

# Outsourced Pattern Matching

Sebastian Faust\*

Carmit Hazay†

Daniele Venturi‡

April 30, 2015

## Abstract

In secure delegatable computation, computationally weak devices (or clients) wish to outsource their computation and data to an *untrusted* server in the cloud. While most earlier work considers the general question of how to securely outsource *any* computation to the cloud server, we focus on *concrete* and *important* functionalities and give the first protocol for the *pattern matching* problem in the cloud. Loosely speaking, this problem considers a text  $T$  that is outsourced to the cloud  $S$  by a sender  $SEN$ . In a query phase, receivers  $REC_1, \dots, REC_l$  run an efficient protocol with the server  $S$  and the sender  $SEN$  in order to learn the positions at which a pattern of length  $m$  matches the text (and nothing beyond that). This is called the *outsourced pattern matching* problem which is highly motivated in the context of delegatable computing since it offers storage alternatives for massive databases that contain confidential data (e.g., health related data about patient history).

Our constructions are *simulation-based secure* in the presence of semi-honest and malicious adversaries (in the random oracle model) and limit the communication in the query phase to  $O(m)$  bits plus the number of occurrences—which is optimal. In contrast to generic solutions for delegatable computation, our schemes do not rely on fully homomorphic encryption but instead use novel ideas for solving pattern matching, based on a reduction to the subset sum problem. Interestingly, we do not rely on the hardness of the problem, but rather we exploit instances that are solvable in polynomial-time. A follow-up result demonstrates that the random oracle is essential in order to meet our communication bound.

---

\*Ruhr-Universität Bochum. Email: [sebastian.faust@gmail.com](mailto:sebastian.faust@gmail.com).

†Faculty of Engineering, Bar-Ilan University, Israel. Email: [carmit.hazay@biu.ac.il](mailto:carmit.hazay@biu.ac.il).

‡Computer Science Department, Sapienza University of Rome, Italy. Email: [daniele.venturi@uniroma1.it](mailto:daniele.venturi@uniroma1.it). Acknowledges funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 644666.

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>	2.3 Oblivious Pseudorandom Function Evaluation . . . . .	9
1.1	Our Contribution . . . . .	4	2.4 Commitment Schemes . . . . .	10
1.1.1	Modeling Outsourced Pattern Matching . . . . .	4	2.5 Collision Resistant Hashing and Merkle Trees . . . . .	10
1.1.2	Semi-Honest Outsourced Pattern Matching from Subset Sum . . . . .	5	2.6 Universal Accumulators . . . . .	11
1.1.3	Malicious Outsourced Pattern Matching . . . . .	6	<b>3 Modeling Outsourced Pattern Matching</b>	<b>12</b>
1.2	Related Work . . . . .	7	<b>4 Security in the Presence of Semi-Honest Adversaries</b>	<b>14</b>
1.3	Limitations and Open Problems . . . . .	7	4.1 Efficiency . . . . .	20
<b>2</b>	<b>Preliminaries</b>	<b>8</b>	<b>5 Security in the Presence of Malicious Adversaries</b>	<b>21</b>
2.1	Basic Notations . . . . .	8	5.1 Dealing with a Malicious Server . . . . .	21
2.2	The Subset Sum Problem . . . . .	8	5.2 Dealing with a Malicious Receiver . . . . .	26

## 1 Introduction

**Background on outsourced secure computation.** The problem of securely outsourcing computation to an *untrusted* server gained momentum with the recent penetration of *cloud computing* services. In cloud computing, clients can lease computing services on demand rather than maintaining their own infrastructure. While such an approach naturally has numerous advantages in cost and functionality, it is crucial that the outsourcing mechanism enforces privacy of the outsourced data and integrity of the computation. Solutions based on cryptographic techniques have been put forward with the concept of *secure delegatable computation* [AIK10, CKV10, GGP10, AJLA<sup>+</sup>12, FGP14], which lately has received broad attention within the cryptographic research community.

In secure delegatable computation, computationally weak devices (or clients) wish to outsource their computation and data to an *untrusted* server in the cloud. The ultimate goal in this setting is to design efficient protocols that minimize the computational overhead of the clients and instead rely on the extended resources of the server. Of course, the amount of work invested by the client in order to verify the correctness of the computation shall be *substantially* smaller than running the computation by itself. Another ambitious challenge of delegatable computation is to design protocols that minimize the *communication* between the cloud and the client, while using an optimal number of rounds. This becomes of particular importance with the proliferation of Smartphone technology and mobile broadband internet connections, as for mobile devices communication and data connectivity is often the more severe bottleneck.

Most recent works in the area of delegatable computation propose solutions to securely outsource *any functionality* to an untrusted server [AIK10, CKV10, GGP10, FGP14]. While this of course is the holy grail in delegatable computation, such generic solutions often suffer from rather poor efficiency and high computation overhead due to the use of fully homomorphic encryption [Gen09]. Furthermore, these solution concepts typically examine a restricted scenario where a *single client* outsources its computation to an *external untrusted server*. Another line of works studies an extended setting with *multiple clients* that mutually distrust each other and wish to securely outsource a joint computation on their inputs with reduced costs [KMR11, KMR12, LATV12, AJLA<sup>+</sup>12, CKKC13, GKL<sup>+</sup>15]. Of course, also in this more complex

setting with multiple clients recent constructions build up on fully homomorphic encryption (or consider a more restricted setting where the clients do not collude with the server, or one of the clients “works harder”).

To move towards more practical schemes, we may give up on outsourcing arbitrary computation to the cloud, but instead focus on particularly efficient constructions for specific important functionalities. This approach has the potential to avoid the use of fully homomorphic encryption (FHE) by exploiting the structure of the particular problem we intend to solve. Some recent works have considered this question and proposed schemes for polynomial evaluation and keywords search [BGV11], set operations [PTT11] and linear algebra [Moh11]. While these schemes are more efficient than the generic constructions mentioned above, they typically only achieve very limited privacy or do not support multiple distrusting clients.

In this paper, we follow this line of work and provide the first protocols for *pattern matching* in the cloud. The problem of *outsourced pattern matching* is highly motivated in the context of delegatable computing since it offers storage alternatives for massive databases that contain confidential data (e.g., health related data about patient history). In contrast to most earlier works, our constructions achieve a high level of security, while avoiding the use of FHE and minimizing the amount of communication between the parties. We emphasize that even with the power of FHE it is not clear how to get down to optimal communication complexity in two rounds.<sup>1</sup>

**Pattern matching in the cloud.** The problem of pattern matching considers a text  $T$  of length  $n$  and a pattern of length  $m$  with the goal to find all the locations where the pattern matches the text. Algorithms for pattern matching have been widely studied for decades due to its broad applicability [KMP77, BM77]. Lately, researchers started to look at this problem in the context of secure two-party computation [TPKC07, HL10, GHS10, KM10, HT10] due to growing interests in private text search. In this setting, one party holds the text whereas the other party holds the pattern and attempts to learn all the locations of the pattern in the text (and only that), while the party holding the text learns nothing about the pattern. Unfortunately, these solutions are not directly applicable in the cloud setting, mostly because the communication overhead per search query grows linearly with the text length. Moreover, the text holder delegates its work to an external untrusted server and cannot control the content of the server’s responses.

To be precise, in the outsourced setting we consider a set of clients comprised from a sender  $SEN$  and a set of receivers  $(REC_1, \dots, REC_l)$  that interact with a server  $S$  in the following way. In a *setup phase* the sender  $SEN$  uploads a preprocessed text to an external server  $S$ . This phase is run only once and may be costly in terms of computation and communication. In a *query phase* the receivers  $REC_1, \dots, REC_l$  query the text by searching patterns and learn the matched text locations. The main two goals are as follows:

1. *Simulation-based security:* We model outsourced pattern matching by a strong simulation-based security definition (cf. Section 3) that captures all security concerns. Namely, we define a new reactive outsourced pattern matching functionality  $\mathcal{F}_{OPM}$  that ensures the secrecy and integrity of the outsourced text and patterns. For instance, a semi-honest server does not gain any information about the text and patterns, except of what it can infer from the answers to the search queries (this is formalized more accurately below). If the server is maliciously corrupted the functionality implies the correctness of the queries’ replies as well. As in the standard secure computation setting, simulation-based modeling is simpler and stronger than game-based definitions, yet harder to achieve. As the simulator must commit first to the data stored on the server, and yet, be able to simulate the server’s view as in the real interaction. We further elaborate on these difficulties below.
2. *Sublinear communication complexity during query phase:* We consider an *amortized* model, where the communication and computational costs of the clients reduce with the number of queries. More

---

<sup>1</sup>Namely, a two-rounds solution based on FHE would need a circuit that tolerates the worst case output size, which in pattern matching implies a maximal number of matches that is proportional to the length of the text (or the database).

concretely, while in the setup phase communication and computation is linear in the length of the text, we want that during the query phase the overall communication and the work put by the receivers is *linear in the number of matches* (which is optimal). Of course, we also require the server running in polynomial time. Notice that our strong efficiency requirement comes at a price: it allows the (possibly corrupted) server to learn the number of matches at the very least, as otherwise there is no way to achieve overall communication complexity that is linear in the number of matches. We model this by giving the server *some leakage* for each pattern query which will be described in detail below.

## 1.1 Our Contribution

In the following we will always talk about a single receiver REC that interacts with a sender SEN and a server S in the query phase. This is only to simplify notation. Our protocols can naturally be applied to a setting with mutually distrusting receivers.

### 1.1.1 Modeling Outsourced Pattern Matching

We follow the standard method for showing security of protocols using the ideal/real world paradigm [Can00]. We give a specification of an ideal execution with a trusted party by defining a reactive outsourced pattern matching functionality  $\mathcal{F}_{\text{OPM}}$ . This functionality works in two phases: In the preprocessing phase sender SEN uploads its preprocessed text  $\tilde{T}$  to the server. Next, in an iterative query phase, upon receiving a search query  $p$  the functionality asks for the approvals of sender SEN (as it may also refuse for this query in the real execution), and the server (as in case of being corrupted it may abort the execution). To model the additional leakage that is required to minimize communication we ask the functionality to forward the server the matched positions in the text upon receiving an approval from SEN. Note that this type of leakage is necessary if S and REC might collude, as in our case. We further note that our functionality returns all matched positions but can be modified so that only the first few matched positions are returned.<sup>2</sup>

**Difficulties with simulating  $\mathcal{F}_{\text{OPM}}$ .** The main challenge in designing a simulator for this functionality is in case when the server is corrupted. In this case the simulator must *commit* to some text in a way that allows it later to produce a sequence of trapdoors that is consistent with the sequence of queries. More precisely, when the simulator commits to a preprocessed text, the leakage that the corrupted server obtains (namely, the positions where the pattern matches the text) has to be consistent with the information that it sees during the query phase. This implies that the simulator must have flexibility when it later matches the committed text to the trapdoors. Due to this inherent difficulty the text must be *encoded* in a way, that given a search query  $p$  and a list of text positions  $(i_1, \dots, i_t)$ , one can produce a trapdoor for  $p$  in such a way that the “search” in the preprocessed text, using this trapdoor, yields  $(i_1, \dots, i_t)$ . We note that alternative (and even simpler) solutions that permute the preprocessed text in order to prevent the server from even learning the matched positions, necessarily require that the server does not collude with the receivers, and even then are not simple to implement. This is because the simulator must have some equivocation mechanism since it cannot distinguish in the preprocessing phase between the last element in a list of matched positions and a non-last element. In contrast, our solutions allow such strong collusion between the server and the receivers. We stress that it is highly non-trivial to achieve collusion for round efficient protocols and that it is impossible when considering general functionalities [GKL<sup>+</sup>15].

---

<sup>2</sup>This definition is more applicable for search engines where the first few results are typically more relevant, whereas the former variant is more applicable for a DNA search where it is important to find all matched positions. For simplicity we only consider the first variant, our solutions support both variants.

We wish to stress that the related problem of keyword search, where the goal is to retrieve a record associated with some keyword given that the keyword appears in the database, is a much simpler problem and easily solvable using pseudorandom functions (PRF). Specifically, keyword search has a fixed bound on the query response and can be solved using oblivious PRF, whereas in pattern matching this bound may be  $O(n)$  and this problem is much harder to solve (even without privacy).

### 1.1.2 Semi-Honest Outsourced Pattern Matching from Subset Sum

Our first construction for outsourced pattern matching is secure against semi-honest adversaries. At the heart of our constructions lies a novel encoding technique that is based on the *subset sum problem*. In this construction a sender SEN generates a vector of random values, conditioned on that the sum of elements in all positions that match the pattern equals some specified value that will be explained below. Namely, SEN builds an instance  $\tilde{\mathbf{T}}$  for the subset sum problem, where given a trapdoor  $R$  the goal is to find whether there exists a subset in  $\tilde{\mathbf{T}}$  that sums to  $R$ . This is in contrast to all prior cryptographic constructions, e.g., [LPS10] that design cryptographic schemes based on the hardness of this problem.

More formally, the subset sum problem is parameterized by two integers  $\ell$  and  $M$ . An instance of the problem is generated by picking random vectors  $\tilde{\mathbf{T}} \leftarrow \mathbb{Z}_M^\ell$ ,  $\mathbf{s} \leftarrow \{0, 1\}^\ell$  and outputting  $(\tilde{\mathbf{T}}, R = \tilde{\mathbf{T}}^\top \cdot \mathbf{s} \bmod M)$ . The problem is to find  $\mathbf{s}$  given  $\tilde{\mathbf{T}}$  and a trapdoor  $R$ . Looking ahead, we will have such a trapdoor  $R_p$  for each pattern  $p$  of length  $m$ , such that if  $p$  matches  $T$  then with overwhelming probability there will be a unique solution to the subset sum instance  $(\tilde{\mathbf{T}}, R_p)$ . This unique solution is placed at exactly the positions where the pattern appears in the text. The receiver REC that wishes to search for a pattern  $p$  obtains this trapdoor from SEN and will hand it to the server. Consequently, we are interested in easy instances of the subset sum problem since we require the server to solve it for each query (we note that using our packaging technique described below, brute force may be applicable as well). We therefore consider low-density instances which can be solved in polynomial time by a reduction to a short vector in a lattice [LO85, Fri86, CJL<sup>+</sup>92] (see Section 2.2 for more details). See Figure 1 for an illustration of our idea.

We further note that the security of the scheme relies heavily on the unpredictability of the trapdoor. Namely, in order to ensure that the server cannot guess the trapdoor for some pattern  $p$  (and thus solve the subset problem and find the matched locations), we require that the trapdoor is unpredictable. We therefore employ a PRF  $F$  on the pattern and fix this value as the trapdoor, where the key  $k$  for the PRF is picked by SEN and the two clients SEN and REC communicate via a secure two-party protocol to compute the evaluation of the PRF. This ensures hardness of guessing the trapdoor for any pattern  $p$ .

**Efficiency.** The scheme described above satisfies the appealing properties that the communication complexity during the setup phase is  $O(\kappa \cdot n)$  and during the query phase is proportional to the number of matches times  $\kappa$ . Furthermore, the space complexity at the server side is  $O(\kappa \cdot n)$  and the construction is round optimal, where the number of messages exchanged by the parties is minimal. A moment of reflection, however, shows that the scheme has very limited usage in practice. Recall that the server is asked to solve subset sum instances of the form  $(\tilde{\mathbf{T}}, R_p)$ , where  $\tilde{\mathbf{T}}$  is a vector of length  $\ell = n - m + 1$  with elements from  $\mathbb{Z}_M$  for some integer  $M$ . In order to ensure correctness we must guarantee that given a subset sum instance, each trapdoor has a unique solution with high probability. In other words, the collision probability, which equals  $2^\ell/M$  (stated also in [IN96]), should be negligible. Fixing  $M = 2^{\kappa+n}$  for a security parameter  $\kappa$ , ensures this for large enough  $\kappa$ , say whenever  $\kappa \geq 80$ . On the other hand, we need the subset sum problem to be solvable in polynomial time. A simple calculation (see analysis in Section 2.2), yields in this case a value of  $\ell \approx \sqrt{\kappa}$ . This poses an inherent limitation on the length of the text to be preprocessed.

