSSZ15]</i>				
Simple Hashing		3.50	27.96	557.54
Cuckoo Hashing		2.43	15.87	391.16
<i>Ours Iterative PSI-CAT</i>				
2D Cuckoo Hashing		6.23	58.90	873.19
Cuckoo Hashing (for two tables with a combined stash)		4.85	31.75	782.32

*Communication (Table 6).* The communication given in Table 6 is measured on the network interface, so these numbers are slightly larger than the theoretical communication (derived from the number of AND gates on the left side in Table 3) due to TCP/IP headers and padding of messages. The lowest communication is achieved by the DH-based PSI-CA protocol of [Sha80, Mea86, DGT12] which is in line with the experiments in [PSSZ15]. Our best protocol for PSI-CAT has between 132x (for  $n = 2^{12}$ ) and 66x (for  $n = 2^{20}$ ) more communication than the DH-based PSI-CA protocol when evaluated with GMW. Recall, however, that our protocol does not leak the cardinality. Our best protocol improves the communication over the PSI-CAT protocol of [PSSZ15] by factor 2.3x (for  $n = 2^{12}$ ) to 2.9x (for  $n = 2^{20}$ ). When using LUT-based evaluation with optimizations of [DKS<sup>+</sup>17], we observe that the communication of all circuit-based PSI-CAT protocols improves over GMW by factor 3.7x which is close to the theoretical upper bound of 4.1x (cf. §6.2). Still, our best LUT-based protocol has more than 20x higher communication than the BaRK-OPRF specialized PSI protocol of [KKRT16], but it is generic.

**Table 6.** Communication in MB for PSI variants on  $n$  elements of bitlength  $\sigma = 32$  bit.

Protocol	Set size $n$	$2^{12}$	$2^{16}$	$2^{20}$
DH/ECC PSI-CA [Sha80, Mea86, DGT12]		0.4	6.6	106.0
BaRK-OPRF PSI [KKRT16]		0.53	8.06	127.20
<i>GMW</i> [GMW87]				
Circuit-Phasing PSI-CAT [PSSZ15]		121.9	1 588.9	20 028.5
Ours Iterative Separate PSI-CAT		72.3	826.1	9 971.4
Ours Iterative Combined PSI-CAT		52.7	638.8	6 950.6
<i>LUT</i> [Gol04, DKS <sup>+</sup> 17]				
Circuit-Phasing PSI-CAT [PSSZ15]		32.6	418.1	—
Ours Iterative Separate PSI-CAT		19.4	221.3	—
Ours Iterative Combined PSI-CAT		14.3	171.3	—

**Application to privacy-preserving ridesharing.** Our PSI-CAT protocol can easily be extended for the privacy-preserving ridesharing functionality of [HOS17], where the intersection is revealed only if the size of the intersection is larger than a threshold. The authors of [HOS17] give a protocol that securely computes this functionality, but has quadratic computation complexity. By slightly extending our circuit for PSI-CAT to encapsulate a key that is released only if the size of the intersection is larger than the threshold and using this key to symmetrically encrypt the last message in any linear complexity PSI protocol (e.g., [PSZ14, PSSZ15, KKRT16, PSZ18]), we get a protocol with almost linear complexity. Our key encapsulation would take less than 3 seconds for  $n = 2^{12}$  elements (cf. our results for PSI-CAT in Table 4), whereas the solution of [HOS17] takes 5 627 seconds, i.e., we improve by factor 1 876x and also asymptotically.

**Acknowledgments.** We thank Oleksandr Tkachenko for his invaluable help with the implementation and benchmarking. We also thank Moni Naor for suggesting the application to achieve differential privacy. Finally, we thank Karn Seth for clarifying the description of the protocol from [IKN<sup>+</sup>17]. This work has been co-funded by the DFG as part of project E4 within the CRC 1119 CROSSING and by the German Federal Ministry of Education and Research (BMBF), the Hessen State Ministry for Higher Education, Research and the Arts (HMWK) within CRISP, and

the BIU Center for Research in Applied Cryptography and Cyber Security in conjunction with the Israel National Cyber Bureau in the Prime Minister’s Office. Calculations for this research were conducted on the Lichtenberg high performance computer of the TU Darmstadt.

