

Combining Private Set-Intersection with Secure Two-Party Computation

MICHELE CIAMPI
University of Edinburgh
UK
mciampi@ed.ac.uk

CLAUDIO ORLANDI
Aarhus University
Denmark
orlandi@cs.au.dk

Abstract

Private Set-Intersection (PSI) is one of the most popular and practically relevant secure two-party computation (2PC) tasks. Therefore, designing special-purpose PSI protocols (which are more efficient than generic 2PC solutions) is a very active line of research. In particular, a recent line of work has proposed PSI protocols based on oblivious transfer (OT) which, thanks to recent advances in OT-extension techniques, is nowadays a very cheap cryptographic building block.

Unfortunately, these protocols cannot be plugged into larger 2PC applications since in these protocols one party (by design) learns the output of the intersection. Therefore, it is not possible to perform secure *post-processing* of the output of the PSI protocol.

In this paper we propose a novel and efficient OT-based PSI protocol that produces an “encrypted” output that can therefore be later used as an input to other 2PC protocols. In particular, the protocol can be used in combination with all common approaches to 2PC including garbled circuits, secret sharing and homomorphic encryption. Thus, our protocol can be combined with the right 2PC techniques to achieve more efficient protocols for computations of the form $z = f(X \cap Y)$ for arbitrary functions f .

1 Introduction

Private Set-Intersection (PSI) is one of the most practically relevant *secure two-party computation* (2PC) tasks. In PSI two parties hold two sets of strings X and Y , respectively. At the end of the protocol one (or both) party should learn the intersection of the two sets $Z = X \cap Y$ and nothing else about the input of the other party.

There are many real-world applications in which PSI is required. As an example, when mobile users install messaging apps, they need to discover whom among their contacts (from their address book) is also using the app, in order to be able to start communicating seamlessly with them. Doing so requires users to learn the intersection of their contact list with the list of registered users of the service which is stored at the server side. This is typically done by having users send their contact list to the server that can then compute the intersection and return the result to the user. Unfortunately this solution is very problematic not only for the privacy of the user, but for the privacy of the users’ contacts as well! In particular, some of the people in the contact list might not want their phone number being transferred and potentially stored by the server, but they have

no control over this.¹ Note that this is not just a theoretically interesting problem and that Signal (one of the most popular end-to-end encrypted messaging app) has recently recognized this as being a real problem and offered partial solutions to it.² PSI has many other applications, including computing intersections of suspect lists, private matchmaking (comparing interests), *testing human genome* [BBC⁺11], *privacy-preserving ride-sharing* [HOS17], *botnet detection* [NMH⁺10], *advertisement conversion rate* [IKN⁺17] and many more.

From a feasibility point of view, PSI is just a special case of 2PC and therefore any generic 2PC protocol (such as [Yao82, GMW87]) could be used to securely evaluate PSI instances as well. However, since PSI is a natural functionality that can be applied in numerous real-world applications, many efficient protocols for this specific functionality have been proposed, with early results dating back to the 80s [Sha80, Mea86]. The problem was formally defined in [FNP04] and follow up work increased the efficiency of PSI protocols to have complexity only *linear* in the inputs of the parties [JL10, CT10]. A very recent work shows how to obtain a protocol where communication complexity is linear in the size of the smaller set and logarithmic in the larger set [CLR17].

However, these protocols still require performing expensive public-key operations (e.g., exponentiations) for every element in the input sets. As public-key operations are orders of magnitudes more expensive than symmetric key operations, these protocols are not practically efficient for large input sets. In the meanwhile, generic techniques for 2PC had improved by several orders of magnitude and the question of whether special purpose protocols or generic protocols were most efficient has been debated in [HEK12, CT12]. Due to its practical relevance, PSI protocols in the *server-aided* model have been proposed as well [KMRS14]. Independent and concurrent works [PSWW18, FNO18] (which were not publicly available at time we first posted our paper on ePrint) consider the problem of using a PSI protocol to construct more complex functionality in an efficient way. More specifically, [PSWW18] provides a way to securely compute many variants of the set intersection functionality using a clever combination of Cuckoo hashing and garbled circuit. The work of Falk et al. [FNO18] focuses on obtaining a PSI protocol that is efficient in terms of communication. In addition, the authors of [FNO18] propose a PSI protocol where the output can be secret shared that has communication complexity of $O(m\lambda \log \log m)$, where λ is the bit-length of the elements and m is the set-size.

The techniques used in our paper significantly differ from the techniques used in [PSWW18, FNO18]. Our solution avoids the use of garbled circuits and rely on the security and the efficiency of OT and symmetric key encryption schemes.

1.1 OT-based PSI

The most efficient PSI protocols today are those following the approach of PSZ [PSSZ15, PSZ14]. These protocols make extensive use of a cryptographic primitive known as *oblivious transfer (OT)*. While OT provably requires expensive public-key operation, OT can be “extended” as shown by [IKNP03, ALSZ13, KK13] i.e., the few necessary expensive public-key operations can be amortized over a very large number of OT instances, and the marginal cost of OT is only a few (faster)

¹Some apps do not transfer the contact list in cleartext, but a hashed version instead. However, since the domain space of phone numbers is small enough to allow for brute forcing of the hashes, this does not guarantee any real privacy guarantee.

²Unfortunately, the Signal team has concluded that current PSI protocols are too inefficient for their application scenario and relied on trusted-hardware instead, in the style of [TLP⁺17]. See <https://signal.org/blog/private-contact-discovery/> for more details on this.

symmetric key operations instead. In particular, improvements in OT-extension techniques directly imply improvements to PSI protocols as shown by e.g., [KKRT16, OOS17].

In a nutshell, the PSZ protocol introduced two important novel ideas to the state of the art of PSI. First, they give an efficient instantiation of the *private set membership protocol* (PSM) introduced in [FIPR05] based on OT. Second, they show how to efficiently implement PSI from PSM using hashing techniques. (An overview of their techniques is given below).

1.2 Our contribution

The main contribution of this paper is to give an efficient instantiation of PSM that provides output in encrypted format and can therefore be combined with further 2PC protocols. Our PSM protocol can be naturally combined with the hashing approach of PSZ to give a PSI protocol with encrypted output achieving the same favourable complexity in the input sizes. This enables the combination the efficiency of modern PSI techniques with the potentials of general 2PC. Combining our protocols with the right 2PC post-processing allows more efficient evaluation of functionalities of the form $Z = f(X \cap Y)$ for any function f . Like in PSZ we only focus on semi-honest security. Using the protocol together with an actively secure OT-extension protocol such as [ALSZ15, KOS15] would result in a protocol with *privacy* but not *correctness* (i.e., the view of the protocol *without* the output can be efficiently simulated), which is a meaningful notion of security in some settings. PSI protocols with security against malicious adversaries have been proposed in e.g., [HL08, RR17a, RR17b]. It is an interesting open problem to design efficient protocols which are both secure against active (or covert) adversaries and that produce encrypted output. Also, like in PSZ, we only focus on the two-party setting. The recent result of [HV17] has shown that multiparty set-intersection can be computed efficiently. Extending our result to the multiparty case is an interesting future research direction.

We also compare the computation complexity of our scheme for PSM with all the circuit-based PSI approaches (which can be combined with further postprocessing) proposed in [PSZ16]. More precisely, in Table 1 we compare our protocol with the protocols of [PSZ16] in terms of number of symmetric key operations, and bits exchanged between the parties. The result of this comparison is that our protocol has better performance, in terms of computational complexity, than all the circuit-based PSI approaches considered for our comparison³. We refer the reader to App. A for more details about this comparison.

1.3 Improving the efficiency of smart contract protocols

Most of the cryptocurrency systems are built on top of blockchain technologies where miners run distributed consensus whose security is ensured as long as the adversary controls only a minority of the miners. Some cryptocurrency systems allow to run complex programs and decentralized applications on the blockchain. In Ethereum⁴ those programs are called smart contracts. Roughly speaking, the aim of a smart contract is to run a protocol and start a transaction to pay a user of the cryptocurrency systems according to the output of the protocol execution. Unfortunately, this

³The complexity of the protocols proposed in [PSZ16] depends upon parameters that are also related to the used hash function. In order to make our comparison fair, we have set these parameters as showed in the first column in Table 10 of [PSZ16]. More precisely, the authors of [PSZ16] show in that table which parameters are adopted for their empirical efficiency comparison for the case where one set is much greater than the other set (which is exactly the case of PSM).

⁴<http://www.ethereum.org>.

interesting feature of the smart contracts does not come for free. Indeed, in order to execute a smart contract, it is required to pay a *gas fee* that depends on the number of instructions of the protocol to be executed. So, higher is the complexity of protocol, higher is the price to pay. In this context a cryptographic protocol that outputs intermediate values in a secret shared way is particularly useful. Suppose that two parties want to securely compute $f(X \cap Y)$ for arbitrary functions f , and reward another party depending on the output of this computation. Instead of writing on a smart contract the entire protocol to compute $f(X \cap Y)$, the two parties could run a sub-protocol Π to obtain a secret share of $\chi = X \cap Y$ without using a smart contract, and then run another sub-protocol Π' to compute $f(\chi)$, this time using a smart contract to enforce the reward policy. Following this approach it is possible to move part of the computation off-chain, thus increasing the performance and, at the same time, decreasing the costs required to execute the smart contract. Moreover, we observe that χ can be reused to compute different functions f' . The scenario described above is particularly interesting if one of the party can be fully malicious, but in this work we will focus on semi-honest security leaving the above as an open question.

	# of sym. key operations	Communication [bits]
Yao SCS [HEK12]	$12\lambda M \log M + 3\lambda M$	$2\lambda M s(1 + 3 \log M)$
GMW SCS [HEK12]	$12\lambda M \log M$	$6\lambda M(s + 2) \log M$
Yao PWC [PSZ16]	$4\lambda M + 6393\lambda$	$\lambda(M3s + 3198s + 15, 6)$
GMW PWC [PSZ16]	$6\lambda M + 9594\lambda$	$\lambda(M4 + 6396 + 2sM + 6396s)$
This work	$4\lambda M + 3\lambda$	$2\lambda M s + M s$

Table 1: Computation and communication complexity comparison for the PSM case. M represents the size of the set, s is the security parameter and λ is the bit-length of each element.

2 Technical overview

2.1 Why PSZ and 2PC do not mix

We start with a quick overview of the PSM protocol in PSZ [PSSZ15, PSZ14], to explain why their protocol inherently reveals the intersection to one of the parties. From a high-level point of view, the protocol is conceptually similar to the PSM protocol from oblivious pseudorandom function (OPRF) of [FIPR05], except that the OPRF is replaced with a similar functionality efficiently implemented using OT. For simplicity, here we will use the OPRF abstraction.

The goal of a PSM protocol is the following: the receiver R has input x , and the sender S has input a set Y ; at the end of the protocol the receiver learns whether $x \in Y$ or not while the sender learns nothing. The protocol starts by using the OPRF subprotocol, so that R learns $x^* = F_k(x)$ (where k is known to S), whereas S learns nothing about x . Now S evaluates the PRF on her own set and sends the set $Y^* = \{y^* = F_k(y) | y \in Y\}$ to R , who checks if $x^* \in Y^*$ and concludes that $x \in Y$ if this is the case. In other words, we map all inputs into pseudorandom strings and then let one of the parties test for membership “in the clear”. Since the party performing the test doesn’t have access to the mapping (e.g., the PRF key), this party can only check for the membership of x and no other points (i.e., all elements in $Y^* \setminus \{x^*\}$ are indistinguishable from random in R ’s view).

From the above description, it should be clear that the PSZ PSM inherently reveals the output to one of the parties. Turning this into a protocol which provides encrypted output is a challenging task. Here is an attempt at a “strawman” solution: we change the protocol such that R still learns

the pseudorandom string $x^* = F_k(x)$ corresponding to x , but now S sends a value *for every element in the universe*. Namely, for each i (in the domain of Y) S sends an encryption of whether $i \in Y$ “masked” using $F_k(i)$ e.g., S sends $c_i = F_k(i) \oplus E(i \in Y)$ ⁵. Now R can compute $c_x \oplus x^* = E(x \in Y)$ i.e., an encrypted version of whether $x \in Y$, which can be then used as input to the next protocol.

While this protocol produces the correct result, its complexity is exponential in the bit-length of $|x|$, which is clearly not acceptable.