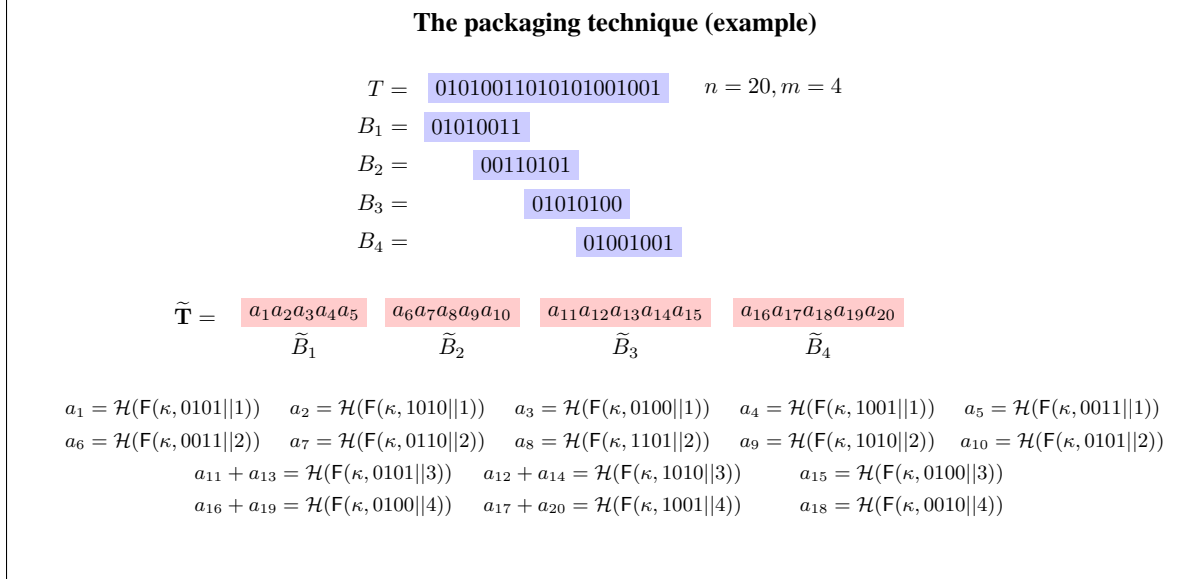


Figure 1: The packaging technique applied to a text of length  $n = 20$  bits, encoded to search for patterns of length  $m = 4$  bits.

**An improved solution using packaging.** To overcome this limitation, we employ an important extension of our construction based on *packaging*. First, the text is partitioned into smaller pieces of length  $2m$  which are handled separately by the protocol, where  $m$  is as usual the pattern length. Moreover, every two consecutive blocks are overlapping in  $m$  positions, so that we do not miss any match in the original text. Even though this approach introduces some overhead, yielding a text  $T'$  of overall length  $2n/m$ , note that now Eq. (2) below yields  $\ell = 2m - m + 1 = m + 1 < \sqrt{\kappa}$ , which is an upper bound on the length of the pattern (and not on the length of the text as before). Namely, we remove the limitation on the text length and consider much shorter block lengths for the subset sum algorithm.

This comes at a price since we now need to avoid using in each block the same trapdoor for some pattern  $p$ , as repetitions allow the server to extract potential valid trapdoors (that have not been queried yet) and figure out information about the text. This is in particular a problem when  $m$  is reasonably short as in this case since the server may just try out all potential trapdoors. One may suggest to generate for each block a completely independently chosen trapdoor. Unfortunately, this does not work as during the query phase the receiver needs to communicate these trapdoors to the server which may require communication linear in the length of the text. We solve this problem by requiring from the function outputting the trapdoors to have some form of “programmability” (which allows to simulate the answers to all queries consistently).

Specifically, we implement this function using the random oracle methodology on top of the PRF, so that a trapdoor now is computed by  $\mathcal{H}(F(k, p)||b)$ , for  $b$  being the block number. Now, the simulator can program the oracle to match with the positions where the pattern appears in each block. Note that using just the random oracle (without the PRF) is not sufficient as well, since an adversary that controls the server and has access to the random oracle can apply it on  $p$  as well.

### 1.1.3 Malicious Outsourced Pattern Matching

We extend our construction to the malicious setting as well, tolerating malicious attacks in the presence of (even colluding) corrupted server and receiver. Our proof ensures that the server returns the correct answers by employing Merkle commitments and accumulators. Informally speaking, Merkle commitments are suc-

cinct commitment schemes for which the commitment size is independent of the length of the committed value (or set). This tool is very useful in ensuring correctness, since now, upon committing to  $\tilde{T}$ , the server decommits the solution to the subset sum trapdoor and the receiver REC can simply verify that the decommitted values correspond to the trapdoor. Nevertheless, this solution does not cover the case of a mismatch since a corrupted server can always return a “no-match” message. In order to avoid it we use accumulators for proving whether an element is in a specified set or not. Specifically, SEN commits to the set of trapdoors for all patterns  $p$  that match  $T$  in at least one position. We then ask the server to prove membership/non-membership relative to this set, which contains at most  $n - m + 1$  elements. This implies that the server cannot declare “no-match” without providing a non-membership proof.

We further note that proving security against a corrupted REC is a straightforward extension of the semi-honest proof using the modifications we made above and the fact that the protocol for implementing the oblivious PRF evaluation is secure against malicious adversaries as well. Combining these proofs in order to prove security in the presence of colluding server and receiver follows easily.

**Efficiency.** Our protocols incur slightly higher overhead for the malicious setting due to the use of Merkle commitments. More specifically, correctness against a corrupted server increases the communication complexity by a factor of  $O(\log n)$  in the query phase. This is because each matched position must be verified against its commitment which costs  $O(\log n)$  using Merkle commitments. On the other hand, we note that our constructions are round optimal and maintain a minimal number of rounds (which is highly non-trivial in the malicious setting). This assumes a two-rounds implementation of oblivious PRF evaluation in the presence of a malicious receiver such as the protocol shown in [HL10].

## 1.2 Related Work

A related line of works regarding symmetric searchable encryption (SSE) [CGKO11, KPR12, KP13, JJK<sup>+</sup>13] allows a party to privately outsource its data to another party while maintaining the ability to search for keywords, where the main goal is typically to minimize the search time. We note that SSE that supports adaptive queries, simulation-based security and query privacy can be applied in a setting with a single sender and multiple receivers, by letting the receivers learn their trapdoors using a secure two-party protocol that is engaged with the sender. As recently demonstrated in [HZ14], there exists a large class of search problems (which includes pattern matching and all its variants), where there exists no SSE with adaptive queries, minimal interaction and communication that is proportional to the query’s response. This lower bound applies to both non-private and private channels scenarios (where in the private channels setting corrupted parties do not see the communication between the honest parties). This implies that our semi-honest construction is tight with respect to the assumptions it requires. Achieving SSE robust against malicious servers requires using additional techniques for ensuring correctness.

In [CS14], Chase and Shen solve outsourced pattern matching by constructing a queryable encryption scheme, which is an encryption scheme that supports search queries between a client and a server with three rounds of communication. Their construction is based on suffix trees (a data structure that solves pattern matching and related problems on unencrypted data). We note that in the setting of [CS14] the server is allowed to learn some relations between the queries, due to the possibility to observe the pattern relative to nodes in the tree that were already visited. In particular, this solution leaks much more information than what we allow in this work.

## 1.3 Limitations and Open Problems

A limitation of our constructions is that they only work for fixed length patterns, as the pre-processing inherently allows to search for patterns of length exactly  $m$ . We stress that most pattern matching algorithms

(in the non-private setting) work under a similar restriction, e.g., [KMP77], where the preprocessing phase depends heavily on the value of  $m$ . The naive way of removing the above limitation is to repeat the preprocessing many times, one for each individual pattern length we want to support. An alternative approach is to use suffix trees as suggested in [CS14]; however this solution does not comply with our ideal functionality for outsourced pattern matching due to the leakage which occurs by performing a search in the tree. We leave it as an open problem to design more efficient solutions for outsourced pattern matching, allowing to handle variable-length patterns without incurring in large leakage.

Our solutions are also limited in the range of values supported for the pattern length  $m$ . While algorithms for low-density subset sum instances would work well even for large values of  $m$  (say  $m \approx 1000$  bits), we additionally need the solution of each subset sum instance to be unique with high probability which yield much smaller values of  $m \approx \sqrt{\kappa}$ . Extending the range of possible values for  $m$  is an interesting direction for future research.

## 2 Preliminaries

### 2.1 Basic Notations

We let  $\mathbb{N}$  be the natural numbers and denote with  $\kappa$  the security parameter. Unless described otherwise, all quantities are implicitly functions of the security parameter, which is usually represented in unary and given as input to all cryptographic algorithms (including the adversary). We let  $\text{poly}(\kappa)$  denote an unspecified function  $O(\kappa^c)$  for some constant  $c$ . A function  $\text{negl}(\kappa)$  is *negligible* (in  $\kappa$ ) if it is  $\kappa^{-\omega(1)}$ . Given a string  $a \in \{0, 1\}^t$  we specify its value in the  $i$ th position by  $a[i]$ ; if  $\mathbf{a}$  is a vector  $\mathbf{a}[i]$  denotes the  $i$ th element of  $\mathbf{a}$ . We write PPT for probabilistic polynomial-time algorithms, i.e., randomized algorithms running in time  $\text{poly}(\kappa)$ . Let  $X = \{X_\kappa\}_{\kappa \in \mathbb{N}}$  and  $Y = \{Y_\kappa\}_{\kappa \in \mathbb{N}}$  be distribution ensembles. We say that  $X$  and  $Y$  are computationally indistinguishable, written  $X \stackrel{c}{\approx} Y$ , if for any PPT algorithm  $D$  we have

$$|\Pr[D(X_\kappa) = 1] - \Pr[D(Y_\kappa) = 1]| \leq \text{negl}(\kappa),$$

where the probability is taken over the random values  $X_\kappa$  and  $Y_\kappa$ , and the randomness of  $D$ .

### 2.2 The Subset Sum Problem

The subset sum problem is parametrized by two integers  $\ell$  and  $M$ . An instance of the problem is generated by picking random vectors  $\mathbf{a} \leftarrow \mathbb{Z}_M^\ell$ ,  $\mathbf{s} \leftarrow \{0, 1\}^\ell$  and outputting  $(\mathbf{a}, R = \mathbf{a}^\top \cdot \mathbf{s} \bmod M)$ . The problem is to find  $\mathbf{s}$  given  $\mathbf{a}$  and  $R$ . We rely on the following simple observation, which already appeared in [IN96].

**Lemma 1 ([IN96])** *Fix  $\ell$  and  $M$ . Let  $\mathbf{a}$  and  $\mathbf{s}$  be chosen uniformly at random and  $R = \mathbf{a}^\top \cdot \mathbf{s} \bmod M$ . Then, the probability that there exists a vector  $\mathbf{s}' \neq \mathbf{s}$  such that  $R = \mathbf{a}^\top \cdot \mathbf{s}' \bmod M$  is upper bounded by  $2^\ell / M$ .*

**Proof:** It is easy to see that the values  $\mathbf{a}^\top \cdot \mathbf{s}$  and  $\mathbf{a}^\top \cdot \mathbf{s}'$  are independent and uniformly distributed for every pair  $\mathbf{s}, \mathbf{s}'$ . Hence,

$$\Pr[\exists \mathbf{s} \neq \mathbf{s}' : \mathbf{a}^\top \cdot \mathbf{s} \bmod M = \mathbf{a}^\top \cdot \mathbf{s}' \bmod M] \leq \sum_{\substack{\mathbf{s}'' \in \{0, 1\}^\ell \\ \mathbf{s}'' \neq \mathbf{0}^\ell}} \Pr[\mathbf{a}^\top \cdot \mathbf{s}'' \bmod M = 0] \leq \frac{2^\ell}{M}. \quad (1)$$

■

The hardness of solving the subset sum problem depends on the ratio between  $\ell$  and  $\log M$ , which is usually referred to as the *density*  $\Delta$  of the subset sum instance. In particular:



1. When  $\Delta < 1/\ell$ , we speak of *low-density* instances which can be solved in polynomial time by a reduction to a short vector in a lattice [LO85, Fri86, CJL<sup>+</sup>92].
2. When  $\Delta > \ell/\log^2 \ell$ , we speak of *high-density* instances which can be solved in polynomial time using dynamic programming, or other sophisticated techniques [CFG89, GM91, FP05, Lyu05, Sha08].

In our protocols, we will need to set the parameters in such a way that the subset sum problem is solvable efficiently. Furthermore, we need the term in Eq. (1) to be negligible in the security parameter  $\kappa$ ; hence we will set  $M = 2^{\kappa+\ell}$ . The latter choice immediately rules out algorithms for high-density subset sum (e.g., algorithms based on dynamic programming, since they usually need to process a matrix of dimension  $M$ ). On the other hand, for low-density instances, Lemma 1 implies  $\ell + \kappa > \ell^2$ , so that we need to choose  $\kappa, \ell$  in such a way that

$$\ell < \frac{1}{2} (\sqrt{1 + 4\kappa} + 1). \quad (2)$$

Algorithms for solving low-density subset sum are based on lattices. In particular, one can show [CJL<sup>+</sup>92] that all low-density subset sum instances with  $\Delta < 0.9408$  can be solved efficiently with overwhelming probability. The concrete run-time depends on the performances of the LLL algorithm as a function of the lattice dimension; values of  $\ell < 1000$  yield to practical performances [CN11, GN08, NS06]. We elaborate more on the impact of the above analysis in our constructions in Section 4.1.

### 2.3 Oblivious Pseudorandom Function Evaluation

Informally speaking, a pseudorandom function (PRF) is an efficiently computable function that looks like a truly random function to any PPT observer.

**Definition 1 (Pseudorandom function)** *Let  $F : \{0, 1\}^\kappa \times \{0, 1\}^m \rightarrow \{0, 1\}^l$  be an efficient, keyed function. We say  $F$  is a pseudorandom function if for all PPT distinguishers  $D$ , there exists a negligible function  $\text{negl}$  such that:*

$$|\Pr[D^{F^{(k, \cdot)}}(1^\kappa) = 1] - \Pr[D^{f_\kappa}(1^\kappa) = 1]| \leq \text{negl}(\kappa),$$

where  $k$  is picked uniformly from  $\{0, 1\}^\kappa$  and  $f_\kappa$  is chosen uniformly at random from the set of functions mapping  $\kappa$ -bit strings into  $l$ -bit strings.

In our protocols, we consider a protocol  $\pi_F$  that *obliviously* evaluates a pseudorandom function in the presence of malicious adversaries. Let  $k \in \{0, 1\}^\kappa$  be a key sampled as above. Then the oblivious PRF evaluation functionality  $\mathcal{F}_{\text{PRF}}$  is defined as  $(k, x) \mapsto (-, F(k, x))$ . Such an oblivious PRF may be instantiated with the Naor-Reingold pseudorandom function [NR97] that is implemented by the protocol presented in [FIPR05] (and proven in the malicious setting in [HL10]). The function is defined by

$$F((a_0, \dots, a_m), x) = g^{a_0 \prod_{i=1}^m a_i^{x[i]}},$$

where  $g$  is a generator for a group  $\mathbb{G}$  of prime order  $p$ ,  $a_i \in \mathbb{Z}_p$  and  $x = (x[1], \dots, x[m]) \in \{0, 1\}^m$ .<sup>3</sup> We remark that both the key and the range are not bit strings, as required by Definition 1, but they can be interpreted as such in a natural way. The protocol involves executing an oblivious transfer (OT) for every bit of the input  $x$ . A two-rounds semi-honest secure implementation can be achieved by using the [FIPR05] protocol combined with any two-rounds semi-honest oblivious transfer. Applying the efficient OT construction of [PVW08] in the malicious setting, implies a two-rounds protocol in the CRS setting with UC security and constant overhead.