## References

- ADN<sup>+</sup>13. N. Asokan, A. Dmitrienko, M. Nagy, E. Reshetova, A.-R. Sadeghi, T. Schneider, and S. Stelle. CrowdShare: Secure mobile resource sharing. In *Applied Cryptography and Network Security (ACNS’13)*, volume 7954 of *LNCS*, pages 432–440. Springer, 2013.
- ALSZ13. G. Asharov, Y. Lindell, T. Schneider, and M. Zohner. More efficient oblivious transfer and extensions for faster secure computation. In *Computer and Communications Security (CCS’13)*, pages 535–548. ACM, 2013.
- ANS10. Y. Arbitman, M. Naor, and G. Segev. Backyard Cuckoo hashing: Constant worst-case operations with a succinct representation. In *Foundations of Computer Science (FOCS’10)*, pages 787–796. IEEE, 2010.
- AP11. R. R. Amossen and R. Pagh. A new data layout for set intersection on GPUs. In *International Symposium on Parallel and Distributed Processing (IPDPS’11)*, pages 698–708. IEEE, 2011.
- BHKR13. M. Bellare, V. Hoang, S. Keelveedhi, and P. Rogaway. Efficient garbling from a fixed-key blockcipher. In *Symposium on Security and Privacy (S&P’13)*, pages 478–492. IEEE, 2013.
- BP06. J. Boyar and R. Peralta. Concrete multiplicative complexity of symmetric functions. In *Mathematical Foundations of Computer Science (MFCS’06)*, volume 4162 of *LNCS*, pages 179–189. Springer, 2006.
- CLR17. H. Chen, K. Laine, and P. Rindal. Fast private set intersection from homomorphic encryption. In *Computer and Communications Security (CCS’17)*, pages 1243–1255. ACM, 2017.
- DC17. A. Davidson and C. Cid. An efficient toolkit for computing private set operations. In *Australian Conference on Information Security and Privacy (ACISP’17)*, volume 10343 of *LNCS*, pages 261–278. Springer, 2017.
- DCW13. C. Dong, L. Chen, and Z. Wen. When private set intersection meets big data: An efficient and scalable protocol. In *Computer and Communications Security (CCS’13)*, pages 789–800. ACM, 2013.
- DD15. S. K. Debnath and R. Dutta. Secure and efficient private set intersection cardinality using Bloom filter. In *Information Security Conference (ISC’15)*, volume 9290 of *LNCS*, pages 209–226. Springer, 2015.
- DGT12. E. De Cristofaro, P. Gasti, and G. Tsudik. Fast and private computation of cardinality of set intersection and union. In *Cryptology and Network Security (CANS’12)*, volume 7712 of *LNCS*, pages 218–231. Springer, 2012.
- DKS<sup>+</sup>17. G. Dessouky, F. Koushanfar, A.-R. Sadeghi, T. Schneider, S. Zeitouni, and M. Zohner. Pushing the communication barrier in secure computation using lookup tables. In *Network and Distributed System Security (NDSS’17)*. The Internet Society, 2017.
- DKT10. E. De Cristofaro, J. Kim, and G. Tsudik. Linear-complexity private set intersection protocols secure in malicious model. In *Advances in Cryptology – ASIACRYPT’10*, volume 6477 of *LNCS*, pages 213–231. Springer, 2010.
- DSMRY09. D. Dachman-Soled, T. Malkin, M. Raykova, and M. Yung. Efficient robust private set intersection. In *Applied Cryptography and Network Security (ACNS’09)*, volume 5536 of *LNCS*, pages 125–142. Springer, 2009.
- DSZ15. D. Demmler, T. Schneider, and M. Zohner. ABY – a framework for efficient mixed-protocol secure two-party computation. In *Network and Distributed System Security (NDSS’15)*. The Internet Society, 2015.
- DT10. E. De Cristofaro and G. Tsudik. Practical private set intersection protocols with linear complexity. In *Financial Cryptography (FC’10)*, volume 6052 of *LNCS*, pages 143–159. Springer, 2010.
- DW07. M. Dietzfelbinger and C. Weidling. Balanced allocation and dictionaries with tightly packed constant size bins. *Theoretical Computer Science*, 380(1-2):47–68, 2007.
- Dwo06. C. Dwork. Differential privacy. In *International Colloquium on Automata, Languages and Programming (ICALP’06)*, volume 4052 of *LNCS*, pages 1–12. Springer, 2006.
- EFG<sup>+</sup>15. R. Egert, M. Fischlin, D. Gens, S. Jacob, M. Senker, and J. Tillmanns. Privately computing set-union and set-intersection cardinality via Bloom filters. In *Australasian Conference on Information Security and Privacy (ACISP’15)*, volume 9144 of *LNCS*, pages 413–430. Springer, 2015.
- EFL12. Y. Eijgenberg, M. Farbstain, M. Levy, and Y. Lindell. SCAPI: The secure computation application programming interface. Cryptology ePrint Archive, Report 2012/629, 2012. <http://ia.cr/2012/629>.