Intuitively, we know that only a polynomial number of c_i ’s will contain encryptions of 1, and therefore we need to find a way to “compress” all the c_i corresponding to $i \notin Y$ into a single one, to bring the complexity of the protocol back to $O(|Y|)$. In the following, after defining some useful notation, we give an intuitive explanation on how to do that.

2.2 Our protocol

We introduce some useful (and informal) notation in order to make easier to understand the ideas behind our construction. We let $Y = \{y_1, \dots, y_M\}$ be the input set of the sender S , and we assume w.l.o.g., that $|Y| = M = 2^m$.⁶ All strings have the same length e.g., $|x| = |y_i| = \lambda$.⁷ We will use a special symbol \perp such that $x \neq \perp \forall x$. We use a function $\text{Prefix}(x, i)$ that outputs the i most significant bits of x ($\text{Prefix}(x, i) \neq \text{Prefix}(x, j)$ when $i \neq j$ independently of the value of x) and for simplicity we define $\text{Prefix}(Y, i)$ to be the set constructed by taking the i most significant bits of every element in Y .

The protocol uses a symmetric key encryption scheme $\text{Sym} = (\text{Gen}, \text{Enc}, \text{Dec})$ with the additional property that given a key $k \leftarrow \text{Gen}(1^s)$ it is possible to efficiently verify if a given ciphertext is in the range of k (see Sec. 3 for a formal definition).

Finally, the output of the protocol will be one of two strings γ_0, γ_1 chosen by S , respectively denoting $x \notin Y$ and $x \in Y$. The exact format of the two strings depends on the protocol used for post-processing. For instance if the post-processing protocol is based on: 1) *garbled circuits*, then γ_0, γ_1 will be the labels corresponding to some input wire; 2) *homomorphic encryption*, then $\gamma_b = \text{Enc}(pk, b)$ for some homomorphic encryption scheme Enc ; 3) *secret-sharing*, then $\gamma_b = s_2 \oplus b$ where s_2 is a uniformly random share chosen by S , so that if R defines its own share as $s_1 = \gamma_b$ then it holds that $s_1 \oplus s_2 = b$.⁸

In order to “compress” the elements of Y we start by considering an undirected graph with a level structure of $\lambda + 1$ levels. The vertices in the last level of this graph will correspond to the elements of Y . More precisely, we associate the secret key $k_{b_\lambda b_{\lambda-1} \dots b_1}$ of a symmetric key encryption scheme Sym to each element $y = b_\lambda b_{\lambda-1} \dots b_1 \in Y$. The main idea is to allow the receiver to obviously navigate this graph in order to get the key $k_{b_\lambda b_{\lambda-1} \dots b_1}$ if $x = y$, for some $y = b_\lambda b_{\lambda-1} \dots b_1 \in Y$, or a special key k^* otherwise. Moreover we allow the receiver to navigate the graph efficiently, that is, every level of the graph is visited only once.

⁵The exact format of the “encryption” $E(\cdot)$ would depend on the subsequent 2PC protocol and is irrelevant for this high-level description.

⁶Sets can always be padded with dummy elements, but the complexity of the protocol can match M that in practice can be $M \approx 2^{m-1}$.

⁷We can assume λ to be smaller than the (statistical) security parameter s and we will denote the bit decomposition of x by $x = x_\lambda \dots x_1$. Otherwise before running the protocol the parties can hash their input down and run the protocol with inputs $h(x)$ and $h(Y) = \{h(y_1), \dots, h(y_M)\}$. Clearly if $x = y_i$ then $h(x) = h(y_i)$, and for correctness we need that $\Pr[h(x) \in h(Y) \wedge x \notin Y] < 2^{-s}$.

⁸Here we use \oplus -secret sharing without loss of generality. Any 2-out-2 secret sharing would work here.

Once a key k is obtained by the receiver, the sender sends $O(|Y|)$ ciphertexts in a such a way that the key obtained by the receiver can decrypt only one ciphertext. Moreover the plaintext of this ciphertext will correspond to γ_0 or γ_1 depending on whether $x \in Y$ or not.

2.2.1 First step: construct the graph G

Each graph level $i \in \{0, \dots, \lambda\}$ has size at most $|\text{Prefix}(Y, i)| + 1$. More precisely, for every $t = b_\lambda b_{\lambda-1} \dots b_{\lambda-i} \in \text{Prefix}(Y, i)$ there is a node in the level i of G that contains a key $k_{b_\lambda b_{\lambda-1} \dots b_{\lambda-i}}$. In addition, in the level i there is a special node, called *sink node* that contains a key k_i^* (which we refer to as *sink key*). The aim of k_i^* is to represent all the values that do not belong to $\text{Prefix}(i, Y)$.

Let us now describe how the graph G is constructed. First, for $i = 1, \dots, \lambda$ the key (for a symmetric key encryption scheme) k_i^* is generated using the generation algorithm $\text{Gen}(\cdot)$. As discussed earlier the aim of these keys is to represent the elements that do not belong to Y . More precisely, the sink key k_i^* , with $i \in \{1, \dots, \lambda\}$ represents all the values that do not belong to $\text{Prefix}(Y, i)$ and the key k_λ^* (the last sink key) will be used to encrypt the output γ_0 corresponding to non-membership in the last step of our protocol. Note that if $\text{Prefix}(x, i) \notin \text{Prefix}(Y, i)$ then $\text{Prefix}(x, j) \notin \text{Prefix}(Y, j)$ for all $j > i$. Therefore, once entered in a sink node, the *sink path* is never abandoned and thus the final sink key k_λ^* , will be retrieved (which allows recovery of γ_0). Let us now give a more formal idea of how G is constructed.

- The root of G is empty, and in the second level there are two vertices k_0 and k_1 where⁹, for $b = 0, 1$

$$k_b = \begin{cases} k \leftarrow \text{Gen}(1^s), & \text{if } b \in \text{Prefix}(Y, 1) \\ k_1^*, & \text{otherwise} \end{cases}$$

- For each vertex k_t in the level $i \in \{1, \dots, \lambda\}$ and for $b = 0, 1$ create the node $k_{t||b}$ as follows (if it does not exists) and connect k_t to it.

$$k_{t||b} = \begin{cases} k \leftarrow \text{Gen}(1^s), & \text{if } t||b \in \text{Prefix}(Y, i+1) \\ k_{i+1}^*, & \text{if } t||b \notin \text{Prefix}(Y, i+1) \\ k_{i+1}^*, & \text{if } k_t = k_i^* \end{cases}$$

We observe that a new node $k_{t||b}$ is generated only when $t||b \in \text{Prefix}(Y, i)$. In the other cases the sink node k_{i+1}^* is used.

In Fig. 1 we show an example of what the graph G looks like for the set $Y = \{010, 011, 110\}$. In this example it is possible to see how, in the 2nd level, all the elements that do not belong to $\text{Prefix}(Y, 2)$ are represented by the sink node k_2^* . Using this technique we have that in the last level of G one node (k_3^* in this example) is sufficient to represent all the elements that do not belong to Y . Therefore, we have that the last level of G contains at most $|Y| + 1$ elements. We also observe that every level of G cannot contain more than $|Y| + 1$ nodes.

⁹In abuse of notation we refer to a vertex using the key represented by the vertex itself.

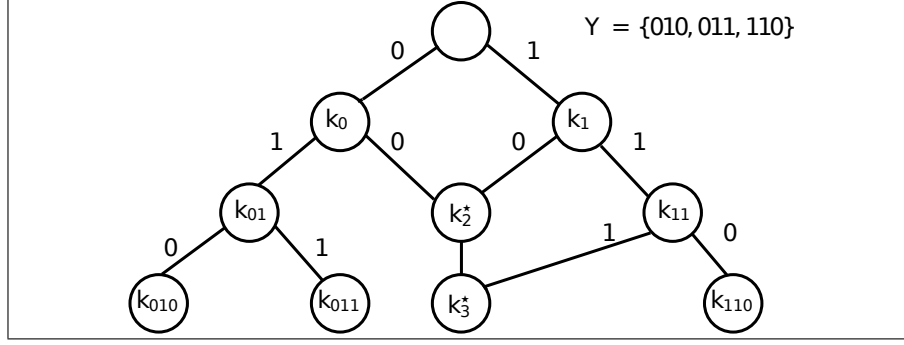


Figure 1: Example of how the graph G appears when the sender holds the set Y .

2.2.2 Second step: oblivious navigation of G

Let $x = x_\lambda x_{\lambda-1} \dots x_1$ be the receiver's (R's) private input and Y be the sender's (S's) private input. After S constructs the graph G we need a way to allow R to obtain $k_{x_\lambda x_{\lambda-1} \dots x_1}$ if $x \in Y$ and the sink key k_λ^* otherwise. All the computation has to be done in such a way that no other information about the set Y is leaked to the receiver, and as well that no information about x is leaked to the sender. In order to do so we use λ executions of 1-out-of-2 OT. The main idea is to allow the receiver to select which branch to explore in G depending on the bits of x . More precisely, in the first execution of OT, R will receive the key k_{x_λ} iff there exists an element in Y with the most significant bit equal to x_λ , the sink key k_1^* otherwise. In the second execution of OT, R uses $x_{\lambda-1}$ as input and S uses (c_0, c_1) where c_0 is computed as follows:

- For each key in the second level of G that has the form $k_{t||0}$, the key $k_{t||0}$ is encrypted using the key k_t .
- For every node v in the first level that is connected to a sink node k_2^* in the second level, compute an encryption of k_2^* using the key contained in v .
- Pad the input with random ciphertexts up to the upper bound for the size of this layer (more details about this step are provided later).
- Randomly permute these ciphertexts.

The procedure to compute the input c_1 is essentially the same (the only difference is that in this case we consider every key with form $k_{t||1}$ and encrypt it using k_t).

Roughly speaking, in this step every key contained in a vertex u of the second level is encrypted using the keys contained in the vertex v of the previous level that is connected to u . For example, following the graph provided in Fig.1, c_0 would be equal to $\{\text{Enc}(k_0, k_2^*), \text{Enc}(k_1, k_2^*)\}$ and c_1 to $\{\text{Enc}(k_0, k_{01}), \text{Enc}(k_1, k_{11})\}$.

Thus, after the second execution of OT R receives $c_{x_{\lambda-1}}$ that contains the ciphertexts described above where only one of these can be decrypted using the key k obtained in the first execution of OT. The obtained plaintext corresponds to the key $k_{x_\lambda x_{\lambda-1}}$ if $\text{Prefix}(x, 2) \in \text{Prefix}(Y, 2)$, to the sink key k_2^* otherwise. The same process is iterated for all the levels of G . More generally, if $\text{Prefix}(x, j) \in \text{Prefix}(Y, j)$ then after the j -th execution of OT R can compute the key $k_{x_\lambda x_{\lambda-1} \dots x_{\lambda-j}}$ using the key obtained in the previous phase. Conversely if $\text{Prefix}(x, j) \notin \text{Prefix}(Y, j)$ then the sink

key k_j^* is obtained by R. We observe that after every execution of OT R does not know which ciphertext can be decrypted using the key obtained in the previous phase, therefore he will try to decrypt all the ciphertext until the decryption procedure is successful. To avoid adding yet more indexes to the (already heavy) notation of our protocol we deal with this using a private-key encryption scheme with efficiently verifiable range. We note that this is not necessary and that when implementing the protocol one can instead use the *point-and-permute technique* [BMR90]. This, and other optimisations and extensions of our protocol, are described in Section 5.

2.2.3 Third step: obtain the correct share

In this step S encrypts the output string γ_0 using the key k_λ^* and uses all the other keys in the last level of G to encrypt the output string γ_1 .¹⁰ At this point the receiver can only decrypt either the ciphertext that contains γ_0 if $x \notin Y$ or one (and only one) of the ciphertexts that contain γ_1 if $x \in Y$. In the protocol that we have described so far R does not know which ciphertext can be decrypted using the key that he has obtained. Also in this case we can use a point-and-permute technique to allow R to identify the only ciphertext that can be decrypted using his key.