<sup>3</sup>We remark that this definition considers a function that is not pseudorandom in the classic sense of it being indistinguishable from a random function whose range is composed of all strings of a given length. Rather, it is indistinguishable from a random function whose range is the group generated by  $g$  as defined below.

## 2.4 Commitment Schemes

A (non-interactive) commitment scheme consists of a triple of efficient algorithms  $(\text{Gen}, \text{Commit}, \text{Open})$  defined as follows. (i) Upon input the security parameter  $\kappa$ , the probabilistic algorithm  $\text{Gen}$  outputs a key  $\text{pk}$ . (ii) Upon input the key  $\text{pk}$  and message  $m \in \{0, 1\}^*$  (and implicit random coins  $r$ ), the probabilistic algorithm  $\text{Commit}$  outputs  $(\gamma, \delta) := \text{Commit}(\text{pk}, m; r)$  where  $\gamma$  is the committed value, while  $\delta$  is the decommitment information needed to open the commitment. Typically  $\delta := (m, r)$ . (iii) Upon input the key  $\text{pk}$ , and a commitment-pair  $(\gamma, \delta)$ , the deterministic algorithm  $\text{Open}$  outputs a value  $m \in \{0, 1\}^*$  or a special symbol  $\perp$ .

A commitment scheme should be complete, i.e., for any security parameter  $\kappa$ ,  $\text{pk} \leftarrow \text{Gen}(1^\kappa)$ , for any message  $m \in \{0, 1\}^*$  and  $(\gamma, \delta) \leftarrow \text{Commit}(\text{pk}, m)$  we have  $\text{Open}(\text{pk}, \delta, \gamma) = m$ . In addition, a commitment scheme has two properties known as binding and hiding:

**Binding Property.** A commitment scheme is *computationally* binding if for any PPT adversary  $\mathcal{A}$  we have

$$\Pr \left[ m \neq m' \wedge m, m' \neq \perp : \begin{array}{l} \text{pk} \leftarrow \text{Gen}(1^\kappa); (\delta, \delta', \gamma) \leftarrow \mathcal{A}(\text{pk}); \\ m \leftarrow \text{Open}(\text{pk}, \delta, \gamma); m' \leftarrow \text{Open}(\text{pk}, \delta', \gamma) \end{array} \right] \leq \text{negl}(\kappa).$$

We say that the scheme is perfectly binding if this property holds even for unbounded adversaries.

**Hiding Property.** For all messages  $m_0, m_1 \in \{0, 1\}^*$ , we have that

$$\{\text{pk}, \text{Commit}(\text{pk}, m_0)\}_{\kappa \in \mathbb{N}} \approx \{\text{pk}, \text{Commit}(\text{pk}, m_1)\}_{\kappa \in \mathbb{N}},$$

where the two ensembles are considered as random variables over the choice of the randomness to generate  $\text{pk} \leftarrow \text{Gen}(1^\kappa)$  and to compute the commitment. A scheme is computationally hiding if these two distributions are computationally indistinguishable (rather than identical).

For ease of notations we will often denote a commitment scheme by  $(\text{Commit}, \text{Open})$ , omitting the algorithm  $\text{Gen}$ .

## 2.5 Collision Resistant Hashing and Merkle Trees

Let in the following  $\{\mathcal{H}_\kappa\}_{\kappa \in \mathbb{N}} = \{H : \{0, 1\}^{p(\kappa)} \rightarrow \{0, 1\}^{p'(\kappa)}\}_{\kappa}$  be a family of hash functions, where  $p(\cdot)$  and  $p'(\cdot)$  are polynomials so that  $p'(\kappa) \leq p(\kappa)$  for sufficiently large  $\kappa \in \mathbb{N}$ . For a hash function  $H \leftarrow \mathcal{H}_\kappa$  a Merkle hash tree [Mer89] is a data structure that allows to commit to  $\ell = 2^d$  messages by a single hash value  $h$  such that revealing any message requires only to reveal  $O(d)$  hash values.

A Merkle hash tree is represented by a binary tree of depth  $d$  where the  $\ell$  messages  $m_1, \dots, m_\ell$  are assigned to the leaves of the tree; the values assigned to the internal nodes are computed using the underlying hash function  $H$ , whereas the value  $h$  that commits to  $m_1, \dots, m_\ell$  is assigned to the root of the tree. To open the commitment to a message  $m_i$ , one reveals  $m_i$  together with all the values assigned to nodes on the path from the root to  $m_i$ , and the values assigned to the siblings of these nodes. We denote the algorithm of committing to  $\ell$  messages  $m_1, \dots, m_\ell$  by  $h := \text{Commit}_M(m_1, \dots, m_\ell)$  and the opening of  $m_i$  by  $(m_i, \text{path}(i)) := \text{Open}_M(h, i)$ . Verifying the opening of  $m_i$  is carried out by essentially recomputing the entire path bottom-up and comparing the final outcome (i.e., the root) to the value given at the commitment phase. For simplicity, we abuse notation and denote by  $\text{path}(i)$  both the values assigned to the nodes in the path from the root to the decommitted value  $m_i$ , together with the values assigned to their siblings.

We often need to talk about the value assigned to a particular node. To this end, we introduce a labeling scheme for the nodes of a tree. We denote the root of the tree by  $\varepsilon$ . For a node  $w \in \bigcup_{i \leq d} \{0, 1\}^i$ , we label its left child by  $w0$  and its right child by  $w1$ . The value that is assigned to a node with a label  $w$  is typically

denoted by  $h_w$ . We also consider *incomplete* Merkle trees. An incomplete Merkle tree is a Merkle tree where some nodes  $w$ , with  $|w| < d$ , have *no* leaves. We say that a (possibly incomplete) Merkle tree  $\mathcal{T}$  with max depth  $d$  is *valid* if for all its nodes  $w$  with two children, we have  $H(h_{w0}||h_{w1}) = h_w$ . We further say that a path  $\text{path}(i)$  is *consistent* with a Merkle tree  $\mathcal{T}$  (or in  $\mathcal{T}$ ) if all the values assigned to the nodes  $w$  in  $\text{path}(i)$  are also assigned to the corresponding nodes in  $\mathcal{T}$ , i.e.,  $h_w = v_w$ , where  $v_w$  denotes the value assigned to node  $w$  in  $\text{path}(i)$ .

The binding property of a Merkle hash tree is due to collision resistance. Intuitively, this says that it is infeasible to efficiently find a pair  $(x, x')$  so that  $H(x) = H(x')$ , where  $H \leftarrow \mathcal{H}_\kappa$  for sufficiently large  $\kappa$ . In fact, one can show that collision resistance of  $\{\mathcal{H}_\kappa\}_{\kappa \in \mathbb{N}}$  carries over to the Merkle hashing. Formally, we say that a family of hash functions  $\{\mathcal{H}_\kappa\}_{\kappa}$  is collision resistant if for any PPT adversary  $\mathcal{A}$  the following experiment outputs 1 with probability  $\text{negl}(\kappa)$ : (i) A hash function  $H$  is sampled from  $\mathcal{H}_\kappa$ ; (ii) The adversary  $\mathcal{A}$  is given  $H$  and outputs  $x, x'$ ; (iii) The experiment outputs 1 if and only if  $x \neq x'$  and  $H(x) = H(x')$ .

We note that Merkle hash tree that does not imply hiding in the natural sense since collision resistance hash functions do not imply privacy. A standard way to obtain privacy is by applying a PRF or an encryption scheme on the messages before applying the construction of the hash tree. In this work we will implicitly use a PRF since the Merkle commitment is computed for the subset sum instance (that is the preprocessed text). Another useful property of Merkle hash tree is *succinctness* where the size of the commitment is independent of the number of committed messages. This property is particularly important since it allows to ensure binding while keeping a small local state, that corresponds to the root of the hash tree.

## 2.6 Universal Accumulators

A (static) universal accumulator is a tuple of algorithms (Init, CreateAcc, CreateWit, Check) specified as follows. (i) Algorithm Init takes as input a security parameter  $\kappa$  and returns a pair  $(\text{pk}_{\text{acc}}, \text{sk}_{\text{acc}})$ . (ii) Algorithm CreateAcc takes as input  $(\text{pk}_{\text{acc}}, \text{sk}_{\text{acc}})$  and a set  $G$  to be accumulated, and it outputs an accumulator  $h_G$  together with some auxiliary information *aux*. (iii) Algorithm CreateWit takes as input a boolean value  $\text{type} \in \{0, 1\}$ , a tuple  $(\text{pk}_{\text{acc}}, G, h_G, \text{aux})$ , and an element  $\gamma$ : if  $\text{type} = 0$ , it outputs a witness  $\zeta_\gamma$  for membership of  $\gamma$  in  $G$  or  $\perp$  (in case  $\gamma \notin G$ ); else, it outputs a witness  $\zeta_\gamma$  for non-membership of  $\gamma$  in  $G$  or  $\perp$  (in case  $\gamma \in G$ ). (iv) Algorithm Check takes as input a boolean value  $\text{type} \in \{0, 1\}$ , and a tuple  $(\text{pk}_{\text{acc}}, h_G, \zeta_\gamma, \gamma)$ : if  $\text{type} = 0$ , it outputs 1 if and only if  $\zeta_\gamma$  is a valid witness w.r.t.  $\gamma \in G$  and  $\text{pk}_{\text{acc}}$  for membership of  $\gamma \in G$ ; else, it outputs 1 if and only if  $\zeta_\gamma$  is a valid witness w.r.t.  $\gamma \notin G$  and  $\text{pk}_{\text{acc}}$  for non-membership of  $\gamma$  in  $G$ .

We say that the accumulator satisfies correctness if for all honestly generated keys, all honestly computed accumulators and witnesses, the Check algorithm returns 1 with overwhelming probability (over the coin tosses of all involved algorithms). As for security, collision freeness (defined below) informally states that it is neither feasible to find a witness for a non-accumulated value nor it is feasible to find a non-membership witness for an accumulated value.

**Definition 2** Let (Init, CreateAcc, CreateWit, Check) be a static, universal accumulator. We say that the accumulator is *collision-free* if for any PPT adversary  $\mathcal{A}$  there exists a negligible function  $\text{negl}(\cdot)$  such that:

$$\Pr \left[ \begin{array}{l} (\text{Check}(0, \text{pk}_{\text{acc}}, h_G, \zeta_{\gamma^*}, \gamma^*) = 1 \wedge \gamma^* \notin G) \\ \vee \\ (\text{Check}(1, \text{pk}_{\text{acc}}, h_G, \zeta_{\gamma^*}, \gamma^*) = 1 \wedge \gamma^* \in G) \end{array} : \begin{array}{l} (\text{pk}_{\text{acc}}, \text{sk}_{\text{acc}}) \leftarrow \text{Init}(1^\kappa) \\ (G, \gamma^*, \zeta_{\gamma^*}) \leftarrow \mathcal{A}^{\mathcal{O}}(\text{pk}_{\text{acc}}) \end{array} \right] \leq \text{negl}(\kappa),$$

where  $h_G \leftarrow \text{CreateAcc}(\text{pk}_{\text{acc}}, \text{sk}_{\text{acc}}, G)$  and  $\mathcal{O} := \text{CreateAcc}(\text{pk}_{\text{acc}}, \text{sk}_{\text{acc}}, \cdot)$ .

The definition above is weaker than the standard definition of collision freeness (see, e.g., [DHS15, Definition 8]), in that the standard definition even lets the adversary choose the randomness used to compute the

value  $h_G$  and gives  $\mathcal{A}$  access to a witness creator oracle for membership and non-membership of arbitrarily chosen elements. The simpler variant above is sufficient for our application.

Universal accumulators exist under a variety of assumptions, including strong RSA [LLX07],  $t$ -SDH [DT08, ATSM09], and based on collision-resistant hashing [BLL00, BLL02, CHKO12], vector commitments [CF13], and zero-knowledge sets [DHS15].

### 3 Modeling Outsourced Pattern Matching

The inputs for the basic pattern matching problem are a text  $T$  of length  $n$  and a pattern  $p$  (i.e., keyword) of length  $m$ ; the goal is to find all the text locations in which the pattern matches the text. A private distributed variant of this problem is defined in the two-party setting, where a sender SEN holds a text  $T$  and a receiver REC holds a pattern  $p$ . The goal of REC is to learn the positions in which  $p$  matches in the text, without revealing anything about the pattern to SEN; at the same time, REC should not learn anything else about the text.<sup>4</sup> In this section we are interested in an *outsourced* variant of the problem, which is specified in two phases. In the *setup phase* a sender SEN uploads a (preprocessed) text to an external server S. This phase is run only once. In the *query phase* receivers  $\text{REC}_1, \text{REC}_2, \dots$  query the text by searching patterns and learn the matched text locations. For simplicity, we focus on a single receiver REC asking multiple queries. However, our model can be easily generalized to the multiple receivers scenario.

The basic idea is to implement the pattern matching functionality using the server as a mediator, answering search queries on behalf of SEN. In order to take some advantage from this modeling, we must allow the server to obtain some leakage about the text, otherwise the communication complexity between the server and receiver REC would be  $O(n)$ . Instead, we are interested in building schemes where the preprocessing phase requires  $O(n)$  workload, but the overall cost of issuing a query grows only linearly with the number of matches (which is as optimal as one can obtain). This optimization comes with the price of revealing some leakage about the text. More precisely, the server learns that for some text positions repetitions occur. Attempts to hide this information from the server by permuting the text fail if the server colludes with REC. For some applications this leakage is tolerable given the improvement of running search queries. To sum up, we aim for  $O(n)$  computation/communication overhead in the preprocessing phase and  $O(t_p)$  in the query phase, where  $t_p$  is the number of occurrences of  $p$  in  $T$ .

We further require that the round complexity of any protocol implemented in this setting is minimal. That is, in the setup phase we require a single message sent from the sender to the server, whereas in the query phase we require clients REC and SEN to exchange only two messages (one in each direction), and one message in each direction between REC and S in order to retrieve the output. We note that our constructions meet this order of rounds, but this may not be the case in general. We denote a scheme with this number of rounds by *round optimal*.

We formalize security using the ideal/real paradigm. Note that, in the context of outsourced computation, the server is a separate entity that does not contribute any input to the computation and is required to run most of the function evaluation. In the ideal setting, such an entity is also communicating with the functionality and, upon corruption, decides whether the functionality sends the outcome of the computation to the prescribed receivers. Denote by  $T_j$  the substring of length  $m$  that starts at text location  $j$  (where  $m$  is the length of the query). The pattern matching ideal functionality in the outsourced setting is depicted in Figure 2. We write  $|T|$  for the bit length of  $T$  and assume that receiver REC asks a number of queries  $p_i$  ( $i \in [\lambda]$ ) where the queries' lengths are always bounded by  $m$ .

---

<sup>4</sup>As specified in the introduction, we can define a search functionality that only returns the first few and most relevant matches. Our discussion and formalization below can be easily adapted for this functionality as well.

**Functionality  $\mathcal{F}_{\text{OPM}}$**

Let  $m, \lambda \in \mathbb{N}$ . Functionality  $\mathcal{F}_{\text{OPM}}$  sets an empty table  $\mathcal{B}$  and proceeds as follows, running with clients SEN and REC, server S and adversary Sim.

