#### **Research Article**

Maura B. Paterson, Douglas R. Stinson\* and Jalaj Upadhyay

# Multi-prover proof of retrievability

https://doi.org/10.1515/jmc-2018-0012 Received March 17, 2018; revised July 3, 2018; accepted August 10, 2018

**Abstract:** There has been considerable recent interest in "cloud storage" wherein a user asks a server to store a large file. One issue is whether the user can verify that the server is actually storing the file, and typically a challenge-response protocol is employed to convince the user that the file is indeed being stored correctly. The security of these schemes is phrased in terms of an extractor which will recover the file given any "proving algorithm" that has a sufficiently high success probability. This forms the basis of proof-of-retrievability (PoR) systems. In this paper, we study multiple server PoR systems. We formalize security definitions for two possible scenarios: (i) A threshold of servers succeeds with high enough probability (worst case), and (ii) the average of the success probability of all the servers is above a threshold (average case). We also motivate the study of confidentiality of the outsourced message. We give MPoR schemes which are secure under both these security definitions and provide reasonable confidentiality guarantees even when there is no restriction on the computational power of the servers. We also show how classical statistical techniques previously used by us can be extended to evaluate whether the responses of the provers are accurate enough to permit successful extraction. We also look at one specific instantiation of our construction when instantiated with the unconditionally secure version of the Shacham-Waters scheme. This scheme gives reasonable security and privacy guarantee. We show that, in the multi-server setting with computationally unbounded provers, one can overcome the limitation that the verifier needs to store as much secret information as the provers.

**Keywords:** Proof of retrievability, multiple users, secret sharing

MSC 2010: 94A60

Communicated by: Spyros Magliveras

### 1 Introduction

In the recent past, there has been a lot of activity on remote storage and the associated cryptographic problem of integrity of the stored data. This question becomes even more important when there are reasons to believe that the remote servers might act maliciously, i.e., one or more servers can delete (whether accidentally or on purpose) a part of the data since there is a good chance that the data will never be accessed, and hence, the client would never find out! In order to assuage such concerns, one would prefer to have a simple auditing system that convinces the client if and only if the server has the data. Such audit protocols, called *proof-of-retrievability* (PoR) systems, were introduced by Juels and Kaliski [11], and closely related *proof-of-data-possession* (PDP) systems were introduced by Ateniese et al. [2].

Maura B. Paterson, Department of Economics, Mathematics and Statistics, Birkbeck, University of London, Malet Street, London WC1E 7HX, United Kingdom, e-mail: m.paterson@bbk.ac.uk

\*Corresponding author: Douglas R. Stinson, David R. Cheriton School of Computer Science, University of Waterloo, Waterloo, ON, N2L 3G1, Canada, e-mail: dstinson@uwaterloo.ca

Jalaj Upadhyay, Department of Computer Science, Johns Hopkins University, Baltimore, MD 21201, USA, e-mail: jalaj@jhu.edu

In a PoR protocol, a client stores a message m on a remote server and keeps only a short private *fingerprint* locally. At some later time, when the client wishes to verify the integrity of its message, it can run an audit protocol in which it acts as a verifier while the server proves that it has the client's data. The formal security of a PoR protocol is expressed in terms of an extractor – there exists an extractor with (black-box or non-blackbox) access to the proving algorithm used by the server to respond to the client's challenge, such that the extractor retrieves the original message given any adversarial server which passes the audits with a threshold probability. Apart from this security requirement, two practical requirements of any PoR system would be to have a reasonable bound on the communication cost of every audit and small storage overhead on both the client and server.

PoR systems were originally defined for the single-server setting. However, in the real world, it is highly likely that a client would store its data on more than one server. This might be due to a variety of reasons. For example, a client might wish to have a certain degree of redundancy if one or more servers fails. In this case, the client is more likely to store multiple copies of the same data. Another possible scenario could be that the client does not trust a single server with all of its data. In this case, the client might distribute the data across multiple servers. Both of these settings have been studied previously in the literature.

The first such study was initiated by Curtmola et al. [9], who considered the first of the above two cases. They addressed the problem of storing copies of a single file on multiple servers. This is an attractive solution considering the fact that replication is a fundamental principle in ensuring the availability and durability of data. Their system allows the client to audit a subset of servers even if some of them collude.

On the other hand, Bowers, Juels and Oprea [8] considered the second of the above two cases. They studied a system where the client's data is distributed and stored on different servers. This ensures that none of the servers has the whole data.

Both of these systems covered one specific instance of the wide spectrum of possibilities when more than one server is involved. For example, none of the works mentioned above addresses the question of the privacy of data. Both of them argue that, for privacy, the client can encrypt its file before storing it on the servers. These systems are secure only in the computational setting and the privacy guarantee is dependent on the underlying encryption scheme. On the other hand, there are known primitives in the setting of distributed systems, like secret sharing schemes, that are known to be unconditionally secure. Moreover, we can also utilize cross-server redundancy to get more practical systems.

#### 1.1 Our contributions

In Section 2, we give the formal description of multi-server PoR (MPoR) systems. We state the definitions for worst-case and the average-case secure MPoR systems. We also motivate the privacy requirement and state the privacy definition for MPoR systems. In Section 3, we define various primitives to the level required to understand this paper.

In Section 4, we give a construction of an MPoR scheme that achieves worst-case security when the malicious servers are computationally unbounded. Our construction is based on ramp schemes and a singleserver PoR scheme. Our construction achieves confidentiality of the message. To exemplify our scheme, we instantiate this scheme with a specific form of ramp scheme.

In Section 5, we give a construction of an MPoR scheme that achieves average-case security against computationally unbounded adversaries. For an MPoR system that affords average-case security, we also show that an extension of classical statistical techniques previously used by us [15] can be used to provide a basis for estimating whether the responses of the servers are accurate enough to allow successful extraction.