- EGMT17. D. Eppstein, M. Goodrich, M. Mitzenmacher, and M. Torres. 2-3 cuckoo filters for faster triangle listing and set intersection. In *Symposium on Principles of Database Systems (PODS'17)*, pages 247–260. ACM, 2017.
- FHNP16. M. J. Freedman, C. Hazay, K. Nissim, and B. Pinkas. Efficient set intersection with simulation-based security. *Journal of Cryptology*, 29(1):115–155, 2016.
- FNP04. M. J. Freedman, K. Nissim, and B. Pinkas. Efficient private matching and set intersection. In *Advances in Cryptology – EUROCRYPT'04*, volume 3027 of *LNCS*, pages 1–19. Springer, 2004.
- GMW87. O. Goldreich, S. Micali, and A. Wigderson. How to play any mental game or a completeness theorem for protocols with honest majority. In *Symposium on Theory of Computing (STOC'87)*, pages 218–229. ACM, 1987.
- Gol04. O. Goldreich. *Foundations of Cryptography*, volume 2: Basic Applications. Cambridge University Press, 2004.
- Gon81. G. H. Gonnet. Expected length of the longest probe sequence in hash code searching. *Journal of the ACM*, 28(2):289–304, 1981.
- HCE11. Y. Huang, P. Chapman, and D. Evans. Privacy-preserving applications on smartphones. In *Hot topics in Security (HotSec'11)*. USENIX, 2011.
- HEK12. Y. Huang, D. Evans, and J. Katz. Private set intersection: Are garbled circuits better than custom protocols? In *Network and Distributed System Security (NDSS'12)*. The Internet Society, 2012.
- HEKM11. Y. Huang, D. Evans, J. Katz, and L. Malka. Faster secure two-party computation using garbled circuits. In *USENIX Security'11*, pages 539–554. USENIX, 2011.
- HL08. C. Hazay and Y. Lindell. Efficient protocols for set intersection and pattern matching with security against malicious and covert adversaries. In *Theory of Cryptography Conference (TCC'08)*, volume 4948 of *LNCS*, pages 155–175. Springer, 2008.
- HN10. C. Hazay and K. Nissim. Efficient set operations in the presence of malicious adversaries. In *Public Key Cryptography (PKC'10)*, volume 6056 of *LNCS*, pages 312–331. Springer, 2010.
- HOS17. P. Hallgren, C. Orlandi, and A. Sabelfeld. PrivatePool: Privacy-preserving ridesharing. In *Computer Security Foundations Symposium (CSF'17)*, pages 276–291. IEEE, 2017.
- HV17. C. Hazay and M. Venkatasubramanian. Scalable multi-party private set-intersection. In *Public Key Cryptography (PKC'17)*, volume 10174 of *LNCS*, pages 175–203. Springer, 2017.
- IKN<sup>+</sup>17. M. Ion, B. Kreuter, E. Nergiz, S. Patel, S. Saxena, K. Seth, D. Shanahan, and M. Yung. Private intersection-sum protocol with applications to attributing aggregate ad conversions. Cryptology ePrint Archive, Report 2017/738, 2017. <http://ia.cr/2017/738>.
- KK13. V. Kolesnikov and R. Kumaresan. Improved OT extension for transferring short secrets. In *Advances in Cryptology – CRYPTO'13 (2)*, volume 8043 of *LNCS*, pages 54–70. Springer, 2013.
- KKRT16. V. Kolesnikov, R. Kumaresan, M. Rosulek, and N. Trieu. Efficient batched oblivious PRF with applications to private set intersection. In *Computer and Communications Security (CCS'16)*, pages 818–829. ACM, 2016.
- KKW17. W. S. Kennedy, V. Kolesnikov, and G. T. Wilfong. Overlaying conditional circuit clauses for secure computation. In *Advances in Cryptology – ASIACRYPT'17*, volume 10625 of *LNCS*, pages 499–528. Springer, 2017.
- KLS<sup>+</sup>17. Á. Kiss, J. Liu, T. Schneider, N. Asokan, and B. Pinkas. Private set intersection for unequal set sizes with mobile applications. *Proceedings on Privacy Enhancing Technologies (PoPETs)*, 2017(4):97–117, 2017.
- KMP<sup>+</sup>17. V. Kolesnikov, N. Matania, B. Pinkas, M. Rosulek, and N. Trieu. Practical multi-party private set intersection from symmetric-key techniques. In *Computer and Communications Security (CCS'17)*, pages 1257–1272. ACM, 2017.
- KMW09. A. Kirsch, M. Mitzenmacher, and U. Wieder. More robust hashing: Cuckoo hashing with a stash. *SIAM Journal on Computing*, 39(4):1543–1561, 2009.
- Kre17. B. Kreuter. Secure multiparty computation at Google. In *Real World Crypto Conference (RWC'17)*, 2017. <http://www.totalwebcasting.com/view/?id=columbia&date=2017-01-04&seq=1>.
- KS08. V. Kolesnikov and T. Schneider. Improved garbled circuit: Free XOR gates and applications. In *International Colloquium on Automata, Languages and Programming (ICALP'08)*, volume 5126 of *LNCS*, pages 486–498. Springer, 2008.
- Mea86. C. Meadows. A more efficient cryptographic matchmaking protocol for use in the absence of a continuously available third party. In *Symposium on Security and Privacy (S&P'86)*, pages 134–137. IEEE, 1986.
- Pan05. R. Panigrahy. Efficient hashing with lookups in two memory accesses. In *ACM-SIAM Symposium on Discrete Algorithms (SODA'05)*, pages 830–839. Society for Industrial and Applied Mathematics, 2005.