1. Upon receiving a message  $(\text{text}, T, m)$  from SEN, send  $(\text{preprocess}, |T|, m)$  to S and Sim, and record  $(\text{text}, T)$ .
2. Upon receiving a message  $(\text{query}, p_i)$  from receiver REC (for  $i \in [\lambda]$ ), where message  $(\text{text}, \cdot)$  has been recorded and  $|p_i| = m$ , it checks if the table  $\mathcal{B}$  already contains an entry of the form  $(p_i, \cdot)$ . If this is not the case then it picks the next available identifier id from  $\{0, 1\}^*$  and adds  $(p_i, \text{id})$  to  $\mathcal{B}$ . It sends  $(\text{query}, \text{REC})$  to SEN and Sim.
  - (a) Upon receiving  $(\text{approve}, \text{REC})$  from sender SEN, read  $(p_i, \text{id})$  from  $\mathcal{B}$  and send  $(\text{query}, \text{REC}, (i_1, \dots, i_t), \text{id})$  to server S, for all text positions  $\{i_j\}_{j \in [t]}$  such that  $T_{i_j} = p_i$ . Otherwise, if no  $(\text{approve}, \text{REC})$  message has been received from SEN, send  $\perp$  to REC and abort.
  - (b) Upon receiving  $(\text{approve}, \text{REC})$  from Sim, read  $(p_i, \text{id})$  from  $\mathcal{B}$  and send  $(\text{query}, p_i, (i_1, \dots, i_t), \text{id})$  to receiver REC. Otherwise, send  $\perp$  to receiver REC.

Figure 2: The outsourced pattern matching functionality

**The definition.** As in the standard static modeling, a corrupted party is either passively or actively controlled by an adversarial entity. In the passive case (a.k.a. semi-honest case) a corrupted party follows the protocol’s instructions and tries to gain additional information about the honest parties’ inputs from its view; in the active case (a.k.a. malicious case) a corrupted party is allowed to follow an arbitrary polynomial-time strategy. In our case, when the server is corrupted we must ensure that the only information leaked by the protocol is about the text positions for which repetitions occur, without disclosing the actual content of the text in these positions or any additional information. A moment of reflection shows that in the security proof the simulator needs to commit to the text *before* given the above leakage from the trusted party. We emphasize that this technicality is not artificial, since even if receiver REC is corrupted at the beginning of the execution, the simulator cannot be given the leakage about the queries in advance since the queries may be asked in an fully adaptive manner. In other words, it may be the case that receiver REC does not know all the queries it will ask in advance.

Formally, denote by  $\text{IDEAL}_{\mathcal{F}_{\text{OPM}}, \text{Sim}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))$  the output of an ideal adversary Sim, server S and clients SEN, REC in the above ideal execution of  $\mathcal{F}_{\text{OPM}}$  upon inputs  $(-, (T, (p_1, \dots, p_\lambda)))$  and auxiliary input  $z$  given to Sim. We note that one can also consider a security definition that captures collusion between the server and one of the clients. Our protocols capture collusion between S and REC. We implement functionality  $\mathcal{F}_{\text{OPM}}$  via three two-party protocols  $\pi = (\pi_{\text{Pre}}, \pi_{\text{Query}}, \pi_{\text{Opm}})$  specified as follows. Protocol  $\pi_{\text{Pre}}$  is run in the preprocessing phase by SEN to preprocess text  $T$  and forwards the outcome to S. During the query phase, protocol  $\pi_{\text{Query}}$  is run between SEN and REC (holding a pattern  $p$ ); this protocol outputs a trapdoor  $R_p$  that depends on  $p$  and will enable the server to search the preprocessed text. Lastly, protocol  $\pi_{\text{Opm}}$  is run by S upon input the preprocessed text and a trapdoor (forwarded by REC); this protocol returns to REC the matched text positions (if any). We denote by  $\text{REAL}_{\pi, \mathcal{A}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))$  the output of adversary  $\mathcal{A}$ , server S and clients SEN, REC in a real execution of  $\pi = (\pi_{\text{Pre}}, \pi_{\text{Query}}, \pi_{\text{Opm}})$  upon inputs  $(-, (T, (p_1, \dots, p_\lambda)))$  and auxiliary input  $z$  given to  $\mathcal{A}$ .

**Definition 3 (Security of outsourced pattern matching)** *We say that  $\pi$  securely implements  $\mathcal{F}_{\text{OPM}}$ , if for any PPT real adversary  $\mathcal{A}$  there exists a PPT ideal adversary (simulator) Sim such that for any tuple of*

inputs  $(T, (p_1, \dots, p_\lambda))$  and auxiliary input  $z$ ,

$$\{\mathbf{IDEAL}_{\mathcal{F}_{\text{OPM}}, \text{Sim}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}} \stackrel{c}{\approx} \{\mathbf{REAL}_{\pi, \mathcal{A}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}}.$$

The schemes described in the next sections implement the ideal functionality  $\mathcal{F}_{\text{OPM}}$  in the random oracle model. As usual in this case, we assume the existence of a publicly available function  $H$  behaving as a truly random function. We stress that all parties are allowed to query the random oracle arbitrarily, regardless the fact that they are malicious or semi-honest.

**The  $\mathcal{F}$ -hybrid model.** Our constructions use secure two-party protocols as subprotocols. The standard way of doing this is to work in a “*hybrid model*” where parties both interact with each other (as in the real model) and use trusted help (as in the ideal model). Specifically, when constructing a protocol  $\pi$  that uses a subprotocol for securely computing some functionality  $\mathcal{F}$ , we consider the case that the parties run  $\pi$  and use “ideal calls” to a trusted party for computing  $\mathcal{F}$ . Upon receiving the inputs from the parties, the trusted party computes  $\mathcal{F}$  and sends all parties their output. Then, after receiving these outputs back from the trusted party the protocol  $\pi$  continues. Let  $\mathcal{F}$  be a functionality and let  $\pi$  be a two-party protocol that uses ideal calls to a trusted party computing  $\mathcal{F}$ . Furthermore, let  $\mathcal{A}$  be a probabilistic polynomial-time machine. Then, the  *$\mathcal{F}$ -hybrid execution of  $\pi$*  on inputs  $(T, (q_1, \dots, q_\lambda))$ , auxiliary input  $z$  to  $\mathcal{A}$  and security parameter  $\kappa$ , denoted  $\mathbf{HYB}_{\pi, \mathcal{A}(z)}(\kappa, (-, T, (q_1, \dots, q_\lambda)))$ , is defined as the output vector of the honest parties and the adversary  $\mathcal{A}$  from the hybrid execution of  $\pi$  with a trusted party computing  $\mathcal{F}$ . By the composition theorem of [Can00] any protocol that securely implements  $\mathcal{F}$  can replace the ideal calls to  $\mathcal{F}$ .

## 4 Security in the Presence of Semi-Honest Adversaries

In this section we present our implementation of the outsourced pattern matching functionality  $\mathcal{F}_{\text{OPM}}$  and prove its security against semi-honest adversaries. A scheme with security against malicious adversaries is described in Section 5, building upon the protocol in this section. Recall first that in the outsourced variant of the pattern matching problem, sender SEN manipulates the text  $T$  and then stores it on the server S in such a way that the preprocessed text can be used later to answer search queries submitted by receiver REC. The challenge is to find a way to hide the text (in order to obtain privacy), while enabling the server to carry out searches on the hidden text whenever it is in possession of an appropriate trapdoor.

We consider a new approach and reduce the pattern matching problem to the subset sum problem (cf. Section 2.2). Namely, consider a text  $T$  of length  $n$ , and assume we want to allow searches for patterns of length  $m$ . For some integer  $M \in \mathbb{N}$ , we assign to each distinct pattern  $p$  that appears in  $T$  a random element  $R_p \in \mathbb{Z}_M$ . Letting  $\ell = n - m + 1$ , the preprocessed text  $\tilde{\mathbf{T}}$  is a vector in  $\mathbb{Z}_M^\ell$  with elements specified as follows. Specifically, for each pattern  $p$  that appears  $t$  times in  $T$ , we sample random values  $a_1, \dots, a_t \in \mathbb{Z}_M$  such that  $R_p = \sum_{j=1}^t a_j$ . Denote with  $i_j \in [\ell]$  the  $j$ th position in  $T$  where  $p$  appears and set  $\tilde{\mathbf{T}}[i_j] = a_j$ . Notice that for each pattern  $p$ , there exists a vector  $\mathbf{s} \in \{0, 1\}^\ell$  such that  $R_p = \tilde{\mathbf{T}}^\top \cdot \mathbf{s}$ . Hence, the positions in  $\tilde{\mathbf{T}}$  where pattern  $p$  matches are identified by a vector  $\mathbf{s}$  and can be viewed as the solution for the subset sum problem instance  $(R_p, \tilde{\mathbf{T}})$ .

Roughly, our protocol works as follows. During protocol  $\pi_{\text{Pre}}$ , we let the sender SEN generate the preprocessed text  $\tilde{\mathbf{T}}$  as described above, and send the result to the server S. Later, when a receiver REC wants to learn at which positions a pattern  $p$  matches the text, clients REC and SEN run protocol  $\pi_{\text{Query}}$ ; at the end of this protocol REC learns the trapdoor  $R_p$  corresponding to  $p$ . Finally, during  $\pi_{\text{OPM}}$  receiver REC sends this trapdoor to S, which can solve the subset sum problem instance  $(R_p, \tilde{\mathbf{T}})$ . The solution to this problem corresponds to the matches of  $p$ , which are forwarded to the receiver REC. To avoid that SEN

**Protocol**  $\pi_{\text{SH}} = (\pi_{\text{Pre}}, \pi_{\text{Query}}, \pi_{\text{Opm}})$

Let  $\kappa \in \mathbb{N}$  be the security parameter and let  $M, m, n, \mu$  be integers, where for simplicity we assume that  $n$  is a multiple of  $2m$ . Further, let  $\mathcal{H} : \{0, 1\}^\mu \rightarrow \mathbb{Z}_M$  be a random oracle and  $F : \{0, 1\}^\kappa \times \{0, 1\}^m \rightarrow \{0, 1\}^\mu$  be a PRF. Protocol  $\pi_{\text{SH}}$  involves a sender SEN holding a text  $T \in \{0, 1\}^n$ , a receiver REC querying for patterns  $p \in \{0, 1\}^m$ , and a server S. The interaction between the parties is specified below.

**Setup phase,  $\pi_{\text{Pre}}$ .** The protocol is invoked between sender SEN and server S. Given input  $T$  and integer  $m$ , sender SEN picks a random key  $k \in \{0, 1\}^\kappa$  and prepares first the text  $T$  for the packaging by writing it as

$$T' := (B_1, \dots, B_u) = ((T[1], \dots, T[2m]), (T[m+1], \dots, T[3m]), \dots, (T[n-2m+1], \dots, T[n])),$$

where  $u = n/m - 1$ . Next, for each block  $B_b$  and each of the  $m + 1$  patterns  $p \in \{0, 1\}^m$  that appear in  $B_b$  we proceed as follows (suppose there are at most  $t$  matches of  $p$  in  $B_b$ ).

1. Sender SEN evaluates  $R_p := \mathcal{H}(F(k, p) || b)$ , samples  $a_1, \dots, a_{t-1} \in \mathbb{Z}_M$  at random and then fixes  $a_t$  such that  $a_t = R_p - \sum_{j=1}^{t-1} a_j \bmod M$ .
2. Set  $\tilde{B}_b[v_j] = a_j$  for all  $j \in [t]$  and  $v_j \in [m + 1]$ . Note that here we denote by  $\{v_j\}_{j \in [t]}$  ( $v_j \in [m + 1]$ ) the set of indexes corresponding to the positions where  $p$  occurs in  $B_b$ . Later in the proof we will be more precise and explicitly denote to which block  $v_j$  belongs by using explicitly the notation  $v_{j_b}$ .

Finally, SEN outsources the text  $\tilde{\mathbf{T}} = (\tilde{B}_1, \dots, \tilde{B}_u)$  to S.

**Query phase,  $\pi_{\text{Query}}$ .** Upon issuing a query  $p \in \{0, 1\}^m$  by receiver REC, clients SEN and REC engage in an execution of protocol  $\pi_{\text{Query}}$  which implements the oblivious PRF functionality  $(k, p) \mapsto (-, F(k, p))$ . Upon completion, REC learns  $F(k, p)$ .

**Oblivious pattern matching phase,  $\pi_{\text{Opm}}$ .** This protocol is executed between server S (holding  $\tilde{\mathbf{T}}$ ) and receiver REC (holding  $F(k, p)$ ). Upon receiving  $F(k, p)$  from REC, the server proceeds as follows for each block  $\tilde{B}_b$ . It interprets  $(\tilde{B}_b, \mathcal{H}(F(k, p) || b))$  as a subset sum instance and computes  $\mathbf{s}$  as the solution of  $\tilde{B}_b \cdot \mathbf{s} = \mathcal{H}(F(k, p) || b)$ . Let  $\{v_j\}_{j \in [t]}$  denote the set of indexes such that  $\mathbf{s}[v_j] = 1$ , then the server S returns the set of indexes  $\{\varphi(b, v_j)\}_{b \in [u], j \in [t]}$  to the receiver REC.

Figure 3: Semi-honest outsourced pattern matching

needs to store all trapdoors, we rely on a PRF to generate the trapdoors itself. More precisely, instead of sampling the trapdoors  $R_p$  uniformly at random, we set  $R_p := F(k, p)$ , where  $F$  is a PRF. Thus, during the query phase REC and SEN run an execution of an oblivious PRF protocol; at the end of this protocol REC learns the output of the PRF, i.e., the trapdoor  $R_p$ .

Although the protocol described above provides a first basic solution for the outsourced pattern matching, it suffers from a strong restriction as only very short texts are supported. We will provide more details on this restriction in Section 4.1.<sup>5</sup> To overcome this severe limitation, we partition the text into smaller pieces each of length  $2m$ , where each such piece is handled as a separate instance of the the protocol. More specifically, for a text  $T = (T[1], \dots, T[n])$  let  $(T[1], \dots, T[2m]), (T[m + 1], \dots, T[3m]), \dots$  be blocks, each of length  $2m$ , such that every two consecutive blocks overlap in  $m$  bits. Then, for each pattern  $p$  that appears in the text the sender SEN computes an individual trapdoor for each block where the pattern  $p$  appears. More precisely, suppose that pattern  $p$  appears in block  $B_b$  then we compute the trapdoor for this

<sup>5</sup>Taking a look ahead, this solution can only support short texts as otherwise either the collision probability (cf. Lemma 1) is large and we cannot achieve correctness, or the subset sum problem cannot be solved efficiently as the instance is not low-density.

block (and pattern  $p$ ) as  $\mathcal{H}(F(k, p) || b)$ . Here,  $\mathcal{H}$  is a cryptographic hash function that will be modeled as a random oracle in our proofs. Given the trapdoors, we apply the preprocessing algorithm to each block individually. The sub-protocols  $\pi_{\text{Query}}$  and  $\pi_{\text{Opm}}$  work as described above with a small change. In  $\pi_{\text{Query}}$  receiver REC learns the output of the PRF  $F(k, p)$  instead of the actual trapdoors and in  $\pi_{\text{Opm}}$  receiver REC forwards directly the result  $F(k, p)$  to S. The server can then compute the actual trapdoor using the random oracle. This is needed to keep the communication complexity of the protocol low. Note that in this case if we let  $\{v_{j_b}\}_{j_b \in [t_b]}$  be the set of indices corresponding to the positions where  $p$  occurs in a given block  $B_b$ , the server needs to map these positions to the corresponding positions in  $T$  (and this has to be done for each of the blocks where  $p$  matches). It is easy to see that such a mapping from a position  $v_{j_b}$  in block  $B_b$  to the corresponding position in the text  $T$  can be computed as  $\varphi(b, v_j) = (b - 1)m + v_j$ . The entire protocol, including the packaging mechanism, is shown in Figure 3; see also Figure 1 for a pictorial representation.

We now prove the following result in the  $\mathcal{F}_{\text{PRF}}$ -hybrid model (where  $\mathcal{F}_{\text{PRF}}$  denotes the oblivious PRF evaluation functionality and is defined in Section 2.3).