On the need for padding As describe earlier, we might need to add some padding to the OT sender’s inputs. To see why we need this we make the following observation. We recall that in the i -th OT execution the sender computes an encryption of the keys in the level i of the artificial graph G using the keys of the previous level $(i - 1)$.¹¹ As a result of this computation the sender obtains the pair (c_0^i, c_1^i) , that will be used as input of the i -th OT execution, where c_0^i (as well as c_1^i) contains a number of encryptions that depends upon the number of vertices on level $(i - 1)$ of G . We observe that this leaks information about the structure of G to the receiver, and therefore leaks information about the elements that belong to Y . Considering the example in Fig. 1, if we allow the receiver to learn that the 2nd level only contains 3 nodes, then the receiver would learn that all the elements of Y have the two most significant bits equal to either t or t' for some $t, t' \in \{0, 1\}^2$ (in Fig.1 for example we have $t = 01$ and $t' = 11$; note however that the receiver would not learn the actual values of t and t').

We note that the technique described in this section can be seen as a special (and simpler) example of securely evaluating a branching program. Secure evaluation of branching programs has previously been considered in [IP07, MN12]: unfortunately these protocols cannot be instantiated using OT-extension and therefore will not lead to practically efficient protocols (the security of these protocols is based on *strong* OT which, in a nutshell, requires the extra property that when executing several OTs in parallel, the receiver should not be able to correlate the answers with the queries beyond correlations which follow from the output).

Finally, we note that the work of Chor et al. [CGN98] uses a data structure similar to the one described here to achieve private information retrieval (PIR) based on keywords. The main difference between keyword based PIR and PSM is that in PSM the receiver should not learn any other information about the data stored by the sender, so their techniques cannot be directly applied to our setting.

¹⁰The key k_λ^* could not exists; e.g. if Y contains all the strings of λ bits.

¹¹The only exception is the first OT execution where just two keys are used as input.

3 Definitions and tools

We denote the security parameter by s and use “ $||$ ” as concatenation operator (i.e., if a and b are two strings then by $a||b$ we denote the concatenation of a and b). For a finite set Q , $x \leftarrow Q$ denotes a sampling of x from Q with uniform distribution. We use the abbreviation PPT that stands for probabilistic polynomial time. We use $\text{poly}(\cdot)$ to indicate a generic polynomial function. We assume the reader to be familiar with standard notions such as *computational indistinguishability* and the *real world/ideal world* security definition for secure two-party computation (see Appendix C for the actual definitions).

3.1 Special private-key encryption

In our construction we use a private-key encryption scheme with two additional properties. The first is that given the key k , it is possible to efficiently verify if a given ciphertext is in the range of k . With the second property we require that an encryption under one key will fall in the range of an encryption under another key with negligible probability

As discussed in [LP09], it is easy to obtain a private-key encryption scheme with the properties that we require. According to [LP09, Definition 2] we give the following definition.

Definition 1. Let $\text{Sym} = (\text{Gen}, \text{Enc}, \text{Dec})$ be a private-key encryption scheme and denote the range of a key in the scheme by $\text{Range}_s(k) = \{\text{Enc}(k, x)\}_{x \in \{0,1\}^s}$. Then

1. We say that Sym has an efficiently verifiable range if there exists a ppt algorithm M such that $M(1^s, k, c) = 1$ if and only if $c \in \text{Range}_s(k)$. By convention, for every $c \notin \text{Range}_s(k)$, we have that $\text{Dec}(k, c) = \perp$.
2. We say that Sym has an elusive range if for every probabilistic polynomial-time machine \mathcal{A} , there exists a negligible function $\nu(\cdot)$ such that $\text{Prob}_{k \leftarrow \text{Gen}(1^s)}[\mathcal{A}(1^s) \in \text{Range}_s(k)] < \nu(s)$.

Most of the the well known techniques used to construct a private-key encryption scheme (e.g. using a PRF) can be used to obtain a *special* private-key encryption scheme as well. The major difference is that a special encryption scheme has (in general) ciphertexts longer than a standard encryption scheme.

4 Our protocol Π^ϵ

In this section we provide the formal description of our protocol $\Pi^\epsilon = (\text{S}, \text{R})$ for the *set-membership* functionality $\mathcal{F}^\epsilon = (\mathcal{F}_S^\epsilon, \mathcal{F}_R^\epsilon)$ where

$$\mathcal{F}_S^\epsilon: (\{\{0,1\}^\lambda\}^M \times (\gamma^0, \gamma^1)) \times \{0,1\}^\lambda \longrightarrow \perp \quad \text{and}$$

$$\mathcal{F}_R^\epsilon: (\{\{0,1\}^\lambda\}^M \times (\gamma^0, \gamma^1)) \times \{0,1\}^\lambda \longrightarrow \{\gamma^0, \gamma^1\}$$

$$(Y, (\gamma^0, \gamma^1), x) \longmapsto \begin{cases} \gamma^1 & \text{if } x \in Y \\ \gamma^0 & \text{otherwise} \end{cases}$$

Where γ^0 and γ^1 are arbitrary strings and are part of the sender’s input. Therefore our scheme protects both Y and γ^{1-b} , when γ^b is received by R.

For the formal description of Π^ϵ , we collapse the first and the second step showed in the information description of Section 2 into a single one. That is, instead of constructing the graph G , the sender only computes the keys at level i in order to feed the i -th OT execution with the correct inputs. The way in which the keys are computed is the same as the vertices for G are computed, we just do not need to physically construct G to allow S to efficiently compute the keys. In our construction we make use of the following tools.

1. A protocol $\Pi_{\mathcal{OT}} = (\mathcal{S}_{\mathcal{OT}}, \mathcal{R}_{\mathcal{OT}})$ that securely (according to Definition 3) computes the following functionality

$$\begin{aligned} \mathcal{F}_{\mathcal{OT}}: (\{0, 1\}^* \times \{0, 1\}^*) \times \{0, 1\} &\longrightarrow \{\perp\} \times \{0, 1\}^* \\ ((c_0, c_1), b) &\longmapsto (\perp, c_b). \end{aligned}$$

2. A symmetric key encryption scheme $\text{Sym} = (\text{Gen}, \text{Enc}, \text{Dec})$ with efficiently verifiable and elusive range.
3. In our construction we make use of the following function:

$$\begin{aligned} \delta: \mathbb{N} &\longrightarrow \mathbb{N} \\ i &\longmapsto \min\{2^i, |Y|\}. \end{aligned}$$

This function computes the maximum number of vertices that can appear in the level i of the graph G . As discussed before, the structure of G leaks information about Y . In order to avoid this information leakage about Y , it is sufficient to add some padding to the OT sender's input so that the input size become $|Y|$. Indeed, as observed above, every level contains at most $|Y|$ vertices. Actually, it is easy to see that $\min\{|Y|, 2^i\}$ represents a better upper bound on the number of vertices that the i -th level can contain. Therefore, in order to compute the size of the padding for the sender's input we use the function δ .

4.1 Formal description

Common input: security parameter s and λ .

S's input: a set Y of size M , $\gamma^0 \in \{0, 1\}^s$ and $\gamma^1 \in \{0, 1\}^s$.

R's input: an element $x \in \{0, 1\}^\lambda$.

First stage

1. For $i=1, \dots, \lambda$, S computes the sink key $k_i^* \leftarrow \text{Gen}(1^s)$.
2. S computes $k_0 \leftarrow \text{Gen}(1^s), k_1 \leftarrow \text{Gen}(1^s)$. For $b = 0, 1$, if $b \notin \text{Prefix}(Y, 1)$ then set $k_b = k_1^{*12}$. Set $(c_0^1, c_1^1) = (k_0, k_1)$.
3. S and R execute $\Pi_{\mathcal{OT}}$, where S acts as the sender $\mathcal{S}_{\mathcal{OT}}$ using (c_0^1, c_1^1) as input and R acts as the receiver $\mathcal{R}_{\mathcal{OT}}$ using x_λ as input. When the execution of $\Pi_{\mathcal{OT}}$ ends R obtains $\kappa_1 := c_{x_\lambda}^1$.

Second stage For $i = 2, \dots, \lambda$:

1. S executes the following steps.

¹²We observe that if Y is not empty (like in our case) then there exists at most one bit b s.t. $b \in \text{Prefix}(Y, 1)$.

- 1.1. Define the empty list c_0^i and for all $t \in \text{Prefix}(Y, i - 1)$ execute the following steps.
 - If $t||0 \in \text{Prefix}(Y, i)$ then compute $k_{t||0} \leftarrow \text{Gen}(1^s)$ and add $\text{Enc}(k_t, k_{t||0})$ to the list c_0^i . Otherwise, if $t||0 \notin \text{Prefix}(Y, i)$ then compute and add $\text{Enc}(k_t, k_i^*)$ to the list c_0^i .
- 1.2. If $|c_0^i| < \delta(i - 1)$ then execute the following steps.
 - Compute and add $\text{Enc}(k_{i-1}^*, k_i^*)$ to the list c_0^i .
 - For $j = 1, \dots, \delta(i - 1) - |c_0^i|$ compute and add $\text{Enc}(\text{Gen}(1^s), 0)$ to c_0^i .¹³
- 1.3. Permute the elements inside c_0^i .
- 1.4. Define the empty¹⁴ list c_1^i and for all $t \in \text{Prefix}(Y, i - 1)$ execute the following step.
 - If $t||1 \in \text{Prefix}(Y, i)$ then compute $k_{t||1} \leftarrow \text{Gen}(1^s)$ and add $\text{Enc}(k_t, k_{t||1})$ to the list c_1^i . Otherwise, if $t||1 \notin \text{Prefix}(Y, i)$ compute and add $\text{Enc}(k_t, k_i^*)$ to the list c_1^i .
- 1.5. If $|c_1^i| < \delta(i - 1)$ then execute the following steps.
 - Compute and add $\text{Enc}(k_{i-1}^*, k_i^*)$ to the list c_1^i .
 - For $j = 1, \dots, \delta(i - 1) - |c_1^i|$ compute and add $\text{Enc}(\text{Gen}(1^s), 0)$ to c_1^i .
- 1.6. Permute the elements inside c_1^i .
2. S and R execute $\Pi_{\mathcal{OT}}$, where S acts as the sender $S_{\mathcal{OT}}$ using (c_0^i, c_1^i) as input and R acts as the receiver $R_{\mathcal{OT}}$ using $x_{\lambda-i+1}$ as input. When the execution of $\Pi_{\mathcal{OT}}$ ends, R obtains $c_{x_{\lambda-i+1}}^i$.

Third stage

1. S executes the following steps.
 - 1.1. Define the empty list l .
 - 1.2. For every $t \in \text{Prefix}(Y, \lambda)$ compute and add $\text{Enc}(k_t, \gamma^1)$ to l .
 - 1.3. If $|l| < 2^\lambda$ then compute and add $\text{Enc}(k_\lambda^*, \gamma^0)$ to l .
 - 1.4. Permute the elements inside l and send l to R.
2. R, upon receiving l acts as follows.
 - 2.1. For $i = 2, \dots, \lambda$ execute the following steps.
 - For every element t in the list $c_{x_{\lambda-i+1}}^i$ compute $\kappa \leftarrow \text{Dec}(\kappa_{i-1}, t)$. If $\kappa \neq \perp$ then set $\kappa_i = \kappa$.
 - 2.2. For all $e \in l$ compute $\text{out} \leftarrow \text{Dec}(\kappa_\lambda, e)$ and output out if and only if $\text{out} \neq \perp$.

Theorem 1. *Suppose $\Pi_{\mathcal{OT}}$ securely computes the 1-out-of-2 OT functionality $\mathcal{F}_{\mathcal{OT}}$ and Sym is a symmetric key encryption scheme with efficiently verifiable range and elusive range, then Π^ϵ securely computes the functionality \mathcal{F}^ϵ .*

We refer the reader to App. B for the formal proof of this theorem.

¹³In this step, as well as in the step 1.e of this stage, the function δ is used to compute the right amount of fake encryption to be added to the list that will we used as input of $R_{\mathcal{OT}}$. The fake encryptions encrypts the value 0, but of course any other value could be used.

¹⁴The following three steps are equal to the previous three steps (1.a, 1.b and 1.c), the only difference is that $t||1$ is considered instead of $t||0$.

4.2 Round complexity: parallelizability of our scheme

In the description of our protocol in Sec 4.1 we have the sender and the receiver engaging λ sequential OT executions. We now show that this is not necessary since the OT executions can be easily parallelized, given that each execution is independent from the other. That is, the output of a former OT execution is not used in a latter execution. For simplicity, we assume that $\Pi_{\mathcal{OT}}$ consists of just two rounds, where the first round goes from the receiver to the sender, and the second goes in the opposite direction. We modify the description of the protocol of Sec 4.1 as follows.