- PL15. M. Pettai and P. Laud. Combining differential privacy and secure multiparty computation. In *Annual Computer Security Applications Conference (ACSAC'15)*, pages 421–430. ACM, 2015.
- PR01. R. Pagh and F. F. Rodler. Cuckoo hashing. In *European Symposium on Algorithms (ESA'01)*, volume 2161 of *LNCS*, pages 121–133. Springer, 2001.
- PR04. R. Pagh and F. F. Rodler. Cuckoo hashing. *Journal of Algorithms*, 51(2):122–144, 2004.
- PSSZ15. B. Pinkas, T. Schneider, G. Segev, and M. Zohner. Phasing: Private set intersection using permutation-based hashing. In *USENIX Security'15*, pages 515–530. USENIX, 2015.
- PSWW18. B. Pinkas, T. Schneider, C. Weinert, and U. Wieder. Efficient circuit-based PSI via cuckoo hashing. In *Advances in Cryptology – EUROCRYPT'18*, volume 10822 of *LNCS*, pages 125–157. Springer, 2018.
- PSZ14. B. Pinkas, T. Schneider, and M. Zohner. Faster private set intersection based on OT extension. In *USENIX Security'14*, pages 797–812. USENIX, 2014. Full version: <http://ia.cr/2014/447>.
- PSZ18. B. Pinkas, T. Schneider, and M. Zohner. Scalable private set intersection based on OT extension. *ACM Transactions on Privacy and Security (TOPS'18)*, 21(2), 2018.
- RR17a. P. Rindal and M. Rosulek. Improved private set intersection against malicious adversaries. In *Advances in Cryptology – EUROCRYPT'17*, volume 10210 of *LNCS*, pages 235–259. Springer, 2017.
- RR17b. P. Rindal and M. Rosulek. Malicious-secure private set intersection via dual execution. In *Computer and Communications Security (CCS'17)*, pages 1229–1242. ACM, 2017.
- Sha80. A. Shamir. On the power of commutativity in cryptography. In *International Colloquium on Automata, Languages and Programming (ICALP'80)*, volume 85 of *LNCS*, pages 582–595. Springer, 1980.
- SZ13. T. Schneider and M. Zohner. GMW vs. Yao? Efficient secure two-party computation with low depth circuits. In *Financial Cryptography (FC'13)*, volume 7859 of *LNCS*, pages 275–292. Springer, 2013.
- Wie16. U. Wieder. Hashing, load balancing and multiple choice. <https://udiwieder.files.wordpress.com/2014/10/hashbook.pdf>, 2016.
- Yao86. A. C. Yao. How to generate and exchange secrets. In *Foundations of Computer Science (FOCS'86)*, pages 162–167. IEEE, 1986.
- Yun15. M. Yung. From mental poker to core business: Why and how to deploy secure computation protocols? In *Computer and Communications Security (CCS'15)*, pages 1–2. ACM, 2015.
- ZRE15. S. Zahur, M. Rosulek, and D. Evans. Two halves make a whole: Reducing data transfer in garbled circuits using half gates. In *Advances in Cryptology – EUROCRYPT'15*, volume 9057 of *LNCS*, pages 220–250. Springer, 2015.

## A Analysis of 2D Cuckoo Hashing

We analyze a setting which corresponds to the basic iterative 2D Cuckoo hashing of §5.1. The tables are of size  $4(1 + \varepsilon)n$ , and each bin can accommodate a single item. We show that in this setting, and when the hash functions are chosen at random, the probability that it is impossible to successfully map all items to bins is at most  $O(1/n)$ . This success probability is similar to the success probability of Cuckoo hashing [AP11].

We note that we currently do not have an analysis of the optimized construction that maps items to bins of capacity 2 (defined at the end of §5.1), and therefore need to rely on simulations for this case. (This issue is not straightforward for standard Cuckoo hashing as well, see [Wie16].) We also do not analyze the size of the stash that is required for bins of capacity 2, and rely on simulations for this purpose.