**Theorem 1** *Let  $\kappa \in \mathbb{N}$  be the security parameter. For integers  $n, m$  we set  $\lambda = \text{poly}(\kappa), \mu = \text{poly}(\kappa), u = n/m - 1, \ell = (m + 1)u$  and  $M = 2^{m+\kappa+1}$ . We furthermore require that  $\kappa$  is such that  $2^{m+1}/M$  is negligible (in  $\kappa$ ). Assume  $\mathcal{H} : \{0, 1\}^\mu \rightarrow \mathbb{Z}_M$  is a random oracle and  $F : \{0, 1\}^\kappa \times \{0, 1\}^m \rightarrow \{0, 1\}^\mu$  is a pseudorandom function. Then, the round optimal protocol  $\pi_{\text{SH}}$  from Figure 3 securely implements the  $\mathcal{F}_{\text{OPM}}$  functionality in the presence of semi-honest adversaries in the  $\mathcal{F}_{\text{PRF}}$ -hybrid model.*

It is easy to verify that if we realize  $\mathcal{F}_{\text{PRF}}$  in two rounds (such as the implementation for [NR97]) then the number of rounds within  $\pi_{\text{SH}}$  is optimal. First, SEN sends only one message to S. Next, we consider a two-rounds oblivious PRF evaluation protocol for  $\pi_{\text{Query}}$ . Finally, REC and S exchange only two messages. We now continue with our security proof.

**Proof:** We first argue about correctness and then prove privacy.

**Correctness.** We say that our construction achieves correctness if with overwhelming probability each pattern query  $p$  issued by REC is answered correctly with respect to the outsourced text  $T$ . More concretely, for each pattern  $p \in \{0, 1\}^m$  and a text  $T \in \{0, 1\}^n$ , let  $\{i_j\}_{j \in [t]}$  be the positions in  $T$  where  $p$  matches the text. Then, our protocol achieves correctness if for any  $T$  and  $p$  it returns the correct matches  $\{i_j\}_{j \in [t]}$ . Suppose that pattern  $p$  appears at position  $i_j$ . We need to show that in this case the algorithm given in Figure 3 returns indeed index  $i_j$ . To this end, suppose that position  $i_j$  lies in block  $B_b$  and  $B_{b+1}$  for some  $b \in [1, v - 1]$  (recall that two consecutive blocks overlap at  $m$  positions, hence the bit  $T[i_j]$  may appear in both blocks; in case we have a match of a pattern exactly in the area that overlaps, then we will always consider only the match in the first block). Wlog. assume that pattern  $p$  appears in block  $B_{b+1}$ . We run the preprocessing algorithm on block  $B_{b+1}$  to obtain the transformed block  $\tilde{B}_{b+1}$  that contains the solution of the subset sum problem for  $(\mathcal{H}(F(k, p) || b + 1), \tilde{B}_{b+1})$  at positions where the pattern  $p$  appears. Hence, when during the execution of protocol  $\pi_{\text{Opm}}$  the server solves this subset sum instance, it will retrieve index  $i_j$  as one of the solutions.

The analysis from above does not hold when one of the following situations occur. First, it may be the case that for two different patterns  $p \neq p'$  we get a collision in the PRF and/or random oracle. The probability that this happens is negligible by the birthday bound. Next, we consider the following two cases when we get a non-unique solution for the subset sum problem:

1. *Two (or more) different subsets that sum to the same trapdoor in some block:* From Lemma 1, it follows that for each subset sum instance collision happens with probability  $2^{2m}/M$ . Taking the union bound we get that the probability of collision for all patterns  $p$  (appearing in some block of length  $2m$ ) is upper bounded by  $2m \cdot 2^{2m}/M$ , which for our choice of parameters is negligible in  $\kappa$ .



2. *There is no match in a block, but a subset in this block sums to trapdoor:* For each trapdoor  $R_p$  the probability that some subset in a block sums to  $R_p$  is upper bounded by  $1/M$ . As there are  $2^m$  possible targets, we get by the union bound that the probability of this event is upper bounded by  $2^m/M$  which is negligible in  $\kappa$ .

As both events above occur with a negligible probability and there are  $u = n/m - 1$  blocks, we can apply the union bound, resulting in a negligible (in  $\kappa$ ) probability of an incorrect result for our protocol. This concludes the correctness proof of our protocol.

**Privacy.** We will show that for any PPT real adversary  $\mathcal{A}$  there exists a PPT ideal adversary (simulator)  $\text{Sim}$  such that for any tuple of inputs  $(T, (p_1, \dots, p_\lambda))$  and auxiliary input  $z$ ,

$$\{\mathbf{IDEAL}_{\mathcal{F}_{\text{OPM}}, \text{Sim}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}} \stackrel{c}{\approx} \{\mathbf{REAL}_{\pi_{\text{SH}}, \mathcal{A}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}}.$$

We prove each corruption setting separately. The final simulator from above can then be obtained by combining the individual simulators.

**A corrupted server.** We begin with the case when the server is corrupted. Let  $\mathcal{A}$  denote the adversary controlling the server  $S$ . We will build a simulator  $\text{Sim}_S$  that simulates the view of  $\mathcal{A}$ . To this end  $\text{Sim}_S$  interacts with the trusted party and uses the leakage from the trusted party, which for each pattern query reveals the matched text positions. We first describe the simulator  $\text{Sim}_S$  with access to  $\mathcal{A}$ .

**Convention:** During the simulation  $\text{Sim}_S$  evaluates queries to the random oracle  $\mathcal{H}$ . Such queries are made by  $\mathcal{A}$  or by  $\text{Sim}_S$  during its simulation of the corrupted server. To evaluate  $\mathcal{H}(x)$ ,  $\text{Sim}_S$  first checks if it has already recorded a pair  $(x, r)$ , in which case  $\mathcal{H}(x)$  evaluates to the value  $r$ . Otherwise,  $\text{Sim}_S$  chooses a random string  $r \in \mathbb{Z}_M$ , records  $(x, r)$  and evaluates  $\mathcal{H}(x)$  to  $r$ .

1. On input the auxiliary input  $z$ ,  $\text{Sim}_S$  invokes  $\mathcal{A}$  (i.e., the corrupted server) on this input. The simulator keeps track of a table  $\mathcal{B}$  that is initially set to the empty table.
2. Upon receiving a (preprocess,  $n, m$ ) message from the trusted party denoting that the honest SEN wants to outsource a text of length  $n$  to the trusted party,  $\text{Sim}_S$  defines text  $\tilde{\mathbf{T}}$  by sampling uniformly at random a vector of length  $\ell := (m + 1)u$  from  $\mathbb{Z}_M^\ell$ . It forwards  $\tilde{\mathbf{T}}$  to adversary  $\mathcal{A}$  and stores it for later usage as well.
3. Upon receiving a (query, REC,  $(i_1, \dots, i_t)$ , id) message from the trusted party indicating that receiver REC submitted a search query that was approved by SEN,  $\text{Sim}_S$  distinguishes two cases.
  - (a) *Pattern queried by REC appears in  $T$ :* In this case  $\{i_j\}_{j \in [t]}$  is not the empty set, and the simulator samples uniformly at random a value  $X_{\text{id}}$  from  $\{0, 1\}^\mu$  (this will take the role of  $F(k, p)$  in the real execution). Then, it proceeds as follows for each of the  $i_j$ 's. It first computes the block number  $b = \lfloor i_j/m \rfloor + 1$  in which the index  $i_j$  occurs and then the starting position  $v_{j_b} = i_j \bmod m$  where the pattern appears in  $\tilde{B}_b$ . Then,  $\text{Sim}_S$  programs the random oracle  $\mathcal{H}(X_{\text{id}}||b)$  to  $\sum_{j_b=1}^{t_b} \tilde{\mathbf{T}}[v_{j_b}]$ ; if  $\mathcal{H}$  has already been programmed to a different value, then we abort.
  - (b) *Pattern does not appear in  $T$ :* In this case  $\{i_j\}_{j \in [t]}$  is the empty set, and we check if table  $\mathcal{B}$  contains a value of the form  $(\text{id}, X_{\text{id}})$ . Otherwise, we pick the value  $X_{\text{id}}$  uniformly at random in  $\{0, 1\}^\mu$  and store  $(\text{id}, X_{\text{id}})$  in  $\mathcal{B}$ .

Finally,  $\text{Sim}_S$  (emulating the role of REC in the real execution) forwards  $X_{\text{id}}$  to the adversary.

4. If  $\mathcal{A}$  does not answer with  $(i_1, \dots, i_t)$ ,  $\text{Sim}_S$  sends  $\perp$  to ideal functionality and abort. Otherwise it sends the trusted party (approve, REC).
5.  $\text{Sim}_S$  outputs whatever  $\mathcal{A}$  does.

We first note that  $\text{Sim}_S$  runs in polynomial time since it only samples a random vector from  $\mathbb{Z}_M^\ell$ , and then calculates the sum of values from a given subset. Next, we show that the distribution produced by  $\text{Sim}_S$  in the ideal world is computationally indistinguishable from the distribution that  $\mathcal{A}$  expects to see in the real world. This is required to hold even given the leakage revealing the positions where the pattern matches the text. We start by defining a hybrid distribution  $\mathbf{HYB}_{\pi_{\text{SH}}, \mathcal{A}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))$  that is defined as the real experiment  $\mathbf{REAL}_{\pi_{\text{SH}}, \mathcal{A}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))$  with the difference that  $X_{\text{id}}$  is not computed by a PRF but rather by a random function  $f_\kappa$ . The following claim formalizes this statement.

**Claim 2** *Let  $F$  be a secure pseudorandom function (cf. Definition 1), then there exists a negligible function  $\text{negl}(\cdot)$  such that for sufficiently large  $\kappa \in \mathbb{N}$  and for any tuple of inputs  $(T, (p_1, \dots, p_\lambda))$  and auxiliary input  $z$ , it holds that*

$$\{\mathbf{REAL}_{\pi_{\text{SH}}, \mathcal{A}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}} \stackrel{c}{\approx} \{\mathbf{HYB}_{\pi_{\text{SH}}, \mathcal{A}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}}. \quad (3)$$

The proof follows by an easy reduction to the pseudorandomness of  $F$ , as a distinguisher for the distributions in Eq. (3) yields a distinguisher for the PRF.

To move to the simulated view, we need to bound the distance between the experiment  $\mathbf{HYB}_{\pi_{\text{SH}}, \mathcal{A}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))$  and the simulated view. To this end, we define the event  $\text{bad}$  that occurs when the simulator aborts in the ideal world.

- Event  $\text{bad}$ : Occurs if the simulation given above is aborted. In this case the corrupted server has asked for a direct query to the random oracle of the form  $(f_\kappa(p) || b)$  before it has seen  $f_\kappa(p)$ , where  $b \in [u]$  and  $p$  is a pattern that occurs in the text  $T$ . Recall that  $f_\kappa$  is a random function as defined in the hybrid world.

We will show that the distribution produced by the simulator  $\text{Sim}_S$  in the ideal world is statistically close to the distribution produced in the hybrid world.

**Claim 3** *For any input text  $T$ , patterns  $p_1, \dots, p_\lambda$ , and auxiliary input  $z$ , it holds that*

$$\{\mathbf{IDEAL}_{\mathcal{F}_{\text{OPM}}, \text{Sim}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}} \equiv_s \{\mathbf{HYB}_{\pi_{\text{SH}}, \mathcal{A}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}}.$$

**Proof:** We show that the view generated by  $\text{Sim}_S$  for the corrupted server in the ideal world has the same distribution as the server expects to see in the hybrid protocol conditioned on the event  $\text{bad}$ . The view contains the preprocessed text, answers to random oracle queries and the trapdoors that are sent by REC during the sub-protocol  $\pi_{\text{OPM}}$ . In the ideal world,  $\text{Sim}_S$  chooses these values in the following way.

- *Answers to direct RO queries:* If the RO has already been asked on this input, then it returns the stored value; otherwise it returns a value chosen uniformly at random,
- *Preprocessed text:* Sampled independently and uniformly from  $\mathbb{Z}_M^\ell$ ,
- *Trapdoors sent by REC:* We first check if table  $\mathcal{B}$  contains an entry of the form  $(\text{id}, X_{\text{id}})$ , in which case we return  $X_{\text{id}}$ . Otherwise, we consider two cases depending on the query  $(\text{query}, \text{REC}, (i_1, \dots, i_t), \text{id})$ :
  1. *Pattern matches the text:*  $X_{\text{id}}$  is chosen uniformly at random from  $\{0, 1\}^\mu$ . For each block  $B_b$  where the pattern appears, we program the random oracle  $\mathcal{H}(X_{\text{id}} || b)$  to  $\sum_{j_b=1}^{t_b} \tilde{\mathbf{T}}[v_{j_b}]$ , where  $v_{j_b}$  are the positions in block  $\tilde{B}_b$  where the pattern appears. We store  $(\text{id}, X_{\text{id}})$  in  $\mathcal{B}$ .

2. *Pattern does not match the text:* We pick  $X_{\text{id}}$  at random from  $\{0, 1\}^\mu$  and store  $(\text{id}, X_{\text{id}})$  in  $B$ .

Notice that in the first case the trapdoor is a uniformly chosen value from  $\mathbb{Z}_M$  as it is the sum of uniformly and independently chosen values.

It is easy to see that individually these values are identically distributed in both the ideal and hybrid execution. For the proof, we need to analyze the joint distribution that the corrupted server sees. The only difference between the joint distribution in the ideal world and in the hybrid world is the way in which the preprocessed text and the trapdoors are sampled. While in the ideal world the preprocessed text is sampled uniformly and independently from  $\mathbb{Z}_M^\ell$  in the hybrid world we prepare it according to the patterns that appear in the text  $T$ . More precisely, in the hybrid world we put at locations where a pattern appears a random additive sharing of a random value. This trivially implies that also in the hybrid world the transformed text is sampled uniformly at random from  $\mathbb{Z}_M^\ell$ . It remains to argue that at the later stage when the receiver REC asks for patterns, the view in the ideal world remains consistent, namely, the sum of the values that are put at the matched positions is equal to the trapdoor. We can achieve this consistency by programming the value of the random oracle to the appropriate value. There is one exception when such programming fails: namely, when the adversary has earlier queried the random oracle on this value. This is exactly when event bad happens and the simulator aborts. It remains to show that the probability that event bad happens is negligible in the security parameter.

Event bad occurs when the adversary asks the random oracle on a value  $X$  with form  $(f_\kappa(p)||b)$  for some  $b \in [u]$  and pattern  $p$  that appears in the text *before* it actually sees  $X$ . As  $f_\kappa$  is a random function this happens with probability at most  $\text{poly}(\kappa)/2^\mu = \text{negl}(\kappa)$ . This concludes the proof. ■

Combining Claims 2 and 3, it holds that the distributions  $\{\mathbf{IDEAL}_{\mathcal{F}_{\text{OPM}}, \text{Sim}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}}$  and  $\{\mathbf{REAL}_{\pi_{\text{SH}}, \mathcal{A}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}}$  are computationally close, which concludes the proof of privacy in case of a corrupted (semi-honest) server.

**A corrupted sender.** Next, we consider the case when SEN is corrupted. Let  $\mathcal{A}$  denote an adversary controlling sender SEN, we build a simulator  $\text{Sim}_{\text{SEN}}$  that generates its view. The simulator  $\text{Sim}_{\text{SEN}}$  needs to emulate the roles of receiver REC and of the server S using the leakage it gets from the trusted party. Let us describe first the simulator  $\text{Sim}_{\text{SEN}}$  that is given access to adversary  $\mathcal{A}$ .