- Step 3 of the *first stage* and step 2 of the *second stage* are moved to the beginning of the *third stage*.
- When S sends the last round of $\Pi_{\mathcal{OT}}$, he also performs the step 1 of the *third stage*. Therefore the list l is sent together with the last rounds of the λ $\Pi_{\mathcal{OT}}$ executions.

Roughly speaking, in this new protocol S first computes all the inputs $(k_0, k_1, c_0^1, c_1^1, \dots, c_0^\lambda, c_1^\lambda)$ for the OTs. Then, upon receiving the λ first rounds of $\Pi_{\mathcal{OT}}$ computed by R using as input the bits of x , S sends λ second round of $\Pi_{\mathcal{OT}}$ together with the list l . We observe that in this case the S's inputs to the λ executions of $\Pi_{\mathcal{OT}}$ can be pre-computed *before* any interaction with R begins.

5 Optimisations and extensions

5.1 Point and permute

In our protocol the receiver must decrypt every ciphertext at every layer to identify the correct one. This is suboptimal both because of the number of decryptions and because encryptions that have efficiently verifiable range necessarily have longer ciphertexts. This overhead can be removed using the standard *point-and-permute technique* [BMR90] which was introduced in the context of garbled circuits. Using this technique we can add to each key in each layer a *pointer* to the ciphertext in the next layer which can be decrypted using this key. This has no impact on security.

5.2 One-time pad

It is possible to reduce the communication complexity of our protocol by using *one-time pad encryption* in the last $\log s$ layers of the graph, in the setting where the output values γ^0, γ^1 are such that $|\gamma^b| < s$. For instance, if the output values are bits (in case we combine our PSM with a GMW-style protocol), then the keys (and therefore the ciphertexts) used in the last layer of the graph only need to be 1 bit long. Unfortunately, since the keys in the second to last layer are used to mask up to two keys in the last layer, the keys in the second to last layer must be of length 2 and so on, which is why this optimisation only gives benefits in the last $\log s$ layer of the graph.

5.3 PSM with secret shared input

Our PSM protocol produces an output which can be post-processed using other 2PC protocols. It is natural to ask whether it is possible to design efficient PSM protocols that also work on encrypted or secret-shared inputs. We note here that our protocol can also be used in the setting in which the input string x is *bit-wise secret-shared* between the sender and the receiver i.e., the receiver knows a share r and the sender knows a share s s.t., $r \oplus s = x$. The protocol does not change for the

receiver, who now inputs the bits of $r = r_\lambda, \dots, r_1$ to the λ one-out-of-two OTs (instead of the bits of x as in the original protocol). The sender, at each layer i , will follow the protocol as described above if $s_i = 0$ and instead swap the inputs to the OT if $s_i = 1$. It can be easily verified that the protocol still produces the correct result and does not leak any extra information.

5.4 Keyword search

Our PSM protocol outputs an encryption of a bit indicating whether $x \in Y$ or not. The protocol can be easily modified to output a value dependent on x itself and therefore implement “encrypted keyword search”. That is, instead of having only two output strings γ^1, γ^0 representing membership and non-membership respectively, we can have $|Y| + 1$ different output strings (one for each element $y \in Y$ and one for non-membership). This can be used for instance in the context where Y is a database containing id’s y and corresponding values $v(y)$, and the output of the protocol should be an encryption of the value $v(x)$ if $x \in Y$ or a standard value $v(\perp)$ if $x \notin Y$. The modification is straightforward: instead of using all the keys in the last layer of the graph to encrypt the same value γ^1 , use each key k_y to encrypt the corresponding value $v(y)$ and the sink key (which is used to encrypt γ^0 in our protocol) to encrypt the value $v(\perp)$.

5.5 PSI from PSM

We can follow the same approach of PSZ [PSSZ15, PSZ14] to turn our PSM protocol into a protocol for PSI. Given a receiver with input X and a sender with input Y the trivial way to construct PSI from PSM is to run $|X|$ copies of PSM, where in each execution the receiver inputs a different x from X and where the sender always inputs her entire set Y . As described above, the complexity of our protocol (as the complexity of the PSM protocol of PSZ) is proportional in the size of $|Y|$, so this naïve approach leads to quadratic complexity $O(|X| \cdot |Y|)$. PSZ deals with this using *hashing* i.e., by letting the sender and the receiver locally preprocess their inputs X, Y before engaging in the PSM protocols. The different hashing techniques are explained and analysed in [PSZ16, Section 3]. We present the intuitive idea and refer to their paper for details: in PSZ the receiver uses *Cuckoo hashing* to map X into a vector X' of size $\ell = O(|X|)$ such that all elements of X are present in X' and such that every $x'_i \in X'$ is either an element of X or a special \perp symbol. The sender instead maps her set Y into $\ell = |X'|$ small buckets Y'_1, \dots, Y'_ℓ such that every element $y \in Y$ is mapped into the “right bucket” i.e., the hashing has the property that if $y = x'_i$ for some i then y will end up in bucket Y'_i (and potentially in a few other buckets). Now PSZ uses the underlying PSM protocol to check whether x'_i is a member of Y'_i (for all i ’s), thus producing the desired result. The overall protocol complexity is now $O(\sum_{i=1}^{\ell} |X'| \cdot |Y'_i|)$ which (by careful choice of the hashing parameters) can be made sub-quadratic. In particular, if one is willing to accept a small (but not negligible) failure probability, the overall complexity becomes only linear in the input size. Since this technique is agnostic of the underlying PSM protocol, we can apply the same technique to our PSM protocol to achieve a PSI protocol that produces encrypted output.

6 Applications

The major advantage provided by Π^ϵ is that the output of the receiver can be an arbitrary value chosen by the sender as a function of x for each value $x \in Y \cup \{\perp\}$. This is in contrast with most of the approaches for set membership, where the value obtained by the receiver is a fixed value (e.g. 0)

when $x \in Y$, or some random value otherwise. We now provide two examples of how our protocol can be used to implement more complex secure set operations. The examples show some guiding principles that can be used to design other applications based on our protocol.

Without loss of generality in the following applications only the receiver will learn the output of the computation. Moreover we assume that the size of X and Y is equal to the same value M .¹⁵ Also for simplicity we will describe our application using the naïve PSI from PSM construction with quadratic complexity, but using the PSZ approach, as described in Sec. 5, it is possible to achieve linear complexity using hashing techniques. Finally, in both our applications we exploit the fact that additions can be performed locally (and for free) using secret-sharing based 2PC. In applications in which round complexity is critical, the protocols can be redesigned using garbled circuits computing the same functionality, since the garbled circuit can be sent from the sender to the other messages of the protocol. However in this case additions have to be performed inside the garbled circuit.

6.1 Computing statistics of the private intersection

Here we want to construct a protocol where sender and receiver have as input two sets, X and Y respectively, and want to compute some statistics on the intersections of their sets. For instance the receiver has a list of id's X and that the sender has a list of id's Y and some corresponding values $v(Y)$ (thus we use the variant of our protocol for *keyword search* described in Section 5). At the end of the protocol the receiver should learn the average of $v(X \cap Y)$ (and not $|X \cap Y|$).

The main idea is the following: the sender and the receiver run M executions of our protocol where the receiver inputs a different x_i from X in each execution. The sender always inputs the same set Y , and chooses the $|Y|+1$ outputs γ_i^y for all $y \in Y \cup \{\perp\}$ for all $i = 1, \dots, M$ in the following way: γ_i^y is going to contain two parts, namely an arithmetic secret sharing of the bit indicating whether $x_i \in Y$ and an arithmetic secret sharing of the value $v(y)$. The arithmetic secret sharing will be performed using a modulo N large enough such that $N > M$ and $N > M \cdot V$ where V is some upper bound on $v(y)$ so to be sure that no modular reduction will happen when performing the addition of the resulting shares. Concretely the sender sets $\gamma_i^y = (-u_i^2 + 1 \pmod N, -v_i^2 + v(y) \pmod N)$ for all $y \in Y$ and $\gamma_i^\perp = (-u_i^2 \pmod N, -v_i^2 \pmod N)$. After the protocol the receiver defines her shares u_i^1, v_i^1 to be the shares contained in her output of the PSM protocol, and then both parties add their shares locally to obtain secret sharing of the size of the intersection and of the sum of the values i.e., $U^1 = \sum_i u_i^1$, $V^1 = \sum_i v_i^1$, $U^2 = \sum_i u_i^2$, and $V^2 = \sum_i v_i^2$. Now the parties check if (U^1, U^2) is a sharing of 0 and, if not, they compute and reveal the result of the computation $\frac{V^1+V^2}{U^1+U^2}$. Both these operations can be performed using efficient two-party protocols for comparison and division such as the one in [T⁺07, DNT12].

6.2 Threshold PSI

In this example we design a protocol $\Pi^t = (P_1^t, P_2^t)$ that securely computes the functionality $\mathcal{F}^t = (\mathcal{F}_{P_1^t}^t, \mathcal{F}_{P_2^t}^t)$ where

$$\mathcal{F}_{P_1^t}^t : \{\{0, 1\}^\lambda\}^M \times \{\{0, 1\}^\lambda\}^M \longrightarrow \perp$$

¹⁵We assume this only to simplify the protocol description, indeed our protocol can be easily instantiated when the two sets have different size.

and

$$\mathcal{F}_{P_2^t}^t : \{\{0, 1\}^\lambda\}^M \times \{\{0, 1\}^\lambda\}^M \longrightarrow \{\{0, 1\}^\lambda\}^*$$

$$(S_1, S_2) \longmapsto \begin{cases} S_1 \cap S_2 & \text{if } |S_1 \cap S_2| \geq t \\ \perp & \text{otherwise} \end{cases}$$

That is, the sender and the receiver have on input two sets, S_1 and S_2 respectively, and the receiver should only learn the intersection between these two sets if the size of the intersection is greater or equal than a fixed (public) threshold value t . In the case that the size of the intersection is smaller than t , then no information about S_1 is leaked to P_2^t and no information about S_2 is leaked to P_1^t . (This notion was recently considered in [HOS17] in the context of privacy-preserving ride-sharing).

As in the previous example, the sender and the receiver run M executions of our protocol where the receiver inputs a different x_i from S_2 in each execution. The sender always inputs the same set S_1 , and chooses the two outputs γ_i^0, γ_i^1 in the following way: γ_i^b is going to contain two parts, namely an arithmetic secret sharing of 1 if $x_i \in Y$ or 0 otherwise, as well as encryption of the same bit using a key k . The arithmetic secret sharing will be performed using a modulus larger than M , so that the arithmetic secret sharings can be added to compute a secret-sharing of the value $|S_1 \cap S_2|$ with the guarantee that no overflow will occur. Then, the sender and the receiver engage in a secure-two party computation of a function that outputs the key k to the receiver if and only if $|S_1 \cap S_2| > t$. Therefore, if the intersection is larger than the threshold now the receiver can decrypt the ciphertext part of the γ values and learn which elements belong to the intersection. The required 2PC is a simple comparison with a known value (the threshold is public) which can be efficiently performed using protocols such as [GSV07, LT13].

7 Acknowledgments

This research received funding from: COST Action IC1306; the Danish Independent Research Council under Grant-ID DFF-6108-00169 (FoCC); the European Union’s Horizon 2020 research and innovation programme under grant agreements No 731583 (SODA) and No 780477 (PRIViLEDGE); “GNCS - INdAM”. The work of 1st author has been done in part while visiting Aarhus University, Denmark.