On the other hand, we present a new proof technique that is interesting by itself, and can be used to analyze the general case of Cuckoo hashing, which stores  $k$  copies of each item.  $k$  is a parameter that can take arbitrary values, where  $k = 1$  corresponds to basic Cuckoo hashing, and  $k = 2$  corresponds to our 2D Cuckoo hashing.

### A.1 The Problem Setting

We model the problem in the following straightforward way: There are  $n$  items to store,  $x_1, \dots, x_n$ , drawn from a universe  $U$ . There are two tables  $T_0, T_1$ , each with  $m$  bins where  $m \geq 4(1 + \varepsilon)n$ .

A bin can contain at most one item. There exist two hash functions  $h_0, h_1$ , where each function takes an item  $x_i$  and returns a **pair** of random indexes in  $[m]$ . In other words, for  $\sigma \in \{0, 1\}$ ,  $h_\sigma : \mathcal{U} \rightarrow [m] \times [m]$ . (Namely, we represent by  $h_\sigma$  the two hash functions that were assigned to table  $T_\sigma$ , in order to simplify the notation.)

A *placement* of the item is a mapping of each  $x_i$  either to  $T_0[h_0(x)]$  or to  $T_1[h_1(x)]$ , such that every bin is assigned at most one item. So if  $x_i$  is assigned to  $T_\sigma$ , it is placed in both bins indexed by  $h_\sigma(x_i)$ . (Note that this setting of two-dimensional Cuckoo hashing is closely related to that of basic Cuckoo hashing. The difference is that in Cuckoo hashing the hash functions return a single index and not a pair, so that each item is placed in one bin, while here each item is placed twice, in two bins, either in  $T_0$  or in  $T_1$ .)

We define a combinatorial structure, the *inference graph*, which represents the constraints on placing items in this setting. We define an item  $x_i$  to be *bad* if it prevents placing another item  $x_j$  in any of the tables. Our main theorem, Theorem 1, shows that if the size of each table is  $> 4n$ , then the probability of a specific item being bad is  $\frac{O(1)}{n^2}$ , and thus the probability there exists a bad item is  $O(1/n)$ .

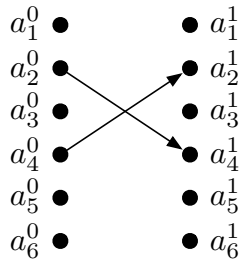
## A.2 The Inference Graph

Given a pair of functions  $h_0, h_1$  and a set of items  $S$ , we build an *Inference Graph*  $G(S, h_0, h_1)$  as follows:

*nodes*: The set of nodes is comprised of two sets  $a_1^0, \dots, a_n^0$  and  $a_1^1, \dots, a_n^1$ . Semantically we think of node  $a_i^\sigma$  as representing the event that  $x_i$  was placed in  $T_\sigma$ .

*edges*: The edges of the graph are directed and represent inferences between the events, so if  $h_\sigma(x_i) \cap h_\sigma(x_j) \neq \emptyset$ , then  $G$  has the directed edges  $(a_i^\sigma, a_j^{1-\sigma})$  and  $(a_j^\sigma, a_i^{1-\sigma})$ . In words, the edge  $(a_i^0, a_j^1)$  means that if  $x_i$  is placed in  $T_0$  then  $x_j$  must be placed in  $T_1$  and vice versa.

Informally, an edge from  $a_i^0$  to  $a_j^1$  represents the fact that placing  $x_i$  in table  $T_0$  makes it impossible to place  $x_j$  in table  $T_0$ , and therefore  $x_j$  must be placed in table  $T_1$ . An example for an inference graph is depicted in Figure 5.



**Fig. 5.** The inference graph for 6 items. Node  $a_i^\sigma$  corresponds to the event that  $x_i$  is mapped to table  $T_\sigma$ . In this example  $h_0(x_2) \cap h_0(x_4) \neq \emptyset$ , and therefore placing  $x_2$  in  $T_0$  induces placing  $x_4$  in  $T_1$ , and vice versa.

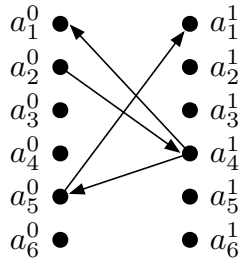
**Placement.** The goal of the following definitions is to form conditions under which an item  $x_i$  could be placed. We then will show that these conditions hold w.h.p. for all items. We denote by

$G(a_i^\sigma)$  the set of nodes reachable from  $a_i^\sigma$  (including  $a_i^\sigma$ ). We may drop the  $\sigma$  in our notation as our claims hold for both  $\sigma = 0$  and  $\sigma = 1$ . When we refer to the items in  $G(a_i)$ , we mean all items associated with nodes in  $G(a_i)$ , that is:  $\{x_j : a_j \in G(a_i)\}$ .