**Convention:** During the simulation  $\text{Sim}_{\text{SEN}}$  evaluates queries to the random oracle  $\mathcal{H}$ . Such queries are made by  $\mathcal{A}$  or by  $\text{Sim}_{\text{SEN}}$  during its simulation of the corrupted SEN. To evaluate  $\mathcal{H}(x)$ ,  $\text{Sim}_{\text{SEN}}$  first checks if it has already recorded a pair  $(x, r)$ , in which case  $\mathcal{H}(x)$  evaluates to the value  $r$ . Otherwise,  $\text{Sim}_{\text{SEN}}$  chooses a random string  $r \in \mathbb{Z}_M$ , records  $(x, r)$  and evaluates  $\mathcal{H}(x)$  to  $r$ .

1. On input text  $T$ , length  $m$ , and auxiliary input  $z$ ,  $\text{Sim}_{\text{SEN}}$  invokes  $\mathcal{A}$  on these inputs.
2. Upon receiving  $\tilde{\mathbf{T}}$  from  $\mathcal{A}$ ,  $\text{Sim}_{\text{SEN}}$  sends a  $(\text{text}, T, m)$  message to the trusted party and receives back  $(\text{preprocess}, |T|, m)$ .
3. Upon receiving a  $(\text{query}, \text{REC})$  message from the trusted party, denoting that the honest REC submitted a query  $p$  of length  $m$  to the trusted party,  $\text{Sim}_{\text{SEN}}$  emulates the role of  $\mathcal{F}_{\text{PRF}}$  and receives a PRF key  $k$  from  $\mathcal{A}$ . It then sends  $(\text{approve}, \text{REC})$  to the trusted party on behalf of SEN (we recall that the functionality expects to receive an approval from both SEN and the ideal adversary). In case  $\mathcal{A}$  refuses to send a PRF key,  $\text{Sim}_{\text{SEN}}$  sends  $\perp$  to the trusted party, aborting the execution.
4.  $\text{Sim}_{\text{SEN}}$  outputs whatever  $\mathcal{A}$  outputs.

We first note that the simulator runs in polynomial time since all it does is emulating the role of  $\mathcal{F}_{\text{PRF}}$ . Next, we note that the receiver's privacy trivially holds in the  $\mathcal{F}_{\text{PRF}}$ -hybrid model since the sender and the receiver do not communicate any message.

**A corrupted receiver.** Finally, we consider the case when receiver REC is corrupted. Let  $\mathcal{A}$  denote an adversary controlling receiver REC, we build a simulator  $\text{Sim}_{\text{REC}}$  that generates its view. Simulator  $\text{Sim}_{\text{REC}}$  needs to emulate the roles of the sender SEN and of the server S using the leakage it gets from the trusted party. Below we describe the simulator  $\text{Sim}_{\text{REC}}$  that is given access to adversary  $\mathcal{A}$ .

**Convention:** During the simulation  $\text{Sim}_{\text{REC}}$  evaluates queries to the random oracle  $\mathcal{H}$ . Such queries are made by  $\mathcal{A}$  or by  $\text{Sim}_{\text{REC}}$  during its simulation of the corrupted REC. To evaluate  $\mathcal{H}(x)$ ,  $\text{Sim}_{\text{REC}}$  first checks if it has already recorded a pair  $(x, r)$ , in which case  $\mathcal{H}(x)$  evaluates to the value  $r$ . Otherwise,  $\text{Sim}_{\text{REC}}$  chooses a random string  $r \in \mathbb{Z}_M$ , records  $(x, r)$  and evaluates  $\mathcal{H}(x)$  to  $r$ .

1. Upon receiving auxiliary input  $z$ ,  $\text{Sim}_{\text{REC}}$  invokes  $\mathcal{A}$  on this input.
2. (preprocess,  $n, m$ ) messages from the trusted party that denote that the honest SEN submitted a text of length  $n$  are ignored.
3. Whenever  $\mathcal{A}$  initializes protocol  $\pi_{\text{Query}}$  to learn the trapdoor corresponding to a given pattern,  $\text{Sim}_{\text{SEN}}$  emulates the role of  $\mathcal{F}_{\text{PRF}}$  and receives an input  $p$  from  $\mathcal{A}$ . It then generates a PRF key  $k$  and forwards  $F_k(p)$  to the adversary. Finally,  $\text{Sim}_{\text{REC}}$  sends (query,  $p$ ) to the trusted party.
4. Upon receiving  $\perp$  from the ideal functionality,  $\text{Sim}_{\text{REC}}$  sends  $\perp$  to  $\mathcal{A}$  and abort. Otherwise, it sends (approve, REC) to the ideal functionality.
5. Upon receiving (query,  $p, (i_1, \dots, i_t), \text{id}$ ) from the trusted party,  $\text{Sim}_{\text{REC}}$  sends  $\{i_j\}_{j \in [t]}$  to  $\mathcal{A}$  and outputs whatever  $\mathcal{A}$  does.

Notice that the simulator runs in polynomial-time since it only emulates the role of  $\mathcal{F}_{\text{PRF}}$ . Next, we claim that the privacy of sender SEN is guaranteed and that REC does not gain any additional information during the execution of the real protocol  $\pi_{\text{SH}}$ . In fact, the hybrid and simulated views are statistically close where the only potential difference is due to an incorrect response by the server in the real execution. Since correctness is obtained with overwhelming probability (as formally stated above), the receiver's view is statistically indistinguishable in both executions.

**Collusion.** So far we considered single corruption cases. However, the proof carries over also in the presence of collusion between the receiver REC and S. In this case, we need to show that the server and the receiver cannot conclude any additional information about the text other than what is obtained by the leakage from the queries. A proof of this statement follows the same argument as in the single corruption case. Note in particular that the only additional information given to the server are the values of the patterns, and it still remains infeasible to guess new trapdoors that enable to obtain more information about the outsourced text.

■

## 4.1 Efficiency

We start by considering the efficiency of our scheme when we do not use the packaging approach. Notice that in this case we do not need the random oracle, but as shown below we have strong limitations on the size of the text. Namely, without packaging the server S is asked to solve subset sum instances of the form

$(\tilde{\mathbf{T}}, R_p)$ , where  $\tilde{\mathbf{T}}$  is a vector of length  $\ell = n - m + 1$  with elements from  $\mathbb{Z}_M$  for some integer  $M$ . To achieve correctness, we require that each subset sum instance has a unique solution with high probability (cf. also the proof of Theorem 1). In order to satisfy this property, one needs to set the parameters in such a way that the quantity  $2^\ell/M$  is negligible. Writing  $M = 2^{\kappa+\ell}$ , we achieve a reasonable correctness level with, e.g., security parameter  $\kappa \geq 80$ . On the other hand, to let  $S$  solve subset sum instances efficiently, we need to consider *low-density* subset sum instances. The analysis of Section 2.2 (see in particular Eq. (2)) yields, in this case, a value of  $\ell \approx \sqrt{\kappa}$ . This poses an apparently inherent limitation on the length of the text to be preprocessed. For instance, even using a higher value  $\kappa \approx 10^4$  (yielding approximately subset sum elements of size 10 KByte) limits the length of the text to only 100 bits.

To overcome this limitations, we can use the packaging approach from our protocol above. Namely, when we structure the text into blocks of length  $2m$  bits, the preprocessed blocks  $\tilde{B}_b$  consist of  $\ell = m + 1$  elements in  $\mathbb{Z}_M$ . As above we can set  $M = 2^{\kappa+\ell}$  to guarantee correctness. For efficiency, we have, however, the advantage that the blocks are reasonably short which yields subset sum instances of the form  $(\tilde{B}_b, R_p)$  that can be solved in polynomial-time. By combining many blocks we can support texts of *any length* polynomial in the security parameter. We further note that for sufficiently small lengths of  $m$  (which are typically some constant), a brute force search in time  $2^m$  per package is sufficient in order to solve the subset sum problem. Finally, we emphasize that the communication/computational complexities of  $\pi_{\text{Query}}$  depend on the underlying oblivious PRF evaluation. This in particular only depends on  $m$  (due to the complexity of the current implementation of the [NR97] PRF). Using improved PRFs can further reduce the communication complexity. Whereas the communication complexity of  $\pi_{\text{OpM}}$  solely depends on the number of matches of  $p$  in  $T$  which is essentially optimal.

## 5 Security in the Presence of Malicious Adversaries

In this section we explain how to modify the semi-honest construction from Section 4 to obtain security against malicious adversaries (that is, corrupted server and receiver). For simplicity, we will consider progressive modifications of protocol  $\pi_{\text{SH}}$ , where each modification deals with a different corruption case. These extensions can then be combined to obtain a construction which supports full malicious security. In order to maintain the presentation of our protocols simple, we will present the protocol in the presence of corrupted server in its basic form without relying on the packaging technique described in Section 4. We stress though that our construction can easily handle packaging as well.

### 5.1 Dealing with a Malicious Server

The underlying idea here is to add an efficient mechanism in protocol  $\pi_{\text{SH}}$  that enables to verify the correctness of the server's answers. Notably, this already provides security against a malicious server, because the server does not have any input/output with respect to protocol  $\pi_{\text{Query}}$ , and we already ensured privacy within the semi-honest proof. To do so, we will rely on Merkle trees [Mer89] commitment schemes (see Section 2.5), that essentially produce a succinct commitment that is independent of the length of the committed message. We modify protocol  $\pi_{\text{Pre}}$  by instructing sender  $\text{SEN}$  to commit to its preprocessed text using Merkle trees. Precisely, let  $\tilde{\mathbf{T}}$  be the preprocessed text outsourced to server  $S$ , as described in the protocol of Figure 3. Sender  $\text{SEN}$  generates a binary tree building on top of leaves  $\{\tilde{\mathbf{T}}[i] \parallel i\}_{i=1}^\ell$ . Then, for every pattern  $p$  such that its trapdoor  $\mathcal{H}(F(k, p))$  corresponds to a subset in  $\tilde{\mathbf{T}}$  with locations  $\{i_1, \dots, i_t\}$ , the server is asked to decommit the paths from the root to the leaves corresponding to these locations. To verify such values, receiver  $\text{REC}$  recomputes the paths all the way back to the root. We emphasize that this approach already ensures that the server cannot return a proper subset of the actual set of matches, since this

would imply that the server has found two different solutions for a given target, which we know occurs only with a negligible probability.

Nevertheless, Merkle trees do not protect against a malicious server declaring that there is no match even though a given pattern is actually matched. To prevent this attack we use universal accumulators (see Section 2.6), so that the server has to provide a witness for non-membership of an element in a set. We remark that applying this technique in a naïve way would not work, since the potential number of elements in the set is exponential in  $n$  (i.e., counting all possible subsets from the preprocessed text) resulting in exponential running time for SEN and S. Instead, we let SEN commit to the set of trapdoors it generated while creating the preprocessed text  $\tilde{\mathbf{T}}$ . This yields a much smaller group  $G$  with at most  $\ell$  elements. (For privacy issues we need to pad  $G$  in such a way that S cannot detect the actual number of trapdoors, which would also reveal the number of distinct substrings of length  $m$  in the text.) During the preprocessing the sender computes an accumulator for this set.

Our final construction is slightly different than this, in that we need to provide REC with the proper information allowing it to verify the values retrieved from the server. First, REC needs to know the commitment to  $\tilde{\mathbf{T}}$ , i.e., the root  $h$  of the Merkle tree, and the accumulator to the set  $G$  (which is denoted  $h_G$ ). These values will be sent from SEN within protocol  $\pi_{\text{Query}}$ . Moreover, in order to verify validity of a witness for (non-)membership, REC also needs the commitment of the trapdoor for pattern  $p$ , denoted by  $\gamma_p$ . Clearly, this value can neither be forwarded by the server (since otherwise it could easily cheat), nor by sender SEN (since otherwise it should keep a state of size equal to the number of trapdoors). Our solution is to split the output of the PRF into two parts  $(R'_p, r_p)$ : The first part is used as input to the random oracle  $\mathcal{H}$  to evaluate the actual trapdoor  $R_p$ ; whereas the second part is used as randomness in the computation of the commitment  $\gamma_p$  corresponding to  $R'_p$  (and thus to  $R_p$ ).<sup>6</sup> Given this randomness, REC can compute the commitment  $\gamma_p$  and thus verify the accumulator's proof.

To sum up, we rely on the following building blocks: (1) a succinct commitment scheme ( $\text{Commit}_M, \text{Open}_M$ ) (implemented via Merkle trees; see Section 2.5); (2) a perfectly binding (and computationally hiding) commitment scheme ( $\text{Commit}, \text{Open}$ ) (see Section 2.4); (3) a static universal accumulator (see Section 2.6). A detailed description of the protocol can be found in Figure 4.

**Theorem 2** *Let  $\kappa \in \mathbb{N}$  be the security parameter. For integers  $n, m, \lambda, \mu, \mu'$  and  $\ell = n - m + 1$ , let  $M = 2^{\ell + \kappa}$ ,  $\mu = \text{poly}(\kappa)$ ,  $\mu' = \text{poly}(\kappa)$ ,  $\lambda = \text{poly}(\kappa)$  and fix  $\kappa$  such that  $2^\ell/M$  is negligible (in  $\kappa$ ). Assume that  $\mathcal{H}$  is a random oracle, and that  $(\text{Commit}_M, \text{Open}_M)$  is binding,  $(\text{Commit}, \text{Open})$  is computationally hiding, and  $(\text{Init}, \text{CreateAcc}, \text{CreateWit}, \text{Check})$  satisfies collision freeness. Then, the round optimal protocol  $\pi_{\text{MalS}}$  from Figure 4 securely implements  $\mathcal{F}_{\text{OPM}}$  in the presence of a malicious server S.*

**Proof:** Proving security follows in two steps. We first prove that protocol  $\pi_{\text{MalS}}$  is correct. Namely, with all but negligible probability the malicious server does not return a false search response. Second, we claim that condition on that  $\pi_{\text{MalS}}$  is correct, we can simulate the view of the server.

**Proof of correctness.** We say that our construction achieves correctness if with overwhelming probability each pattern query  $p$  issued by REC is answered correctly with respect to the outsourced text  $T$ . More concretely, for each pattern  $p \in \{0, 1\}^m$  and a text  $T \in \{0, 1\}^n$ , let  $\{i_j\}_{j \in [t]}$  be the positions in  $T$  where  $p$  matches the text. Then, our protocol is correct if it returns for any such  $p$  and  $T$  the same matches  $\{i_j\}_{j \in [t]}$ . As in Theorem 1, we neglect the event that a solution of subset sum is non-unique due to a collision in the PRF and/or random oracle (since the probability that this happens is negligible in the security parameter  $\kappa$ ).

<sup>6</sup>In the basic form of the protocol it does not make any difference whether we commit to  $R_p$  or to  $R'_p$ . However, in the extension based on packaging committing to  $R'_p$  yields a better communication complexity in the setup phase.

**Protocol**  $\pi_{\text{MalS}} = (\pi_{\text{Pre}}, \pi_{\text{Query}}, \pi_{\text{Opm}})$

Let  $\kappa \in \mathbb{N}$  be the security parameter and let  $M, \lambda, m, n, \mu, \mu'$  be integers. Further, let  $\mathcal{H} : \{0, 1\}^\mu \rightarrow \mathbb{Z}_M$  be a random oracle and  $F : \{0, 1\}^\kappa \times \{0, 1\}^m \rightarrow \{0, 1\}^{\mu'}$  be a PRF. Consider a succinct commitment scheme  $(\text{Commit}_M, \text{Open}_M)$ , a computationally hiding commitment scheme  $(\text{Commit}, \text{Open})$  and a static universal accumulator  $(\text{Init}, \text{CreateAcc}, \text{CreateWit}, \text{Check})$ . Protocol  $\pi_{\text{MalS}}$  involves a sender SEN holding a text  $T \in \{0, 1\}^n$ , a receiver REC querying patterns  $p \in \{0, 1\}^m$  and a server S. The interaction is specified below.