A Complexity analysis

We focus our analysis of the protocol described in Sec 4.1 without taking into account the many possible optimisations showed in Sec. 5. In Π^ϵ , sender and receiver run λ executions of a 1-out-of-2 OT; in addition, they perform some symmetric key operations. More precisely, in order to compute the inputs for the i -th OT executions, with $i \in \{2, \dots, \lambda\}$, S computes $2 \cdot \min\{2^{i-1}, |Y|\}$ encryptions using the private-key encryption scheme Sym . We now observe that each encryption could contain a different key, and that this key needs to be generated by running $\text{Gen}(\cdot)$.¹⁶ This means that $4M$ represents an upper bound on the number of symmetric key operations performed by S to compute the input of one OT execution. Moreover, in the last interaction with R , S computes M encryptions. Therefore, an upper bound on the number symmetric key operations performed by S is $(\lambda - 1) \cdot 4M + M + 2 \approx \lambda \cdot 4M$, where 2 represents the cost of running $\text{Gen}(\cdot)$ twice in order to compute the two keys required to feed the first OT execution¹⁷. In every OT execution i , with $i \in \{2, \dots, \lambda\}$, R receives $\min\{2^{i-1}, |Y|\}$ encryptions, and tries to decrypt all of them. Moreover, in the last interaction with S , R receives M encryptions and tries to decrypt all of them as well. This means that the upper bound on the number of symmetric key operations made by R is $(\lambda - 1) \cdot M + M = \lambda \cdot M$. Following [PSZ16] we assume that 3 symmetric key operations are required for one OT execution. Therefore the total amount of symmetric key operations is $\lambda M 4 + 3\lambda$ for the sender and $\lambda M + 3\lambda$ for the receiver. In order to compare the efficiency of our protocol with the PSI protocols provided in [PSZ16] and to be consistent with their complexity analysis, we consider only the computation complexity for the party with the majority of the workload in the comparison. In Table 1 of Sec. 1 we have compared the computation (and the communication) complexity of our protocol with the circuit-based PSI approaches (which can be combined with further postprocessing) considered in [PSZ16]. More precisely, we compare the sort-compare-shuffle (SCS) circuit of [HEK12] and the pairwise-comparison (PWC) circuit proposed in [PSZ16] with our approach for PSM.

As showed in Table 1, our protocol has better performance than all the circuit-based PSI approaches (which can be combined with further postprocessing) considered in [PSZ16]. We note that, as described in Sec. 4.4 of [PSZ16], the approach based on evaluating the OPRF inside circuit is faster than any other PSI protocols if one set is much smaller than the other (like in the case of PSM), but in this case the output will necessarily leak to the receiver, which prevents composition with further 2PC protocol. We refer the reader to Table 7 of [PSZ16] for a detailed efficiency comparison between different PSI protocols. Finally, we observe that the complexities analysis proposed in [PSZ16] is related to PSI protocols, while in this section we have only compared the efficiency of the PSM subprotocol.

A.1 Communication complexity

The communication complexity of our protocol is dominated by the communication complexity of the underlying OT protocol $\Pi_{OT} = (S_{OT}, R_{OT})$. Let $\text{sOT}(D)$ be the amount of data exchanged between S_{OT} and R_{OT} when S_{OT} uses an input of size D , and let $\text{sSYM}(A)$ be the size of a ciphertext for the encryption scheme Sym when a plaintext of size A is used. Then the communication

¹⁶We recall that $|Y| = M$ and that λ is the bit size of a set element.

¹⁷In this section, without loss of generality, we assume that to encrypt a message of size λ it is sufficient to run the encryption algorithm Enc only once.

complexity of our protocol is

$$\lambda \cdot \text{sOT}(2 \cdot M \cdot \text{sSYM}(\lambda)) + M \cdot \text{sSYM}(\lambda)$$

where $2 \cdot M$ is the number of ciphertexts used as input of OT and M is the amount of ciphertexts that are sent in the last interaction between \mathbf{S} and \mathbf{R} . If we assume that a ciphertext for \mathbf{Sym} is roughly of size s , and that $\Pi_{\mathcal{OT}}$ has a communication complexity that is approximately close to the size of the input used¹⁸, we obtain that the overall communication complexity of our protocol is well approximated by $\lambda Ms2 + Ms$, that is comparable to the communication complexity of the approaches proposed in [PSZ16].

B Security proof of Theorem 1

Proof. In order to prove the security of Π^ϵ , according to Def. 3 we need show two probabilistic polynomial-time algorithms $\mathcal{S}_\mathbf{S}$ and $\mathcal{S}_\mathbf{R}$ called simulators, such that the following two conditions hold:

$$\begin{aligned} \{(\mathcal{S}_\mathbf{S}(1^s, Y, \gamma^0, \gamma^1, \mathcal{F}_\mathbf{S}^\epsilon(Y, \gamma^0, \gamma^1, x)), \mathcal{F}^\epsilon(Y, \gamma^0, \gamma^1, x))\}_{\{Y, x, s\}} \approx \\ \{\text{view}_\mathbf{S}^{\Pi^\epsilon}(1^s, Y, \gamma^0, \gamma^1, x), \text{output}^{\Pi^\epsilon}(1^s, Y, \gamma^0, \gamma^1, x)\}_{\{Y, x, s\}} \end{aligned} \quad (1)$$

$$\begin{aligned} \{(\mathcal{S}_\mathbf{R}(1^s, x, \mathcal{F}_\mathbf{R}^\epsilon(Y, \gamma^0, \gamma^1, x)), \mathcal{F}^\epsilon(Y, \gamma^0, \gamma^1, x))\}_{\{Y, x, s\}} \approx \\ \{\text{view}_\mathbf{R}^{\Pi^\epsilon}(1^s, Y, \gamma^0, \gamma^1, x), \text{output}^{\Pi^\epsilon}(1^s, Y, \gamma^0, \gamma^1, x)\}_{\{Y, x, s\}} \end{aligned} \quad (2)$$

where $Y \in \{\{0, 1\}^*\}^*$, $x \in \{0, 1\}^*$, and $s \in \mathbb{N}$.

Therefore we divide our proof in two parts. In the former we show a PPT algorithm $\mathcal{S}_\mathbf{S}$ that satisfies the property of the first point, and then a PPT algorithm $\mathcal{S}_\mathbf{R}$ that satisfies the requirement of the second point. Moreover, in order to make the security proof of our scheme easier, without loss of generality we assume \mathbf{Sym} to be secure in the setting where the challenge messages \mathbf{m}_0 and \mathbf{m}_1 are lists of λ values. That is $\mathbf{m}_0 = \{m_0^1, \dots, m_0^\lambda\}$ and $\mathbf{m}_1 = \{m_1^1, \dots, m_1^\lambda\}$. The challenger, upon receiving these lists picks $b \leftarrow \{0, 1\}$, defines an empty list \mathbf{cx} and for $i = 1, \dots, \lambda$ acts as follows:

1. computes $k_i \leftarrow \text{Gen}(1^s)$;
2. computes $\text{Enc}(k_i, m_b^i)$ and adds it to \mathbf{cx} .

The aim of the adversary is to guess the bit b having on input just \mathbf{m}_0 , \mathbf{m}_1 , \mathbf{cx} and an auxiliary input z .

B.1 $\mathcal{S}_\mathbf{S}$ description and proof of indistinguishability

$\mathcal{S}_\mathbf{S}$ runs \mathbf{S} with some randomness r and the input Y . At this point $\mathcal{S}_\mathbf{S}$ needs a strategy to act as a receiver of OT in all the λ OT executions (without the receiver's input x). In order to do that, $\mathcal{S}_\mathbf{S}$ runs the simulator of $\Pi_{\mathcal{OT}}$, that we call $\mathcal{S}_{\mathbf{S}_{\mathcal{OT}}}$ (and that exists by assumption), in every OT execution. We observe that in order to run $\mathcal{S}_{\mathbf{S}_{\mathcal{OT}}}$ in the i -th OT execution the inputs c_0^i and c_1^i

¹⁸This is actually true for the most common implementations of OT (OT extension).

need to be known. Clearly those values can be efficiently computed since the randomness r and the input Y used to run \mathcal{S} are known.

We now show more formally how $\mathcal{S}_{\mathcal{S}}$ works. Let $\mathcal{S}_{\mathcal{S}_{\mathcal{OT}}}$ be such that

$$\{\mathcal{S}_{\mathcal{S}_{\mathcal{OT}}}(1^s, (c_0, c_1), \perp), \mathcal{F}_{\mathcal{OT}}((c_0, c_1), b)\}_{\{c_0, c_1, b, s\}} \approx \{\text{view}_{\mathcal{S}_{\mathcal{OT}}}^{\Pi_{\mathcal{OT}}}(1^s, (c_0, c_1), b), \text{output}^{\Pi_{\mathcal{OT}}}(1^s, (c_0, c_1), b)\}_{\{c_0, c_1, b, s\}}$$

where $c_0, c_1 \in \{0, 1\}^*$, $b \in \{0, 1\}$, and $s \in \mathbb{N}$. $\mathcal{S}_{\mathcal{S}}$, on input Y and 1^s executes the following steps.

1. pick a $r \leftarrow \{0, 1\}^s$ and run \mathcal{S} on input $1^s, Y$ using r as a randomness.
2. For every OT execution i , with $i = 1, \dots, \lambda$, run $\mathcal{S}_{\mathcal{S}_{\mathcal{OT}}}$ on input $1^s, c_0^i$ and c_1^i , where c_0^i and c_1^i are computed using the same procedure that \mathcal{S} uses.
3. Continue the execution against \mathcal{S} as \mathcal{R} would do.

In order to conclude this first part of the proof we just need to prove the following lemma.

Lemma 1.

$$\{(\mathcal{S}_{\mathcal{S}}(1^s, Y, \mathcal{F}_{\mathcal{S}}^{\in}(Y, \gamma^0, \gamma^1, x)), \mathcal{F}^{\in}(Y, \gamma^0, \gamma^1, x))\} \approx \{\text{view}_{\mathcal{S}}^{\Pi^{\in}}(1^s, Y, \gamma^0, \gamma^1, x), \text{output}^{\Pi^{\in}}(1^s, Y, \gamma^0, \gamma^1, x)\}$$

where $Y \in \{\{0, 1\}^*\}^*$, $x \in \{0, 1\}^*$, and $s \in \mathbb{N}$.¹⁹

Proof. The proof goes through hybrid arguments starting from the real execution of Π^{\in} . We gradually modify the execution until the input of \mathcal{R} is not needed anymore in such a way that the final hybrid represents the simulator $\mathcal{S}_{\mathcal{S}}$. We denote with $\text{OUT}_{\mathcal{S}}^{\mathcal{H}_i}(1^s)$ the view of \mathcal{S} in the hybrid experiment \mathcal{H}_i with $i \in \{0, \dots, \lambda\}$. The hybrid experiments that we consider are the following.

1. \mathcal{H}_0 is identical to the real execution of Π^{\in} . More precisely \mathcal{H}_0 runs \mathcal{S} using fresh randomness and interacts with him as \mathcal{R} would do on input x .
2. \mathcal{H}_i proceeds according to \mathcal{H}_0 with the difference that in the first i OT executions $\mathcal{S}_{\mathcal{S}_{\mathcal{OT}}}$ is used.

Since \mathcal{F}^{\in} is a deterministic function we have that

$$\{\text{view}_{\mathcal{S}}^{\Pi^{\in}}(1^s, Y, \gamma^0, \gamma^1, x), \mathcal{F}^{\in}(Y, \gamma^0, \gamma^1, x)\} \equiv \{\text{view}_{\mathcal{S}}^{\Pi^{\in}}(1^s, Y, \gamma^0, \gamma^1, x), \text{output}^{\Pi^{\in}}(1^s, Y, \gamma^0, \gamma^1, x)\}.$$

Moreover we observe that

$$\{\text{OUT}_{\mathcal{H}_0}(1^s), \mathcal{F}^{\in}(Y, \gamma^0, \gamma^1, x)\} = \{\text{view}_{\mathcal{S}}^{\Pi^{\in}}(1^s, Y, \gamma^0, \gamma^1, x), \mathcal{F}^{\in}(Y, \gamma^0, \gamma^1, x)\}$$

and that

$$\{\text{OUT}_{\mathcal{H}_\lambda}(1^s), \mathcal{F}^{\in}(Y, \gamma^0, \gamma^1, x)\} = \{(\mathcal{S}_{\mathcal{S}}(1^s, Y, \gamma^0, \gamma^1, \mathcal{F}_{\mathcal{S}}^{\in}(Y, \gamma^0, \gamma^1, x)), \mathcal{F}^{\in}(Y, \gamma^0, \gamma^1, x))\}.$$

¹⁹To avoid overburdening the notation, here and in the rest of this paper, we omit to specify the inputs domain when it is clear from the context.

Therefore the only thing that remains to argue is that

$$\{\text{OUT}_{\mathcal{H}_{i-1}}(1^s), \mathcal{F}^\epsilon(Y, \gamma^0, \gamma^1, x)\} \approx \{\text{OUT}_{\mathcal{H}_i}(1^s), \mathcal{F}^\epsilon(Y, \gamma^0, \gamma^1, x)\}$$

for $i = 1, \dots, \lambda$. We now show that if this statement does not hold then we can construct an adversary \mathcal{A}^{SOT} that breaks the security of Π_{OT} against malicious sender. Let \mathcal{C}^{SOT} be the challenger for the security game w.r.t. the security of Π_{OT} against malicious sender; the reduction works as follows.