If a node  $a_j^\gamma \in G(a_i^\sigma)$ , then a placement of  $x_i$  in table  $T_\sigma$  implies that  $x_j$  must be placed in table  $T_\gamma$ .

**Definition 1.** A node  $a_i$  in the Inference Graph is called bad if there is a  $j$  such that both  $a_j^0, a_j^1 \in G(a_i)$ . An item  $x_i$  is bad if both  $a_i^0$  and  $a_i^1$  are bad.

Namely, a node  $a_i^\sigma$  is bad if there is an item  $x_j$  such that placing  $x_i$  in table  $T_\sigma$  prevents a placement of  $x_j$  in either table  $T_0$  or table  $T_1$ . (This is depicted in Figure 6.) An item  $x_i$  is bad if placing it in either table  $T_0, T_1$  prevents finding a placement for other items.



**Fig. 6.** Paths coming out of  $a_2^0$  demonstrating that  $a_2^0$  is bad. Namely, placing  $x_2$  in  $T_0$  prevents placing  $x_4$  in  $T_0$ , and therefore  $x_4$  must be placed in  $T_1$ . This prevents placing  $x_1$  in  $T_1$ . In addition, this prevents placing  $x_5$  in  $T_1$ , and therefore it must be placed in  $T_0$ , but this prevents placing  $x_1$  in  $T_0$ . Therefore, if  $x_2$  is placed in  $T_1$ , then  $x_1$  cannot be placed in any of the tables.

Clearly, a bad item implies that not all items could be placed. The following lemma states the converse is also true.

**Lemma 1.** If node  $a_i^\sigma$  is not bad, then all items of  $G(a_i^\sigma)$  could be placed.

*Proof.* We first place  $x_i$  in  $T_\sigma$ . Then we place all its neighbors in  $T_{1-\sigma}$  and continue iteratively. If item  $x_j$  cannot be placed, then it must intersect items both in  $T_0$  and  $T_1$ , which means both  $a_j^0$  and  $a_j^1$  are in  $G(a_i)$ , which is a contradiction.

**Lemma 2.** If none of the items are bad, then all items could be placed.

*Proof.* The algorithm that places all the items is now straightforward: Let  $S$  be the set of currently unplaced items. Pick an item  $x_i \in S$  and, since it is not bad, either  $a_i^\sigma$  is not bad for some  $\sigma \in \{0, 1\}$ . Now, according to Lemma 1, all items in  $G(a_i)$  could be placed successfully. Let  $S'$  be the remaining items, i.e.,  $S' = S \setminus G(x_i)$ . Given  $S'$  and all the free locations in  $T_0, T_1$ , we compute the new inference graph  $G' = G(h_0, h_1, S')$  and continue inductively. The only thing remaining to observe is that if there were no bad items in  $G$ , then there are no bad items in  $G'$ . To see this, observe that  $G'$  is a subgraph of  $G$ . Indeed, let  $x_j$  be an item in  $G'$ . Note that all slots of  $h_\sigma(x_j)$  must be free, otherwise  $a^\sigma \in G(a_i)$ , so every inference made in  $G'$  is true also for  $G$ .

### A.3 Main Result

The goal now is to calculate the probability that an item is bad.

**Theorem 1.** *If the size of each table is greater than  $4(1+\varepsilon)n$ , then for every item  $x_i$  the probability of being bad is at most  $\left(\frac{1+\varepsilon}{\varepsilon}\right)^3 \cdot \frac{2}{m^2}$ , where the probability is taken over the choice of the hash functions.*

Taking a union bound over all  $x_i$  proves the correctness of the scheme. The remainder of the section is dedicated to the proof of the theorem.

Denote by  $m$  the size of  $T_\sigma$ . Our approach is to show that it is unlikely that an item is bad. For that to happen, both its nodes need to be bad, and we would like to count how many bad graphs are there and show that they are unlikely to appear. In order to facilitate this bound, we need to constrain further the exact notion of a bad graph, captured by the next definition.

**Definition 2.** *A bad path rooted at  $a_i^\sigma$  is a simple path from  $a_i^\sigma$  to  $a_i^{1-\sigma}$  in which no item appears twice. That is, for each  $j \neq i$  at most one of  $\{a_j^0, a_j^1\}$  can appear in the path.*

The next lemma shows that bad paths are the only type of bad subgraphs we need to care about.