**Setup phase,  $\pi_{\text{Pre}}$ .** The protocol is invoked between sender SEN and server S.

1. *Preprocessing the text.* Given input  $T$  and integer  $m$ , sender SEN defines a vector  $\tilde{\mathbf{T}}$  of length  $\ell = n - m + 1$  and picks a random PRF key  $k \in \{0, 1\}^\kappa$ . For any  $p \in \{0, 1\}^m$  denote by  $\{i_j\}_{j \in [t]}$  the set of indexes corresponding to the positions where  $p$  occurs in  $T$ . Sender SEN evaluates  $(R'_p, r_p) = F(k, p)$  (for  $R'_p \in \{0, 1\}^\mu$ ), computes  $R_p = \mathcal{H}(R'_p)$ , samples  $a_1, \dots, a_{t-1} \in \mathbb{Z}_M$  at random and then fixes  $a_t$  such that  $a_t = R_p - \sum_{j=1}^{t-1} a_j \bmod M$ . It then sets  $\tilde{\mathbf{T}}[i_j] = a_j$  for all  $j \in [t]$ .
2. *Committing to the text.* In addition, SEN commits to  $\tilde{\mathbf{T}}$  by letting  $h = \text{Commit}_M(\tilde{\mathbf{T}}[1] \parallel (1), \dots, (\tilde{\mathbf{T}}[\ell] \parallel (\ell)))$  (for  $\{\tilde{\mathbf{T}}[i] \parallel (i)\}_i$  the committed values).
3. *Committing to trapdoors.* Denote by  $p_1, \dots, p_\nu$  all *distinct* patterns of length  $m$  in  $T$ , for some  $\nu \leq \ell$ . Then, SEN computes  $\gamma_i = \text{Commit}(R'_{p_i}; r_{p_i})$  for all  $i = 1, \dots, \nu$  and sets  $G := \{\gamma_1, \dots, \gamma_\nu, \gamma_{\nu+1}, \dots, \gamma_\ell\}$  for randomly chosen  $\gamma_{\nu+1}, \dots, \gamma_\ell$ .
4. *Setup for the accumulator.* Finally, SEN samples  $(\text{pk}_{\text{acc}}, \text{sk}_{\text{acc}}) \leftarrow \text{Init}(1^\kappa)$ , computes  $(h_G, \text{aux}) \leftarrow \text{CreateAcc}(\text{pk}_{\text{acc}}, \text{sk}_{\text{acc}}, G)$ , outsources  $(\tilde{\mathbf{T}}, h, G, \text{pk}_{\text{acc}}, h_G, \text{aux})$  to S and keeps  $(k, h, \text{pk}_{\text{acc}}, h_G)$ .

**Query phase,  $\pi_{\text{Query}}$ .** Upon issuing a query  $p \in \{0, 1\}^m$  by receiver REC, client SEN and REC engage in an execution of protocol  $\pi_{\text{Query}}$  which implements the oblivious PRF functionality  $(k, p) \mapsto (\perp, F(k, p))$ , such that  $(R'_p, r_p) = F(k, p)$ . Receiver REC also receives from SEN the commitments  $h$  and  $h_G$ , and the value  $\text{pk}_{\text{acc}}$ . Upon completion, REC computes trapdoor  $R_p = \mathcal{H}(R'_p)$ .

**Oblivious pattern matching phase,  $\pi_{\text{Opm}}$ .** This protocol is engaged between server S that inputs  $(\tilde{\mathbf{T}}, h, G, \text{pk}_{\text{acc}}, h_G, \text{aux})$  and receiver REC that inputs  $(h, \text{pk}_{\text{acc}}, h_G, R_p, R'_p, r_p)$ .

1. *Solving subset sum.* Upon receiving  $(R'_p, \gamma_p)$  from REC, where  $\gamma_p = \text{Commit}(R'_p; r_p)$ , the server views  $(\tilde{\mathbf{T}}, R_p = \mathcal{H}(R'_p))$  as a subset sum instance and solves this instance. Let  $\mathbf{s}$  denote the solution, and denote with  $\{i_j\}_{j \in [t]}$  the set of indexes such that  $\mathbf{s}[i_j] = 1$ .
2. *Opening the Merkle tree.* The server runs  $(\tilde{\mathbf{T}}[i_j], \text{path}(i_j)) = \text{Open}_M(h, i_j)$  for all  $j \in [t]$ .
3. *Proving (non-)membership.* The server computes a witness for (non-)membership  $\zeta_{\gamma_p} \leftarrow \text{CreateWit}(\text{type}, \text{pk}_{\text{acc}}, G, h_G, \text{aux}, \gamma_p)$  (where  $\text{type} = 0$  if  $\gamma_p \in G$  and  $\text{type} = 1$  otherwise). Hence, S returns REC the set  $\{i_j\}_{j \in [t]}$  together with  $\{\tilde{\mathbf{T}}[i_j], \text{path}(i_j)\}_{j \in [t]}$  and  $\zeta_{\gamma_p}$ .

Receiver REC verifies the openings  $\{\tilde{\mathbf{T}}[i_j], \text{path}(i_j)\}_{j \in [t]}$ , checks that  $\sum_{j=1}^t \tilde{\mathbf{T}}[i_j] = R_p$  and that  $(\text{Check}(0, \text{pk}_{\text{acc}}, h_G, \zeta_{\gamma_p}, \gamma_p)) \vee (\text{Check}(1, \text{pk}_{\text{acc}}, h_G, \zeta_{\gamma_p}, \gamma_p)) = 1$ .

Figure 4: Outsourced pattern matching resisting a malicious server

For each pattern  $p$  we consider two subcases: the case that  $p$  matches the text in at least one position, and the case that  $p$  does not match the text. In both cases we reduce correctness to the collision intractability of Merkle trees and the collision freeness of the universal accumulator.

**Claim 4** For any pattern  $p^* \in \{0, 1\}^m$  such that  $p^*$  matches  $T$  in at least one position, protocol  $\pi_{\text{Query}}$  is correct assuming the binding property of Merkle trees and collision freeness of the accumulator.

**Proof:** Consider a pattern  $p^*$  and assume that it matches the text in positions  $\tau = \{i_j\}_{j \in [t]}$ . We consider two bad events and bound their probability.

- $\text{bad}_1$ : Occurs in the real execution whenever S convinces REC that  $p^*$  does not match the text.
- $\text{bad}_2$ : Occurs in the real execution whenever S convinces REC that the solution to the subset sum instance  $(\tilde{\mathbf{T}}, R_{p^*})$  is  $\tau' = \{i'_j\}_{j \in [t]}$  for  $\{i'_j\}_j \neq \{i_j\}_j$ .

Intuitively,  $\text{bad}_1$  happens with negligible probability since provoking this event means that the server is able to prove that the commitment to a trapdoor  $R_p$  does not belong to the set  $G$  of all commitments. This violates collision freeness of the accumulator.

We start by considering a malicious  $\mathcal{A}$  controlling S in the real world and provoking event  $\text{bad}_1$  for text  $T$  and pattern  $p^*$ . From such an adversary, we build an efficient machine  $\mathcal{A}$  breaking collision freeness of  $(\text{Init}, \text{CreateAcc}, \text{CreateWit}, \text{Check})$ , as follows. Adversary  $\mathcal{A}$  is given  $\text{pk}_{\text{acc}}$ , and processes the text  $T$  as described in the protocol of Figure 4; in particular, it computes  $(\tilde{\mathbf{T}}, G, h)$  and forwards  $G$  to its own oracle  $\mathcal{O}$  obtaining  $(h_G, \text{aux})$  such that  $(h_G, \text{aux}) \leftarrow \text{CreateAcc}(\text{pk}_{\text{acc}}, \text{sk}_{\text{acc}}, G)$ . The values  $(\tilde{\mathbf{T}}, h, \text{pk}_{\text{acc}}, G, h_G, \text{aux})$  are forwarded to S. Next,  $\mathcal{A}$  simulates a query for pattern  $p^*$  (matching  $T$  by hypothesis), as described in protocol  $\pi_{\text{MalS}}$ ; this yields a pair  $(R'_{p^*}, r_{p^*})$  which is forwarded to the server. S returns the matching positions  $\{i_j\}_{j \in [t]}$  together with a witness  $\zeta_{p^*}$ . Finally,  $\mathcal{A}$  outputs  $(G, \zeta_{p^*}, \gamma_{p^*})$ .

Note that the above simulation is perfect from the point of view of S; in particular, since we are assuming S provokes event  $\text{bad}_1$  with non-negligible probability, we can conclude that  $\text{Check}(1, \text{pk}_{\text{acc}}, h_G, \zeta_{p^*}, \gamma_{p^*}) = 1$  which contradicts collision freeness as  $\gamma_{p^*} \in G$ .

Now, consider the second event  $\text{bad}_2$ . Intuitively,  $\text{bad}_2$  happens with a negligible probability since provoking this event means that the server is able to decommit at least one position in the text to a non-consistent value. This violates the binding property of Merkle trees. Recall that, with overwhelming probability, there exists at most one solution to the subset sum instance  $(\tilde{\mathbf{T}}, R_{p^*})$ . Denote this “non-collision event” by  $\text{nocol}$ . We can write:

$$\Pr[\text{bad}_2] = \Pr[\text{bad}_2 | \text{nocol}] \cdot \Pr[\text{nocol}] + \Pr[\text{bad}_2 | \overline{\text{nocol}}] \cdot \Pr[\overline{\text{nocol}}] \leq \Pr[\text{bad}_2 | \text{nocol}] + \Pr[\overline{\text{nocol}}].$$

Due to the low collision probability, we have  $\Pr[\overline{\text{nocol}}] \leq \text{negl}(\kappa)$ . It remains to prove that  $\Pr[\text{bad}_2 | \text{nocol}] \leq \text{negl}(\kappa)$ . However, conditioned on  $\text{nocol}$ , the only way to provoke event  $\text{bad}_2$  is by violating the binding property of Merkle trees: There exists an index  $i_j \in \tau'$  such that  $\text{Open}_M(h, i_j) \neq \tilde{\mathbf{T}}[i'_j]$ . In particular, an adversary provoking this event can be turned into a concrete polynomial-time machine breaking collision resistance of  $(\text{Commit}_M, \text{Open}_M)$ . Thus,  $\Pr[\text{bad}_2 | \text{nocol}]$  must also be negligible. Together with the previous equation, this concludes the proof of the claim. ■

**Claim 5** For any pattern  $p^* \in \{0, 1\}^m$  such that  $p^*$  does not match  $T$ , protocol  $\pi_{\text{Query}}$  is correct assuming the binding property of Merkle trees and collision freeness of the accumulator.

This proof follows similarly to the proof of Claim 4 and is therefore omitted. The only difference is that, in the reduction to collision freeness, adversary  $\mathcal{A}$  is able to produce a witness of membership for a value  $\gamma_{p^*} \notin G$ .



**Privacy.** Finally, we note that conditioned on correctness, simulating the server's view reduces almost immediately to the semi-honest case; so essentially the same simulator as in the proof of Theorem 1 will do. We only need to specify how the simulator can simulate the additional messages which are included in protocol  $\pi_{\text{MalS}}$  but not in  $\pi_{\text{SH}}$ . Specifically, during the setup phase Sim needs also to produce the values  $(h, h_G, aux)$ . To do so, the simulator will simply define  $G$  by committing to  $\ell$  random values and then compute  $(h_G, aux)$  as in a real execution. Note that the value  $h$  can be computed consistently with the sampled pre-processed text  $\tilde{\mathbf{T}}$ . Finally, during the oblivious pattern matching phase, Sim will first compute the value  $R_p$  as described in the proof of Theorem 1. If the pattern matches, then  $\gamma_p$  is chosen to be a random element in  $G$ ; otherwise a fresh commitment of a random message is used. It follows from the computationally hiding property of (Commit, Open) that the above strategy is not detectable by the adversary but with a negligible probability.

A formal description of the simulator  $\text{Sim}_S$  (with access to  $\mathcal{A}$ ) follows.

**Convention:** During the simulation  $\text{Sim}_S$  evaluates queries to the random oracle  $\mathcal{H}$ . Such queries are made by  $\mathcal{A}$  or by  $\text{Sim}_S$  during its simulation of the corrupted server. To evaluate  $\mathcal{H}(x)$ ,  $\text{Sim}_S$  first checks if it has already recorded a pair  $(x, r)$ , in which case  $\mathcal{H}(x)$  evaluates to the value  $r$ . Otherwise,  $\text{Sim}_S$  chooses a random string  $r \in \mathbb{Z}_M$ , records  $(x, r)$  and evaluates  $\mathcal{H}(x)$  to  $r$ .

1. On input the auxiliary input  $z$ ,  $\text{Sim}_S$  invokes  $\mathcal{A}$  (i.e., the corrupted server) on this input. The simulator keeps track of a table  $\mathcal{B}$  that is initially set to the empty table.
2. Upon receiving a (preprocess,  $n, m$ ) message from the trusted party denoting that the honest SEN wants to outsource a text of length  $n$  to the trusted party,  $\text{Sim}_S$  defines text  $\tilde{\mathbf{T}}$  by sampling uniformly at random a vector of length  $\ell := n - m + 1$  from  $\mathbb{Z}_M^\ell$ .

Next,  $\text{Sim}_S$  computes  $h = \text{Commit}_M((\tilde{\mathbf{T}}[1] \parallel (1)), \dots, (\tilde{\mathbf{T}}[\ell] \parallel (\ell)))$  (for  $\{\tilde{\mathbf{T}}[i] \parallel (i)\}_i$  the committed values), lets  $G = \{\gamma_1, \dots, \gamma_\ell\}$  for randomly chosen  $\gamma_1, \dots, \gamma_\ell$ , and sets  $(h_G, aux) \leftarrow \text{CreateAcc}(\text{pk}_{\text{acc}}, \text{sk}_{\text{acc}}, G)$  where  $(\text{pk}_{\text{acc}}, \text{sk}_{\text{acc}}) \leftarrow \text{Init}(1^\kappa)$ . Hence, it forwards  $(\tilde{\mathbf{T}}, h, \text{pk}_{\text{acc}}, h_G, aux)$  to adversary  $\mathcal{A}$  and stores  $(\tilde{\mathbf{T}}, G)$  for later usage as well.

3. Upon receiving a (query, REC,  $(i_1, \dots, i_t)$ , id) message from the trusted party indicating that receiver REC submitted a search query that was approved by SEN,  $\text{Sim}_S$  distinguishes two cases.
  - (a) *Pattern queried by REC appears in  $T$ :* In this case  $\{i_j\}_{j \in [t]}$  is not the empty set, and the simulator samples uniformly at random a value  $X_{\text{id}}$  from  $\{0, 1\}^\mu$ . Then, it programs the random oracle  $\mathcal{H}(X_{\text{id}})$  to  $\sum_{j=1}^t \tilde{\mathbf{T}}[i_j]$ ; if  $\mathcal{H}$  has already been programmed to a different value, then we abort. Furthermore,  $\text{Sim}_S$  picks an element  $\gamma \leftarrow G$ .
  - (b) *Pattern does not appear in  $T$ :* In this case  $\{i_j\}_{j \in [t]}$  is the empty set, and we check if table  $\mathcal{B}$  contains a value of the form  $(\text{id}, X_{\text{id}})$ . Otherwise, we pick the value  $X_{\text{id}}$  uniformly at random in  $\{0, 1\}^\mu$  and store  $(\text{id}, X_{\text{id}})$  in  $\mathcal{B}$ . Furthermore,  $\text{Sim}_S$  picks an element  $\gamma$  by committing to a random message.

Finally,  $\text{Sim}_S$  (emulating the role of REC in the real execution) forwards  $(X_{\text{id}}, \gamma)$  to the adversary.

4. If  $\mathcal{A}$  does not answer with  $(i_1, \dots, i_t)$ ,  $\text{Sim}_S$  sends  $\perp$  to ideal functionality and abort. Otherwise it sends the trusted party (approve, REC).
5.  $\text{Sim}_S$  outputs whatever  $\mathcal{A}$  does.