1. \mathcal{A}^{SOT} runs S with randomness r and interacts with him according to \mathcal{H}_{i-1} (\mathcal{H}_i) until the i -th OT execution.
2. At this point \mathcal{A}^{SOT} computes (c_0^i, c_1^i) and sends $((c_0^i, c_1^i), x_{\lambda-i+1})$ to \mathcal{C}^{SOT} .
3. \mathcal{A}^{SOT} then acts as a proxy between \mathcal{C}^{SOT} and S .
4. When the interaction between \mathcal{C}^{SOT} and S is over, \mathcal{A} continues the execution with S according to \mathcal{H}_{i-1} (\mathcal{H}_i).

The security proof ends with the observation that if \mathcal{C}^{SOT} has used the simulator \mathcal{S}_{SOT} then the joint distribution of the view of S and $\mathcal{F}^\epsilon(Y, \gamma^0, \gamma^1, x)$ corresponds to $\{\text{OUT}_{\mathcal{H}_i}^{\text{R}}(1^s), \mathcal{F}^\epsilon(Y, \gamma^0, \gamma^1, x)\}$, to $\{\text{OUT}_{\mathcal{H}_{i-1}}^{\text{R}}(1^s), \mathcal{F}^\epsilon(Y, \gamma^0, \gamma^1, x)\}$ otherwise. \square

\square

B.2 \mathcal{S}_{R} description and proof of indistinguishability

At a very high level, \mathcal{S}_{R} runs R with some randomness r and the input x . \mathcal{S}_{R} then needs a strategy to acts as a sender of OT in all the λ OT executions (without sender's input Y). In order to do that, \mathcal{S}_{R} runs the simulator of Π_{OT} , that we call $\mathcal{S}_{\text{R}_{\text{OT}}}$ in every OT execution²⁰. Moreover we need to feed the OT simulator with the correct input, depending on the value x . More precisely in the first OT execution $\mathcal{S}_{\text{R}_{\text{OT}}}$ is run by using as input a key k_1 . In the i -th OT execution (for $i = 2, \dots, \lambda$) the simulator will run using $x_{\lambda-i+1}$ and c^i . The ciphertext c^i contains encryptions of a fixed value, let us say 0, computed using a fresh secret key (different for every ciphertext) and one encryption of the key k_i using the key k_{i-1} . After the λ OT executions \mathcal{S}_{R} sends to R M encryptions of 0 using a randomly generated secret key (also in this case a different secret key is used for each encryption of 0) and the encryption of the message $\text{out} = \mathcal{F}_S^\epsilon(Y, \gamma^0, \gamma^1, x)$ using the key k_λ . We now show more formally how \mathcal{S}_{R} works. Let $\mathcal{S}_{\text{R}_{\text{OT}}}$ be such that

$$\{\mathcal{S}_{\text{R}_{\text{OT}}}(1^s, b, c_b), \mathcal{F}_{\text{OT}}((c_0, c_1), b)\}_{\{c_0, c_1, b, s\}} \approx \{\text{view}_{\text{R}_{\text{OT}}}^{\Pi_{\text{OT}}}(1^s, (c_0, c_1), b), \text{output}^{\Pi_{\text{OT}}}(1^s, (c_0, c_1), b)\}_{\{c_0, c_1, b, s\}}$$

where $c_0, c_1 \in \{0, 1\}^*$, $b \in \{0, 1\}$, and $s \in \mathbb{N}$.

\mathcal{S}_{R} , on input x , out and 1^s executes the following steps.

1. Compute $k_1 \leftarrow \text{Gen}(1^s)$ and run $\mathcal{S}_{\text{R}_{\text{OT}}}$ on input $(1^s, x_\lambda, k_1)$.

²⁰We recall that $\mathcal{S}_{\text{R}_{\text{OT}}}$ exists by assumption.

2. For $i = 2, \dots, \lambda$ execute the following steps.
 - 2.1. Define the empty list c^i . For $j = 1, \dots, \min\{2^i, |Y|\} - 1$ compute and add $\text{Enc}(\text{Gen}(1^s), 0)$ to c^i .
 - 2.2. Compute $k_i \leftarrow \text{Gen}(1^s)$, and add $\text{Enc}(k_{i-1}, k_i)$ to the list c^i .
 - 2.3. Permute the elements inside c^i .
 - 2.4. Run \mathcal{S}_{ROT} on input $(1^s, x_{\lambda-i+1}, c^i)$.
3. Define an empty list l .
4. For $i = 1, \dots, M - 1$ compute and add $\text{Enc}(\text{Gen}(1^s), 0)$ to l .
5. Add $\text{Enc}(k_\lambda, \text{out})$ to l .
6. Permute the elements inside l and send it.
7. Continue the execution according to R 's description.

In order to conclude this latter part of the proof we need to prove the following lemma.

Lemma 2.

$$\begin{aligned} & \{(\mathcal{S}_R(1^s, x, \text{out}, \mathcal{F}_R^\infty(Y, \gamma^0, \gamma^1, x)), \mathcal{F}^\infty(Y, \gamma^0, \gamma^1, x))\}_{\{Y, x, s\}} \approx \\ & \{\text{view}_R^{\Pi^\infty}(1^s, Y, \gamma^0, \gamma^1, x), \text{output}^{\Pi^\infty}(1^s, Y, \gamma^0, \gamma^1, x)\}_{\{Y, x, s\}} \end{aligned}$$

where $Y \in \{\{0, 1\}^*\}^*$, $x \in \{0, 1\}^*$, and $s \in \mathbb{N}$.

Proof. The proof goes through hybrid arguments starting from the real execution of Π^∞ . We gradually modify the execution until the input of S (Y) is not needed anymore such that the final hybrid would represent the simulator \mathcal{S}_R . We denote with $\text{OUT}_{\mathcal{H}_i}^R(1^s)$ the view of R in the hybrid experiment \mathcal{H}_i with $i \in \{0, \dots, \lambda\}$.

1. \mathcal{H}_0 is identical to the real execution of Π^∞ . More precisely \mathcal{H}_0 runs R using fresh randomness and interacts with him as S would do on input Y .
2. \mathcal{H}_1 proceeds according to \mathcal{H}_0 with the difference that in the first OT executions \mathcal{S}_{ROT} is used on input $(1^s, x_\lambda, k_1 \leftarrow \text{Gen}(1^s))$.
3. \mathcal{H}_i proceeds according to \mathcal{H}_1 with the difference that in the j -th OT executions, with $2 \leq j \leq i$, \mathcal{S}_{ROT} is run on input $(1^s, x_{\lambda-j+1}, c^j = c_{x_{\lambda-j+1}}^j)$.
4. \mathcal{H}^* proceeds according to \mathcal{H}_λ with the difference that in each OT execution i , with $2 \leq i \leq \lambda$, the input c^i for the simulator \mathcal{S}_{ROT} is computed as follows.
 - For $j = 1, \dots, \min\{2^i, |Y|\} - 1$ compute and add $\text{Enc}(\text{Gen}(1^s), 0)$ to c^i .
 - Compute $k_i \leftarrow \text{Gen}(1^s)$, and add $\text{Enc}(k_{i-1}, k_i)$ to the list c^i .
 - Permute the elements inside c^i .

Moreover the first step of the third stage is performed as follows.

- Define an empty list l .
- For $i = 1, \dots, M - 1$ compute and add $\text{Enc}(\text{Gen}(1^s), 0)$ to l .
- Add $\text{Enc}(k_\lambda, \text{out})$ to l .
- Permute the element inside l and send it to R.

Since \mathcal{F}^\in is deterministic we have that

$$\{\text{view}_R^{\Pi^\in}(1^s, Y, \gamma^0, \gamma^1, x), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\} \equiv \{\text{view}_R^{\Pi^\in}(1^s, Y, \gamma^0, \gamma^1, x), \text{output}^{\Pi^\in}(1^s, Y, \gamma^0, \gamma^1, x)\} .$$

Moreover we observe that

$$\{\text{OUT}_{\mathcal{H}_0}^R(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\} = \{\text{view}_R^{\Pi^\in}(1^s, Y, \gamma^0, \gamma^1, x), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\}$$

and that

$$\{\text{OUT}_{\mathcal{H}^*}^R(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\} = \{(\mathcal{S}_R(1^s, x, \text{out}, \mathcal{F}_R^\in(Y, \gamma^0, \gamma^1, x)), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x))\} .$$

Therefore there are two things that remain to argue:

1. $\{\text{OUT}_{\mathcal{H}_{i-1}}^R(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\} \approx \{\text{OUT}_{\mathcal{H}_i}^R(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\}$ for $i = 1, \dots, \lambda$ and
2. $\{\text{OUT}_{\mathcal{H}_\lambda}^R(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\} \approx \{\text{OUT}_{\mathcal{H}^*}^R(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\}$.

We now start by showing that if the first statement does not hold for $i = 1$, then we can construct an adversary $\mathcal{A}^{\mathcal{S}_{\mathcal{O}\mathcal{T}}}$ that breaks the security of $\Pi_{\mathcal{O}\mathcal{T}}$ against malicious receiver. Let $\mathcal{C}^{\mathcal{R}_{\mathcal{O}\mathcal{T}}}$ be the challenger for the security game w.r.t. the security of $\Pi_{\mathcal{O}\mathcal{T}}$ against malicious receiver. The reduction works as follows.

1. $\mathcal{A}^{\mathcal{R}_{\mathcal{O}\mathcal{T}}}$ runs R with randomness r , computes $k_0 \leftarrow \text{Gen}(1^s)$, $k_1 \leftarrow \text{Gen}(1^s)$ and sends $((k_0, k_1), x_\lambda)$ to $\mathcal{C}^{\mathcal{R}_{\mathcal{O}\mathcal{T}}}$.
2. $\mathcal{A}^{\mathcal{R}_{\mathcal{O}\mathcal{T}}}$ then acts as a proxy between $\mathcal{C}^{\mathcal{R}_{\mathcal{O}\mathcal{T}}}$ and R.
3. When the interaction between $\mathcal{C}^{\mathcal{S}_{\mathcal{O}\mathcal{T}}}$ and R is over, $\mathcal{A}^{\mathcal{R}_{\mathcal{O}\mathcal{T}}}$ continues the execution with R according to \mathcal{H}_0 (\mathcal{H}_1).

This part of the security proof ends with the observation that if $\mathcal{C}^{\mathcal{R}_{\mathcal{O}\mathcal{T}}}$ has used the simulator $\mathcal{S}_{\mathcal{R}_{\mathcal{O}\mathcal{T}}}$ then the joint distribution of the view of R and $\mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)$ corresponds to $\{\text{OUT}_{\mathcal{H}_0}^R(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\}$, to $\{\text{OUT}_{\mathcal{H}_1}^R(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\}$ otherwise.

The proof that

$$\{\text{OUT}_{\mathcal{H}_{i-1}}^R(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\} \approx \{\text{OUT}_{\mathcal{H}_i}^R(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\}$$

for $i = 2, \dots, \lambda$ follows the same arguments.

In order to prove that

$$\{\text{OUT}_{\mathcal{H}_\lambda}^{\text{R}}(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\} \approx \{\text{OUT}_{\mathcal{H}^*}^{\text{R}}(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\},$$

thus concluding the lemma's security proof, we need to consider the following intermediate hybrid experiment \mathcal{H}_y^* with $y \in \{1, \dots, \lambda\}$. The description of the hybrid experiment follows.