**Lemma 3.** *If a node is bad, then it is the root of a bad path.*

*Proof.* Assume  $a_i^\sigma$  is bad, then there must be at least one  $j$  for which  $a_i^\sigma$  is connected to both  $a_j^0, a_j^1$ . Further, we can assume that there is no  $k \neq j$  such that both  $a_k^0$  and  $a_k^1$  appear on the paths from  $a_i^\sigma$  to  $a_j^0, a_j^1$ . We can make this assumption because if there is, we may take the pair  $a_k^0, a_k^1$  instead.

Now recall that by construction, if an edge  $(a_k^\sigma, a_\ell^{1-\sigma})$  appears in the inference graph, then so does the edge  $(a_\ell^\sigma, a_k^{1-\sigma})$ . A simple induction shows that if there is a path  $a_i^\sigma \rightsquigarrow a_j^{1-\sigma}$ , then there is a path  $a_j^\sigma \rightsquigarrow a_i^{1-\sigma}$ . Thus, we can construct a path  $a_i^\sigma \rightsquigarrow a_j^\sigma \rightsquigarrow a_i^{1-\sigma}$ . Further, there is no item that is repeated in the path.

We now need to show that the probability of an item being bad is small. We start by bounding the probability of a node being the root of a bad path.

**Lemma 4.** *If  $m \geq 4(1+\varepsilon)n$ , then for every item  $x_i$ , the probability of  $a_i^0$  being the root of a bad path is at most  $\frac{1+\varepsilon}{\varepsilon m}$ .*

*Proof.* We count the labeled bad paths rooted at  $a_i^0$ . Let  $k$  be the number of edges on the path, and  $k+1$  the number of nodes. The head of the path is already set to be  $a_i^0$ , and the tail is  $a_i^1$ , so there are at most  $n^{k-1}$  ways of choosing the nodes in between. We conclude:

$$\#\text{possible labeled bad paths rooted at } a_i^0 \leq \sum_k n^{k-1}. \quad (1)$$

Given a labeling of a bad path, we calculate the probability that it actually appears in the graph. Consider an edge  $a_k^\sigma, a_\ell^{1-\sigma}$ . This edge appears iff  $h_\sigma(x_k)$  has a non empty intersection with  $h_\sigma(x_\ell)$ , which happens with probability  $\leq 4/m$ . Here we use the fact that no item on the path is repeated, so the occurrences of the  $k$  edges are independent events. The probability that a given path appears in the graph is therefore  $\leq 4^k/m^k$ .

Combined with (1) we have that the probability of  $a_i^0$  being the root of a bad path is at most

$$\frac{1}{m} \sum_{k \geq 1} \left(\frac{4n}{m}\right)^{k-1} \leq \frac{1}{m} \sum_{k \geq 1} \left(\frac{1}{1+\varepsilon}\right)^{k-1} \leq \frac{1+\varepsilon}{\varepsilon m}. \quad (2)$$



*Proof (of Theorem 1).* Note that we bounded the probability that node  $a_i^0$  is bad. For item  $x_i$  to be bad, node  $a_i^1$  has to be bad as well. Again, it is enough to bound the case that  $a_i^1$  is the root of a bad path. Now, if the bad path covers different items than the path from  $a_i^0$ , then the same calculation holds and the probability that item  $x_i$  is bad is at most  $\left(\frac{1+\varepsilon}{\varepsilon m}\right)^2$ . Otherwise there is some item  $x_j$  which appears on both bad paths rooted at  $a_i^0$  and  $a_i^1$ . Let  $k_1$  be the length of the bad path starting at  $a_i^0$  and  $k_2$  be the length of the prefix of the bad path starting at  $a_i^1$  and ending at  $a_j^\sigma$ . Note that there are at most  $k_1$  possibilities for choosing  $x_j$ . The probability that such a path exists is therefore at most

$$\begin{aligned} & \frac{1}{m} \sum_{k_1 \geq 1} \left(\frac{1}{1+\varepsilon}\right)^{k_1-1} \frac{2k_1}{m} \sum_{k_2 \geq 1} \left(\frac{1}{1+\varepsilon}\right)^{k_2-1} \\ & \leq \frac{2(1+\varepsilon)}{\varepsilon m^2} \sum_{k_1 \geq 1} k_1 \left(\frac{1}{1+\varepsilon}\right)^{k_1-1} \\ & \leq \left(\frac{1+\varepsilon}{\varepsilon}\right)^3 \cdot \frac{2}{m^2}. \end{aligned}$$

Taking a union bound over all items we have:

**Corollary 1.** *If  $m \geq 4(1+\varepsilon)n$ , then the probability that there is a failure is at most  $\frac{2(1+\varepsilon)^2}{\varepsilon^3 d^2} \cdot \frac{1}{n}$ .*