We first note that  $\text{Sim}_S$  runs in polynomial time since it only samples a random vector from  $\mathbb{Z}_M^\ell$ , and then calculates the sum of values from a given subset. Next, we show that the distribution produced by  $\text{Sim}_S$  in

the ideal world is computationally indistinguishable from the distribution that  $\mathcal{A}$  expects to see in the real world. This is required to hold even given the leakage revealing the positions where the pattern matches the text. We start by defining a hybrid distribution  $\mathbf{HYB}_{\pi_{\text{MalS}}, \mathcal{A}(z)}^1(\kappa, (-, T, (p_1, \dots, p_\lambda)))$  that is defined as the real experiment  $\mathbf{REAL}_{\pi_{\text{MalS}}, \mathcal{A}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))$  with the difference that  $X_{\text{id}}$  is not computed by a PRF but rather by a random function  $f_\kappa$ . We have the following claim:

**Claim 6** *Let  $F$  be a secure pseudorandom function (cf. Definition 1), then for all sufficiently large  $\kappa \in \mathbb{N}$  and for any tuple of inputs  $(T, (p_1, \dots, p_\lambda))$  and auxiliary input  $z$ , it holds that*

$$\{\mathbf{REAL}_{\pi_{\text{MalS}}, \mathcal{A}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}} \stackrel{c}{\approx} \{\mathbf{HYB}_{\pi_{\text{MalS}}, \mathcal{A}(z)}^1(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}}.$$

The proof of the above claim is similar to the proof of Claim 2, and consists of a simple reduction to the pseudorandomness property of the PRF.

Next we consider a hybrid distribution  $\mathbf{HYB}_{\pi_{\text{MalS}}, \mathcal{A}(z)}^2(\kappa, (-, T, (p_1, \dots, p_\lambda)))$  that is defined as the previous hybrid, with the difference that the elements in the set  $G$  are replaced by commitments to random messages. A standard hybrid argument yields the following claim:

**Claim 7** *Let (Commit, Open) be computationally hiding, then for all sufficiently large  $\kappa \in \mathbb{N}$  and for any tuple of inputs  $(T, (p_1, \dots, p_\lambda))$  and auxiliary input  $z$ , it holds that*

$$\{\mathbf{HYB}_{\pi_{\text{MalS}}, \mathcal{A}(z)}^1(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}} \stackrel{c}{\approx} \{\mathbf{HYB}_{\pi_{\text{MalS}}, \mathcal{A}(z)}^2(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}}.$$

Finally, we need to bound the distance between the experiment  $\mathbf{HYB}_{\pi_{\text{MalS}}, \mathcal{A}(z)}^2(\kappa, (-, T, (p_1, \dots, p_\lambda)))$  and the simulated view. We do this in the next claim.

**Claim 8** *For any input text  $T$ , patterns  $p_1, \dots, p_\lambda$ , and auxiliary input  $z$ , it holds that*

$$\{\mathbf{IDEAL}_{\mathcal{F}_{\text{OPM}}, \text{Sim}(z)}(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}} \equiv_s \{\mathbf{HYB}_{\pi_{\text{MalS}}, \mathcal{A}(z)}^2(\kappa, (-, T, (p_1, \dots, p_\lambda)))\}_{\kappa \in \mathbb{N}}.$$

**Proof:** Define the following event  $\text{bad}_3$  that occurs when the simulator aborts in the ideal world.

- Event  $\text{bad}_3$ : Occurs if the simulation given above is aborted. In this case the corrupted server has asked for a direct query to the random oracle of the form  $(f_\kappa(p)|_\mu)$  (where  $y|_\mu$  denotes the truncation of bit-string  $y$  to the first  $\mu$  bits) *before* it has seen  $f_\kappa(p)$ , where  $p$  is a pattern that occurs in the text  $T$ . Recall that  $f_\kappa$  is a random function as defined in the first hybrid world.

An argument similar to the proof of Theorem 1 shows that the simulated view and the output distribution in  $\mathbf{HYB}_{\pi_{\text{MalS}}, \mathcal{A}(z)}^2(\kappa, (-, T, (p_1, \dots, p_\lambda)))$  are identical conditioned on the event  $\text{bad} := \text{bad}_1 \wedge \text{bad}_2 \wedge \text{bad}_3$  not happening.

On the other hand, since  $f_\kappa$  is a random function, it is easy to see that  $\text{bad}_3$  occurs with probability at most  $\text{poly}(\kappa)/2^\mu = \text{negl}(\kappa)$ . Now, Claim 4 and Claim 5 together with a union bound imply that the probability of  $\text{bad}$  is negligible. This concludes the proof. ■

Putting together Claim 6, Claim 7 and Claim 8 concludes the proof. ■

## 5.2 Dealing with a Malicious Receiver

We note that the protocol of Figure 3 is already secure against a malicious REC (with a small change). In fact, the only inputs receiver REC provides is an execution of protocol  $\pi_{\text{SH}}$  (sub-protocol  $\pi_{\text{Query}}$ ) are the search queries  $p_i$ 's. To this end, we claim that the only way REC can attack the protocol is by guessing a trapdoor of which it did not submit a query for. We further claim that the probability of this event is

negligible in the security parameter (where this event occurs with negligible probability even conditioned on the fact that the receiver knows the Merkle commitment and the accumulator, as these preserve the privacy of the text as proven above). Loosely speaking, this follows from the security of the PRF. Namely, given that this event occurs with a non-negligible probability a distinguisher can detect this trapdoor in the queries submitted to the random oracle. The reduction to the pseudorandomness of the PRF follows the same outlines as the semi-honest reduction for the case that the server is corrupted and is therefore omitted. Hence, it suffices to ensure that protocol  $\pi_{\text{Query}}$  is secure in the presence of malicious adversaries so that a simulator can extract these queries. (An example is the protocol of [HL10], taken from [FIPR05], which implements the Naor-Riengold function [NR97] in the presence of malicious adversaries.)

## References

- [AIK10] Benny Applebaum, Yuval Ishai, and Eyal Kushilevitz. From secrecy to soundness: Efficient verification via secure computation. In *ICALP (I)*, pages 152–163, 2010.
- [AJLA<sup>+</sup>12] Gilad Asharov, Abhishek Jain, Adriana López-Alt, Eran Tromer, Vinod Vaikuntanathan, and Daniel Wichs. Multiparty computation with low communication, computation and interaction via threshold fhe. In *EUROCRYPT*, pages 483–501, 2012.
- [ATSM09] Man Ho Au, Patrick P. Tsang, Willy Susilo, and Yi Mu. Dynamic universal accumulators for DDH groups and their application to attribute-based anonymous credential systems. In *CT-RSA*, pages 295–308, 2009.
- [BGV11] Siavosh Benabbas, Rosario Gennaro, and Yevgeniy Vahlis. Verifiable delegation of computation over large datasets. In *CRYPTO*, pages 111–131, 2011.
- [BLL00] Ahto Buldas, Peeter Laud, and Helger Lipmaa. Accountable certificate management using undeniable attestations. In *CCS*, pages 9–17, 2000.
- [BLL02] Ahto Buldas, Peeter Laud, and Helger Lipmaa. Eliminating counterevidence with applications to accountable certificate management. *Journal of Computer Security*, 10(3):273–296, 2002.
- [BM77] Robert S. Boyer and J. Strother Moore. A fast string searching algorithm. *Commun. ACM*, 20(10):762–772, 1977.
- [Can00] Ran Canetti. Security and composition of multi-party cryptographic protocols. *Journal of Cryptology*, 13:143–202, 2000.
- [CF13] Dario Catalano and Dario Fiore. Vector commitments and their applications. In *PKC*, pages 55–72, 2013.
- [CFG89] Mark Chaimovich, Gregory Freiman, and Zvi Galil. Solving dense subset-sum problems by using analytical number theory. *J. Complexity*, 5(3):271–282, 1989.
- [CGKO11] Reza Curtmola, Juan A. Garay, Seny Kamara, and Rafail Ostrovsky. Searchable symmetric encryption: Improved definitions and efficient constructions. *Journal of Computer Security*, 19(5):895–934, 2011.
- [CHKO12] Philippe Camacho, Alejandro Hevia, Marcos A. Kiwi, and Roberto Opazo. Strong accumulators from collision-resistant hashing. *Int. J. Inf. Sec.*, 11(5):349–363, 2012.
- [CJL<sup>+</sup>92] Matthijs J. Coster, Antoine Joux, Brian A. LaMacchia, Andrew M. Odlyzko, Claus-Peter Schnorr, and Jacques Stern. Improved low-density subset sum algorithms. *Computational Complexity*, 2:111–128, 1992.
- [CKKC13] Seung Geol Choi, Jonathan Katz, Ranjit Kumaresan, and Carlos Cid. Multi-client non-interactive verifiable computation. In *TCC*, pages 499–518, 2013.
- [CKV10] Kai-Min Chung, Yael Tauman Kalai, and Salil P. Vadhan. Improved delegation of computation using fully homomorphic encryption. In *CRYPTO*, pages 483–501, 2010.

- [CN11] Yuanmi Chen and Phong Q. Nguyen. Bkz 2.0: Better lattice security estimates. In *ASIACRYPT*, pages 1–20, 2011.
- [CS14] Melissa Chase and Emily Shen. Pattern matching encryption. *IACR Cryptology ePrint Archive*, 2014:638, 2014.
- [DHS15] David Derler, Christian Hanser, and Daniel Slamanig. Revisiting cryptographic accumulators, additional properties and relations to other primitives. *IACR Cryptology ePrint Archive*, 2015:87, 2015.
- [DT08] Ivan Damgård and Nikos Triandopoulos. Supporting non-membership proofs with bilinear-map accumulators. *IACR Cryptology ePrint Archive*, 2008:538, 2008.
- [FGP14] Dario Fiore, Rosario Gennaro, and Valerio Pastro. Efficiently verifiable computation on encrypted data. In *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security, Scottsdale, AZ, USA, November 3-7, 2014*, pages 844–855, 2014.
- [FIPR05] Michael J. Freedman, Yuval Ishai, Benny Pinkas, and Omer Reingold. Keyword search and oblivious pseudorandom functions. In *TCC*, pages 303–324, 2005.
- [FP05] Abraham Flaxman and Bartosz Przydatek. Solving medium-density subset sum problems in expected polynomial time. In *STACS*, pages 305–314, 2005.
- [Fri86] Alan M. Frieze. On the lagarias-odlyzko algorithm for the subset sum problem. *SIAM J. Comput.*, 15(2):536–539, 1986.
- [Gen09] Craig Gentry. Fully homomorphic encryption using ideal lattices. In *STOC*, pages 169–178, 2009.
- [GGP10] Rosario Gennaro, Craig Gentry, and Bryan Parno. Non-interactive verifiable computing: Outsourcing computation to untrusted workers. In *CRYPTO*, pages 465–482, 2010.
- [GHS10] Rosario Gennaro, Carmit Hazay, and Jeffrey S. Sorensen. Text search protocols with simulation based security. In *Public Key Cryptography*, pages 332–350, 2010.
- [GKL<sup>+</sup>15] S. Dov Gordon, Jonathan Katz, Feng-Hao Liu, Elaine Shi, and Hong-Sheng Zhou. Multi-client verifiable computation with stronger security guarantees. In *TCC*, pages 144–168, 2015.
- [GM91] Zvi Galil and Oded Margalit. An almost linear-time algorithm for the dense subset-sum problem. In *ICALP*, pages 719–727, 1991.
- [GN08] Nicolas Gama and Phong Q. Nguyen. Predicting lattice reduction. In *EUROCRYPT*, pages 31–51, 2008.
- [HL10] Carmit Hazay and Yehuda Lindell. Efficient protocols for set intersection and pattern matching with security against malicious and covert adversaries. *J. Cryptology*, 23(3):422–456, 2010.
- [HT10] Carmit Hazay and Tomas Toft. Computationally secure pattern matching in the presence of malicious adversaries. In *ASIACRYPT*, pages 195–212, 2010.
- [HZ14] Carmit Hazay and Hila Zarosim. The feasibility of outsourced database search in the plain model. In *Manuscript*, 2014.
- [IN96] Russell Impagliazzo and Moni Naor. Efficient cryptographic schemes provably as secure as subset sum. *J. Cryptology*, 9(4):199–216, 1996.
- [JJK<sup>+</sup>13] Stanislaw Jarecki, Charanjit S. Jutla, Hugo Krawczyk, Marcel-Catalin Rosu, and Michael Steiner. Outsourced symmetric private information retrieval. In *ACM Conference on Computer and Communications Security*, pages 875–888, 2013.
- [KM10] Jonathan Katz and Lior Malka. Secure text processing with applications to private dna matching. In *ACM Conference on Computer and Communications Security*, pages 485–492, 2010.
- [KMP77] Donald E. Knuth, James H. Jr. Morris, and Vaughan R. Pratt. Fast pattern matching in strings. *SIAM J. Comput.*, 6(2):323–350, 1977.
- [KMR11] Seny Kamara, Payman Mohassel, and Mariana Raykova. Outsourcing multi-party computation. *IACR Cryptology ePrint Archive*, 2011:272, 2011.

- [KMR12] Seny Kamara, Payman Mohassel, and Ben Riva. Salus: a system for server-aided secure function evaluation. In *ACM Conference on Computer and Communications Security*, pages 797–808, 2012.
- [KP13] Seny Kamara and Charalampos Papamanthou. Parallel and dynamic searchable symmetric encryption. In *Financial Cryptography*, pages 258–274, 2013.
- [KPR12] Seny Kamara, Charalampos Papamanthou, and Tom Roeder. Dynamic searchable symmetric encryption. In *ACM Conference on Computer and Communications Security*, pages 965–976, 2012.
- [LATV12] Adriana López-Alt, Eran Tromer, and Vinod Vaikuntanathan. On-the-fly multiparty computation on the cloud via multikey fully homomorphic encryption. In *STOC*, pages 1219–1234, 2012.
- [LLX07] Jiangtao Li, Ninghui Li, and Rui Xue. Universal accumulators with efficient nonmembership proofs. In *ACNS*, pages 253–269, 2007.
- [LO85] J. C. Lagarias and Andrew M. Odlyzko. Solving low-density subset sum problems. *J. ACM*, 32(1):229–246, 1985.
- [LPS10] Vadim Lyubashevsky, Adriana Palacio, and Gil Segev. Public-key cryptographic primitives provably as secure as subset sum. In *TCC*, pages 382–400, 2010.
- [Lyu05] Vadim Lyubashevsky. The parity problem in the presence of noise, decoding random linear codes, and the subset sum problem. In *APPROX-RANDOM*, pages 378–389, 2005.
- [Mer89] Ralph C. Merkle. A certified digital signature. In *CRYPTO*, pages 218–238, 1989.
- [Moh11] Payman Mohassel. Efficient and secure delegation of linear algebra. *IACR Cryptology ePrint Archive*, 2011:605, 2011.
- [NR97] Moni Naor and Omer Reingold. Number-theoretic constructions of efficient pseudo-random functions. In *FOCS*, pages 458–467, 1997.
- [NS06] Phong Q. Nguyen and Damien Stehlé. LLL on the average. In *ANTS*, pages 238–256, 2006.
- [PTT11] Charalampos Papamanthou, Roberto Tamassia, and Nikos Triandopoulos. Optimal verification of operations on dynamic sets. In *CRYPTO*, pages 91–110, 2011.
- [PVW08] Chris Peikert, Vinod Vaikuntanathan, and Brent Waters. A framework for efficient and composable oblivious transfer. In *CRYPTO*, pages 554–571, 2008.
- [Sha08] Andrew Shallue. An improved multi-set algorithm for the dense subset sum problem. In *ANTS*, pages 416–429, 2008.
- [TPKC07] Juan Ramón Troncoso-Pastoriza, Stefan Katzenbeisser, and Mehmet Utku Celik. Privacy preserving error resilient dna searching through oblivious automata. In *ACM Conference on Computer and Communications Security*, pages 519–528, 2007.