1. Compute $k_1 \leftarrow \text{Gen}(1^s)$ and run $\mathcal{S}_{\text{RO}\mathcal{T}}$ on input $(1^s, x_\lambda, k_1)$.
2. For $i = 2, \dots, y$ execute the following steps.
 - 2.1. Define the empty list c^i . For $j = 1, \dots, \min\{2^i, |Y|\} - 1$ compute and add $\text{Enc}(\text{Gen}(1^s), 0)$ to c^i .
 - 2.2. Compute $k_i \leftarrow \text{Gen}(1^s)$, and add $\text{Enc}(k_{i-1}, k_i)$ to the list c^i .
 - 2.3. Permute the elements inside c^i .
 - 2.4. Run $\mathcal{S}_{\text{RO}\mathcal{T}}$ on input $(1^s, x_{\lambda-i+1}, c^i)$.
3. For each $t \in \text{Prefix}(Y, y) - \{x_\lambda \dots x_{\lambda-y+1}\}$ compute $k_t \leftarrow \text{Gen}(1^s)$.
If $x_\lambda \dots x_{\lambda-y+1} \in \text{Prefix}(Y, y)$ then set $k_{x_\lambda \dots x_{\lambda-y+1}} = k_y$, otherwise $k_y^* = k_y$.
4. For $i = y + 1, \dots, \lambda$ execute the following steps.
 - 4.1. Define the empty list c^i and for each $t \in \text{Prefix}(Y, i - 1)$ execute the following steps.
 - If $t || x_{\lambda-i+1} \in \text{Prefix}(Y, i)$ then compute $k_{t || x_{\lambda-i+1}} \leftarrow \text{Gen}(1^s)$ and add $\text{Enc}(k_t, k_{t || x_{\lambda-i+1}})$ to the list c^i . Otherwise, if $t || x_{\lambda-i+1} \notin \text{Prefix}(Y, i)$ then compute and add $\text{Enc}(k_t, k_i^*)$ to the list c^i .
 - 4.2. If $|c^i| < \delta(i - 1)$ then execute the following steps.
 - Compute and add $\text{Enc}(k_{i-1}^*, k_i^*)$ to the list c^i .
 - For $i = 1, \dots, \delta(i - 1) - |c^i|$ compute and add $\text{Enc}(\text{Gen}(1^s), 0)$ to c^i .
 - 4.3. Permute the elements inside c^i .
 - 4.4. Run $\mathcal{S}_{\text{RO}\mathcal{T}}$ on input $(1^s, x_{\lambda-i+1}, c^i)$.
5. For every $t \in \text{Prefix}(Y, \lambda)$ compute and add $\text{Enc}(k_t, \gamma^1)$ to l .
6. If $|l| < 2^\lambda$ then compute and add $\text{Enc}(k_\lambda^*, \gamma^0)$ to l .
7. Permute the elements inside l and send l to R.

We now prove that

$$\{\text{OUT}_{\mathcal{H}_{y-1}^*}^{\text{R}}(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\} \approx \{\text{OUT}_{\mathcal{H}_y^*}^{\text{R}}(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\}$$

for $y = 2, \dots, \lambda$. The proof proceeds by contradiction. Suppose that there exists some $y \in \{2, \dots, \lambda\}$ such that

$$\{\text{OUT}_{\mathcal{H}_{y-1}^*}^{\text{R}}(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\} \approx \{\text{OUT}_{\mathcal{H}_y^*}^{\text{R}}(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\}$$

then we can construct an adversary \mathcal{A}^{Sym} that breaks the security of the encryption scheme Sym . Let \mathcal{C}^{Sym} be the challenger for the security game w.r.t to Sym . Our adversary runs R with randomness r and executes the following steps.

1. Compute $k_1 \leftarrow \text{Gen}(1^s)$ and run $\mathcal{S}_{\mathcal{R}_{\mathcal{OT}}}$ on input $(1^s, x_\lambda, k_1)$.
2. For $i = 2, \dots, y - 1$ execute the following steps.
 - 2.1. Define the empty list c^i . For $j = 1, \dots, \min\{2^i, |Y|\} - 1$ compute and add $\text{Enc}(\text{Gen}(1^s), 0)$ to c^i .
 - 2.2. Compute $k_i \leftarrow \text{Gen}(1^s)$, and add $\text{Enc}(k_{i-1}, k_i)$ to the list c^i .
 - 2.3. Permute the elements inside c^i .
 - 2.4. Run $\mathcal{S}_{\mathcal{R}_{\mathcal{OT}}}$ on input $(1^s, x_{\lambda-i+1}, c^i)$.
3. Define two empty lists \mathbf{m}_0 and \mathbf{m}_1 that will represent the challenge messages to be sent to \mathcal{C}^{Sym} .
4. For each $t \in \text{Prefix}(Y, y) - \{x_\lambda \dots x_{\lambda-y+1}\}$ compute $k_t \leftarrow \text{Gen}(1^s)$ and add it to the list \mathbf{m}_0 .
5. For $j = 1, \dots, |\text{Prefix}(Y, y) - \{x_\lambda \dots x_{\lambda-y+1}\}|$ compute and add 0 to \mathbf{m}_1 .
6. Send the challenge messages to \mathcal{C}^{Sym} .
7. Upon receiving the challenge ciphertext \mathbf{cx} , set $c^y = \mathbf{cx}$.
8. For $j = 1, \dots, \min\{2^y, |Y|\} - |\mathbf{cx}| - 1$ compute and add $\text{Enc}(\text{Gen}(1^s), 0)$ to c^y .
9. Compute $k_y \leftarrow \text{Gen}(1^s)$, and add $\text{Enc}(k_{y-1}, k_y)$ to the list c^y .
10. Permute the elements inside c^y .
11. Run $\mathcal{S}_{\mathcal{R}_{\mathcal{OT}}}$ on input $(1^s, x_{\lambda-i+1}, c^y)$.
12. If $x_\lambda \dots x_{\lambda-y+1} \in \text{Prefix}(Y, y)$ then set $k_{x_\lambda \dots x_{\lambda-y+1}} = k_y$, otherwise set $k_y^* = k_y$.
13. For $i = y + 1, \dots, \lambda$ execute the following steps.
 - 13.1. Define the empty list c^i and for each $t \in \text{Prefix}(Y, i - 1)$ execute the following steps.

If $t || x_{\lambda-i+1} \in \text{Prefix}(Y, i)$ then compute $k_{t || x_{\lambda-i+1}} \leftarrow \text{Gen}(1^s)$ and add $\text{Enc}(k_t, k_{t || x_{\lambda-i+1}})$ to the list c^i . Otherwise, if $t || x_{\lambda-i+1} \notin \text{Prefix}(Y, i)$ then compute and add $\text{Enc}(k_t, k_i^*)$ to the list c^i .
 - 13.2. If $|c^i| < \delta(i - 1)$ then execute the following steps.
 - Compute and add $\text{Enc}(k_{i-1}^*, k_i^*)$ to the list c^i .
 - For $i = 1, \dots, \delta(i - 1) - |c^i|$ compute and add $\text{Enc}(\text{Gen}(1^s), 0)$ to c^i .
 - 13.3. Permute the elements inside c^i .
 - 13.4. Run $\mathcal{S}_{\mathcal{R}_{\mathcal{OT}}}$ on input $(1^s, x_{\lambda-i+1}, c^i)$.
14. For every $t \in \text{Prefix}(Y, \lambda)$ compute and add $\text{Enc}(k_t, \gamma^1)$ to l .
15. If $|l| < 2^\lambda$ then compute and add $\text{Enc}(k_\lambda^*, \gamma^0)$ to l .
16. Permute the elements inside l and send l to \mathcal{R} .

This part of the security proof ends with the observation that if \mathcal{C}^{Sym} has used \mathbf{m}_0 then the joint distribution of the view of \mathbf{R} and $\mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)$ corresponds to $\{\text{OUT}_{\mathcal{H}_{y-1}^*}^{\mathbf{R}}(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\}$, to $\{\text{OUT}_{\mathcal{H}_y^*}^{\mathbf{R}}(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\}$ otherwise.

Since the following two distributions coincide

$$\{\text{OUT}_{\mathcal{H}_\lambda}^{\mathbf{R}}(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\} = \{\text{OUT}_{\mathcal{H}_\lambda^*}^{\mathbf{R}}(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\}$$

to complete the entire security proof we just need to prove that $\{\text{OUT}_{\mathcal{H}_\lambda^*}^{\mathbf{R}}(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\} \approx \{\text{OUT}_{\mathcal{H}_\lambda^*}^{\mathbf{R}}(1^s), \mathcal{F}^\in(Y, \gamma^0, \gamma^1, x)\}$. The indistinguishability between the two distributions can be proved by using arguments similar to the one used lately. That is, by proceedings by contradiction and constructing adversary that breaks the security of the encryption scheme Sym . \square

C Standard definitions

C.1 Computational indistinguishability definition

Definition 2 (Computational indistinguishability). *Let $X = \{X_s\}_{s \in \mathbb{N}}$ and $Y = \{Y_s\}_{s \in \mathbb{N}}$ be ensembles, where X_s and Y_s are probability distribution over $\{0, 1\}^l$, for some $l = \text{poly}(s)$. We say that $X = \{X_s\}_{s \in \mathbb{N}}$ and $Y = \{Y_s\}_{s \in \mathbb{N}}$ are computationally indistinguishable, denote $X \approx Y$, if for every PPT distinguisher \mathcal{D} there exists a negligible function ν such that for sufficiently large $s \in \mathbb{N}$,*

$$\left| \text{Prob}(t \leftarrow X_s : \mathcal{D}(1^s, t) = 1) - \text{Prob}(t \leftarrow Y_s : \mathcal{D}(1^s, t) = 1) \right| < \nu(s).$$

We note that in the usual case where $|X_s| = \Omega(s)$ and s can be derived from a sample of X_s , it is possible to omit the auxiliary input 1^s . In this paper we also use the definition of *Statistical Indistinguishability*. This definition is the same as Definition 2 with the only difference that the distinguisher \mathcal{D} is unbounded. In this case we use $X \equiv Y$ to denote that two ensembles are statistically indistinguishable.

C.2 Two party computation

A two-party protocol problem is cast by specifying a random process that maps pairs of inputs to pairs of outputs (one for each party). We refer to such a process as a functionality and denote it as $F = (F_1, F_2)$. That is, for every pair of inputs $x, y \in \{0, 1\}^s$, the output-pair is a random variable $(F_1(x, y), F_2(x, y))$ ranging over pairs of strings. The first party (with input x) wishes to obtain $F_1(x, y)$ and the second party (with input y) wishes to obtain $F_2(x, y)$. We often denote such a functionality by $(x, y) \rightarrow (F_1(x, y), F_2(x, y))$.

Let $F = (F_1, F_2)$ be a probabilistic polynomial-time functionality and let $\Pi = (P_1, P_2)$ be a two-party protocol for computing F where P_1 and P_2 denote the two parties. The view of the party P_i ($i \in \{1, 2\}$) during an execution of Π on (x, y) and security parameter s is denoted by $\text{view}_{P_i}^\Pi(x, y, 1^s)$.

The output of the party P_i ($i \in \{1, 2\}$) during an execution of Π on (x, y) and security parameter s is denoted by $\text{output}_{P_i}^\Pi(1^s, x, y)$ and can be computed from its own view of the execution. We denote the joint output of both parties by $\text{output}^\Pi(1^s, x, y) = (\text{output}_{P_1}^\Pi(1^s, x, y), \text{output}_{P_2}^\Pi(1^s, x, y))$.

Definition 3 (Secure two-party computation [HL10]). Let $F = (F_1, F_2)$ be a functionality. We say that Π securely computes F in the presence of static semi-honest adversaries if there exist probabilistic polynomial-time algorithms \mathcal{S}_{P_1} and \mathcal{S}_{P_2} called simulators, such that

$$\{(\mathcal{S}_{P_1}(1^s, x, F_1(x, y)), F(x, y))\}_{\{x, y, s\}} \approx \{\text{view}_{P_1}^\Pi(1^s, x, y), \text{output}^\Pi(1^s, x, y)\}_{\{x, y, s\}}$$

and

$$\{(\mathcal{S}_{P_2}(1^s, y, F_2(x, y)), F(x, y))\}_{\{x, y, s\}} \approx \{\text{view}_{P_2}^\Pi(1^s, x, y), \text{output}^\Pi(1^s, x, y)\}_{\{x, y, s\}}$$

where $x, y \in \{0, 1\}^*$ such that $|x| = |y|$, and $s \in \mathbb{N}$.