#### A.4 Generalizations

*Standard Cuckoo hashing.* The proof above is a strict generalization of standard Cuckoo hashing, and in fact could be used to prove tight bounds for Cuckoo hashing. To see this, observe that the only place where the proof uses the assumption that  $h_0(x)$  and  $h_1(x)$  sample *pairs* of indexes is when computing the probability that an edge appears in Lemma 4, which is bounded by  $4/m$ . Indeed, if  $h_0(x)$  and  $h_1(x)$  sample  $d$  locations this probability is  $d^2/m$ . If they are standard hash functions  $\mathcal{U} \rightarrow [m]$  as is in Cuckoo hashing, then  $d = 1$  and the probability of a collision is  $1/m$ . This implies Theorem 1 holds for space  $m > (1+\varepsilon)n$  for each table. This is the tight space requirement as was shown in [PR04], see also [Wie16].

*Placing  $k$  copies of each item.* We can analyze in a similar way what happens when each hash function samples a  $k$ -tuple, so  $k$  copies of each item must be placed either in  $T_0$  or in  $T_1$ . In this case a collision probability of  $k^2/m$  implies the space of each table must be greater than  $k^2n$ . It is not hard to see that 3-party set intersection, as described in §5.3, could be implemented with  $k = 4$ , where  $T_0, T_1$  are each of size  $> 16n$ .

## B Buckets with more than a Single Item

In our experiments in §6.1, we assume that each slot in the table has room for storing two items rather than just one. The experiments showed the following two observations: (1) The total space that is required is smaller. (2) The failure probability decreases and seems to behave like  $O(1/n^3)$ . Both observations are expected to some extent, because they are similar to the behavior of standard Cuckoo hashing. In the following, there is an heuristic explanation for the  $O(1/n^3)$  behavior of the failure probability.

Let's try to construct the smallest example of a scenario where insertion is impossible and the scheme fails. This is a sensible approach, since typically the smallest example is the one most likely to happen (though this is by no means a proof), so the probability that this type of example occurs is a guideline to the real error probability.

Say there are 5 items  $x_1, \dots, x_5$  such that  $\cap_{1 \leq i \leq 5} h_0(x_i) \neq \emptyset$  and  $\cap_{1 \leq i \leq 5} h_1(x_i) \neq \emptyset$ . Clearly, in this case an insertion is not possible, as two items are placed in  $T_0$ , two are placed in  $T_1$ , and there is no placement for the fifth item. Let's calculate the probability that this event happens. For a given set of 5 items, the probability that they intersect in  $T_0$  is roughly  $d^5/m^4$ , where  $d$  is the number of sample locations of the hash functions. To see this, note that the probability that the first two items intersect is  $d^2/m$  and then the probability of each additional intersection is  $d/m$ . So the probability of 5 items intersecting on both sides is  $d^{10}/m^8$ . Now we need to multiply this by  $n^5$  possibilities to pick the 5 items, and, since  $m \approx d^2n$ , we have that this is roughly  $1/(d^6n^3)$ .

## C Cuckoo Hashing with a Stash of Non-Constant Size

Previous research on the size of the stash in standard Cuckoo hashing, e.g., [KMW09], only considered a stash of constant size. Our asymptotic construction in §4 needs a stash of non-constant size in order to make the failure probability negligible.

One approach for proving bounds on the size of the stash is examining the Cuckoo graph and bounding the probability that its connected components would have many cycles. This is done combinatorially by counting all possible such graphs and then computing, for each graph, the probability that it is instantiated in the actual Cuckoo graph. See Theorem 5.5 in [Wie16]. A careful examination of the proof of Theorem 5.5 reveals that in fact there is no assumption that the size of the stash is constant. So the probability that the number of stash edges in the Cuckoo graph exceeds  $s$  is at most

$$\frac{1}{n^s} \cdot \left( \frac{2}{1 + \varepsilon} \right)^s \cdot \sum_t \frac{t^{8s}}{(1 + \varepsilon)^t} \quad (3)$$

(where  $t$  goes over all bad edges in the graph). Setting (say)  $s = \log \log n$ , Equation (3) evaluates to less than  $O(\frac{1}{n^{s-1}})$ , which is sub-polynomial.

It is important to note that the running time of the insertion algorithm is bounded by the size of a connected component in the Cuckoo graph, and that is dominated by a Geometric distribution. So with negligible probability it is not ruled out that connected components would be super-logarithmic in size. Any insertion algorithm needs to take this into account. Still, the insertion time of each element is expected to be  $O(1)$ , and the running time for inserting all items is tightly concentrated around  $O(n)$ .