References

- [ALSZ13] Gilad Asharov, Yehuda Lindell, Thomas Schneider, and Michael Zohner. More efficient oblivious transfer and extensions for faster secure computation. In *2013 ACM SIGSAC Conference on Computer and Communications Security, CCS'13, Berlin, Germany, November 4-8, 2013*, pages 535–548, 2013.
- [ALSZ15] Gilad Asharov, Yehuda Lindell, Thomas Schneider, and Michael Zohner. More efficient oblivious transfer extensions with security for malicious adversaries. In *Advances in Cryptology - EUROCRYPT 2015 - 34th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Sofia, Bulgaria, April 26-30, 2015, Proceedings, Part I*, pages 673–701, 2015.
- [BBC⁺11] Pierre Baldi, Roberta Baronio, Emiliano De Cristofaro, Paolo Gasti, and Gene Tsudik. Countering GATTACA: efficient and secure testing of fully-sequenced human genomes. In *Proceedings of the 18th ACM Conference on Computer and Communications Security, CCS 2011, Chicago, Illinois, USA, October 17-21, 2011*, pages 691–702, 2011.
- [BMR90] Donald Beaver, Silvio Micali, and Phillip Rogaway. The round complexity of secure protocols (extended abstract). In *Proceedings of the 22nd Annual ACM Symposium on Theory of Computing, May 13-17, 1990, Baltimore, Maryland, USA*, pages 503–513, 1990.
- [CGN98] Benny Chor, Niv Gilboa, and Moni Naor. Private information retrieval by keywords. *IACR Cryptology ePrint Archive*, 1998:3, 1998. Appeared in the THEORY OF CRYPTOGRAPHY LIBRARY.
- [CLR17] Hao Chen, Kim Laine, and Peter Rindal. Fast private set intersection from homomorphic encryption. In Bhavani M. Thuraisingham, David Evans, Tal Malkin, and Dongyan Xu, editors, *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, CCS 2017, Dallas, TX, USA, October 30 - November 03, 2017*, pages 1243–1255. ACM, 2017.
- [CT10] Emiliano De Cristofaro and Gene Tsudik. Practical private set intersection protocols with linear complexity. In *Financial Cryptography and Data Security, 14th International Conference, FC 2010, Tenerife, Canary Islands, January 25-28, 2010, Revised Selected Papers*, pages 143–159, 2010.
- [CT12] Emiliano De Cristofaro and Gene Tsudik. Experimenting with fast private set intersection. In *Trust and Trustworthy Computing - 5th International Conference, TRUST 2012, Vienna, Austria, June 13-15, 2012. Proceedings*, pages 55–73, 2012.
- [DNT12] Morten Dahl, Chao Ning, and Tomas Toft. On secure two-party integer division. In *Financial Cryptography and Data Security - 16th International Conference, FC 2012, Kralendijk, Bonaire, February 27-March 2, 2012, Revised Selected Papers*, pages 164–178, 2012.
- [FIPR05] Michael J. Freedman, Yuval Ishai, Benny Pinkas, and Omer Reingold. Keyword search and oblivious pseudorandom functions. In *Theory of Cryptography, Second Theory of*

Cryptography Conference, TCC 2005, Cambridge, MA, USA, February 10-12, 2005, Proceedings, pages 303–324, 2005.

- [FNO18] Brett Hemenway Falk, Daniel Noble, and Rafail Ostrovsky. Private set intersection with linear communication from general assumptions. *IACR Cryptology ePrint Archive*, 2018:238, 2018.
- [FNP04] Michael J. Freedman, Kobbi Nissim, and Benny Pinkas. Efficient private matching and set intersection. In *Advances in Cryptology - EUROCRYPT 2004, International Conference on the Theory and Applications of Cryptographic Techniques, Interlaken, Switzerland, May 2-6, 2004, Proceedings*, pages 1–19, 2004.
- [GMW87] Oded Goldreich, Silvio Micali, and Avi Wigderson. How to play any mental game or A completeness theorem for protocols with honest majority. In *Proceedings of the 19th Annual ACM Symposium on Theory of Computing, 1987, New York, New York, USA*, pages 218–229, 1987.
- [GSV07] Juan A. Garay, Berry Schoenmakers, and José Villegas. Practical and secure solutions for integer comparison. In *Public Key Cryptography - PKC 2007, 10th International Conference on Practice and Theory in Public-Key Cryptography, Beijing, China, April 16-20, 2007, Proceedings*, pages 330–342, 2007.
- [HEK12] Yan Huang, David Evans, and Jonathan Katz. Private set intersection: Are garbled circuits better than custom protocols? In *19th Annual Network and Distributed System Security Symposium, NDSS 2012, San Diego, California, USA, February 5-8, 2012*, 2012.
- [HL08] Carmit Hazay and Yehuda Lindell. Efficient protocols for set intersection and pattern matching with security against malicious and covert adversaries. In *Theory of Cryptography, Fifth Theory of Cryptography Conference, TCC 2008, New York, USA, March 19-21, 2008.*, pages 155–175, 2008.
- [HL10] Carmit Hazay and Yehuda Lindell. *Efficient Secure Two-Party Protocols - Techniques and Constructions*. Information Security and Cryptography. Springer, 2010.
- [HOS17] Per Hallgren, Claudio Orlandi, and Andrei Sabelfeld. Privatepool: Privacy-preserving ridesharing. In *IEEE 30th Computer Security Foundations Symposium, CSF 2017, Santa Barbara, CA, USA, August 21-25, 2017*, pages 276–291, 2017.
- [HV17] Carmit Hazay and Muthuramakrishnan Venkitasubramaniam. Scalable multi-party private set-intersection. In *Public-Key Cryptography - PKC 2017 - 20th IACR International Conference on Practice and Theory in Public-Key Cryptography, Amsterdam, The Netherlands, March 28-31, 2017, Proceedings, Part I*, pages 175–203, 2017.
- [IKN⁺17] Mihaela Ion, Ben Kreuter, Erhan Nergiz, Sarvar Patel, Shobhit Saxena, Karn Seth, David Shanahan, and Moti Yung. Private intersection-sum protocol with applications to attributing aggregate ad conversions. *Cryptology ePrint Archive*, Report 2017/738, 2017. <http://eprint.iacr.org/2017/738>.

- [IKNP03] Yuval Ishai, Joe Kilian, Kobbi Nissim, and Erez Petrank. Extending oblivious transfers efficiently. In *Advances in Cryptology - CRYPTO 2003, 23rd Annual International Cryptology Conference, Santa Barbara, California, USA, August 17-21, 2003, Proceedings*, pages 145–161, 2003.
- [IP07] Yuval Ishai and Anat Paskin. Evaluating branching programs on encrypted data. In *Theory of Cryptography, 4th Theory of Cryptography Conference, TCC 2007, Amsterdam, The Netherlands, February 21-24, 2007, Proceedings*, pages 575–594, 2007.
- [JL10] Stanislaw Jarecki and Xiaomin Liu. Fast secure computation of set intersection. In *Security and Cryptography for Networks, 7th International Conference, SCN 2010, Amalfi, Italy, September 13-15, 2010. Proceedings*, pages 418–435, 2010.
- [KK13] Vladimir Kolesnikov and Ranjit Kumaresan. Improved OT extension for transferring short secrets. In *Advances in Cryptology - CRYPTO 2013 - 33rd Annual Cryptology Conference, Santa Barbara, CA, USA, August 18-22, 2013. Proceedings, Part II*, pages 54–70, 2013.
- [KKRT16] Vladimir Kolesnikov, Ranjit Kumaresan, Mike Rosulek, and Ni Trieu. Efficient batched oblivious PRF with applications to private set intersection. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, Vienna, Austria, October 24-28, 2016*, pages 818–829, 2016.
- [KMRS14] Seny Kamara, Payman Mohassel, Mariana Raykova, and Seyed Saeed Sadeghian. Scaling private set intersection to billion-element sets. In *Financial Cryptography and Data Security - 18th International Conference, FC 2014, Christ Church, Barbados, March 3-7, 2014, Revised Selected Papers*, pages 195–215, 2014.
- [KOS15] Marcel Keller, Emanuela Orsini, and Peter Scholl. Actively secure OT extension with optimal overhead. In *Advances in Cryptology - CRYPTO 2015 - 35th Annual Cryptology Conference, Santa Barbara, CA, USA, August 16-20, 2015, Proceedings, Part I*, pages 724–741, 2015.
- [LP09] Yehuda Lindell and Benny Pinkas. A proof of security of yao’s protocol for two-party computation. *J. Cryptology*, 22(2):161–188, 2009.
- [LT13] Helger Lipmaa and Tomas Toft. Secure equality and greater-than tests with sublinear online complexity. In *Automata, Languages, and Programming - 40th International Colloquium, ICALP 2013, Riga, Latvia, July 8-12, 2013, Proceedings, Part II*, pages 645–656, 2013.
- [Mea86] Catherine A. Meadows. A more efficient cryptographic matchmaking protocol for use in the absence of a continuously available third party. In *Proceedings of the 1986 IEEE Symposium on Security and Privacy, Oakland, California, USA, April 7-9, 1986*, pages 134–137, 1986.
- [MN12] Payman Mohassel and Salman Niksefat. Oblivious decision programs from oblivious transfer: Efficient reductions. In *Financial Cryptography and Data Security - 16th International Conference, FC 2012, Kralendijk, Bonaire, February 27-March 2, 2012, Revised Selected Papers*, pages 269–284, 2012.

- [NMH⁺10] Shishir Nagaraja, Prateek Mittal, Chi-Yao Hong, Matthew Caesar, and Nikita Borisov. Botgrep: Finding P2P bots with structured graph analysis. In *19th USENIX Security Symposium, Washington, DC, USA, August 11-13, 2010, Proceedings*, pages 95–110, 2010.
- [OOS17] Michele Orrù, Emanuela Orsini, and Peter Scholl. Actively secure 1-out-of-n OT extension with application to private set intersection. In *Topics in Cryptology - CT-RSA 2017 - The Cryptographers' Track at the RSA Conference 2017, San Francisco, CA, USA, February 14-17, 2017, Proceedings*, pages 381–396, 2017.
- [PSSZ15] Benny Pinkas, Thomas Schneider, Gil Segev, and Michael Zohner. Phasing: Private set intersection using permutation-based hashing. In *24th USENIX Security Symposium, USENIX Security 15, Washington, D.C., USA, August 12-14, 2015.*, pages 515–530, 2015.
- [PSWW18] Benny Pinkas, Thomas Schneider, Christian Weinert, and Udi Wieder. Efficient circuit-based PSI via cuckoo hashing. In Jesper Buus Nielsen and Vincent Rijmen, editors, *Advances in Cryptology - EUROCRYPT 2018 - 37th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Tel Aviv, Israel, April 29 - May 3, 2018 Proceedings, Part III*, volume 10822 of *Lecture Notes in Computer Science*, pages 125–157. Springer, 2018.
- [PSZ14] Benny Pinkas, Thomas Schneider, and Michael Zohner. Faster private set intersection based on OT extension. In *Proceedings of the 23rd USENIX Security Symposium, San Diego, CA, USA, August 20-22, 2014.*, pages 797–812, 2014.
- [PSZ16] Benny Pinkas, Thomas Schneider, and Michael Zohner. Scalable private set intersection based on ot extension. Cryptology ePrint Archive, Report 2016/930, 2016. <http://eprint.iacr.org/2016/930>.
- [RR17a] Peter Rindal and Mike Rosulek. Improved private set intersection against malicious adversaries. In *Advances in Cryptology - EUROCRYPT 2017 - 36th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Paris, France, April 30 - May 4, 2017, Proceedings, Part I*, pages 235–259, 2017.
- [RR17b] Peter Rindal and Mike Rosulek. Malicious-secure private set intersection via dual execution. In Bhavani M. Thuraisingham, David Evans, Tal Malkin, and Dongyan Xu, editors, *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, CCS 2017, Dallas, TX, USA, October 30 - November 03, 2017*, pages 1229–1242. ACM, 2017.
- [Sha80] Adi Shamir. On the power of commutativity in cryptography. In *Automata, Languages and Programming, 7th Colloquium, Noordwijkerhout, The Netherland, July 14-18, 1980, Proceedings*, pages 582–595, 1980.
- [T⁺07] Tomas Toft et al. Primitives and applications for multi-party computation. *PhD Thesis, University of Aarhus, Denmark*, 2007.

- [TLP⁺17] Sandeep Tamrakar, Jian Liu, Andrew Paverd, Jan-Erik Ekberg, Benny Pinkas, and N. Asokan. The circle game: Scalable private membership test using trusted hardware. In *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, AsiaCCS 2017, Abu Dhabi, United Arab Emirates, April 2-6, 2017*, pages 31–44, 2017.
- [Yao82] Andrew Chi-Chih Yao. Protocols for secure computations (extended abstract). In *23rd Annual Symposium on Foundations of Computer Science, Chicago, Illinois, USA, 3-5 November 1982*, pages 160–164, 1982.