One of the benefits of an MPoR system is that it provides cross-server redundancy. In the past, this feature has been used by Bowers, Juels and Oprea [8] to propose a multi-server system called HAIL. We first note that the constructions in Section 4 and Section 5 do not provide any improvement on the storage overhead of the server or the client. In Section 6, we give a construction based on the Shacham–Waters protocol [16] that allows significant reduction of the storage overhead of the client in the multi-server setting.

#### 1.2 Related works

The concept of *proof of retrievability* is due to Juels and Kaliski [11]. A PoR scheme incorporates a challengeresponse protocol in which a verifier can check that a message is being stored correctly, along with an extractor that will actually reconstruct the message, given the algorithm of a "prover" who is able to correctly respond to a sufficiently high percentage of challenges.

There are also papers that describe the closely related (but slightly weaker) idea of a proof-of-datapossession scheme (PDP scheme), e.g., [2]. A PDP scheme permits the possibility that not all of the message blocks can be reconstructed. Ateniese et al. [2] also introduced the idea of using homomorphic authenticators to reduce the communication complexity of the system. This scheme was improved in a follow-up work by Ateniese et al. [4]. Shacham and Waters [16] later showed that the scheme of Ateniese et al. [1] can be transformed into a PoR scheme by constructing an extractor that extracts the file from the responses of the prover

Bowers, Juels and Oprea [8] extended the idea of Juels and Kaliski [11] and used error-correcting codes. The main difference in their construction is that they use the idea of an "outer" and an "inner" code (in the same vein as concatenated codes), to get a good balance between the extra storage overhead and computational overhead in responding to the audits. Dodis, Vadhan and Wichs [10] provided the first example of an unconditionally secure PoR scheme, also constructed from an error-correcting code, with extraction performed through *list decoding* in conjunction with the use of an *almost-universal hash function*. They also give different constructions depending on the computational capabilities of the server. Previously [15], we studied PoR schemes in the setting of unconditional security and showed some close connections to error-correcting codes.

Recently, Ateniese, Kamara and Katz [5] defined the framework of proof-of-storage systems to understand PDP and PoR system in a unified manner. They argue that existing PoR [16] and PDP [2] schemes can be seen as an instantiation of their framework. They used homomorphic identification schemes to give efficient proofof-storage systems in the random-oracle model. They further exhibited that existing constructions of PoR and PDP schemes are specific instantiation of their construction. Wang et al. [19] gave the first privacy preserving public auditable proof-of-storage systems. We refer the readers to the survey by Kamara and Lauter [12] regarding the architecture of proof-of-storage systems.

**Distributed cloud computing.** All the constructions mentioned above considered single server system; however, such systems are prone to failure leading to catastrophic problems [20]. However, proof-of-storage systems have been also studied in the setting where there is more than one server or more than one client. The first such setting was studied by Curtmola et al. [9]. They studied a multiple-replica PDP system, which is the natural generalization of a single-server PDP system to t servers.

Bowers, Juels and Oprea [8] introduced a distributed system that they called HAIL. Their system allows a set of provers to prove the integrity of a file stored by a client. The idea in HAIL is to exploit the cross-prover redundancy. They considered an active and mobile adversary that can corrupt the whole set of provers.

Recently, Ateniese et al. [3] considered the problem from the client side, where n clients store their respective files on a single prover in a manner such that the verification of the integrity of a single client's file simultaneously gives the integrity guarantee of the files of all the participating clients. They called such a system an entangled cloud storage.

#### 1.3 Comparison with Bowers, Juels and Oprea

The focus of this paper is PoR systems in the distributed setting; therefore, we only compare our work with existing works in the distributed setting. The scheme of Curtmola et al. [9] only considers multiple replica of the same underlying PDP systems, while the construction of Ateniese et al. [3] is for the multiple clients setting. In other words, the scheme of Bowers, Juels and Oprea [8] is closest to ours. However, there are a few key differences.

- (i) The construction of Bowers, Juels and Oprea [8] is secure only in the computational setting, while we provide security in the setting of unconditional security.
- (ii) Bowers, Juels and Oprea [8] use various tools and algorithms to construct their systems, including error-correcting codes, pseudo-random functions, message authentication codes and universal hash function families. On the other hand, we only use ramp schemes in our constructions, making our schemes easier to state and analyze, and arguably simpler to implement.
- (iii) We consider two types of security guarantees, namely, the worst-case scenario and the average-case scenario. On the other hand, Bowers, Juels and Oprea [8] only consider the worst-case scenario.
- (iv) The construction of Bowers, Juels and Oprea [8] only aims to protect the integrity of the message, while we consider both the privacy and integrity of the message. Privacy of data has emerged as an important requirement in cloud storage due to recent attacks [21].
- We work under a stronger requirement than [8] we require extraction to succeed with probabil-(v) ity equal to 1, whereas in [8], extraction succeeds with probability close to 1, depending in part on properties of a certain class of hash functions used in the protocol.

We use the term Prover to identify any server that stores the file of a client. We use the term Verifier for any entity that verifies whether the file of a client is stored properly or not by the server. We also assume that a file is composed of message blocks of an appropriate fixed length. If the file consists of single block, we simply call it the file.

# 2 Security model of multi-server PoR systems

The essential components of multi-server PoR (MPoR) systems are natural generalizations of single-server PoR systems. The first difference is that there are  $\rho$  provers and the Verifier might store different messages on each of them. Also, during an audit phase, the Verifier can pick a subset of provers on which it runs the audits. The last crucial difference is that the Extractor has (black-box or non-black-box) access to a subset of proving algorithms corresponding to the provers that the Verifier picked to audit. We detail them below for the sake of completeness.

Let Prover<sub> $\rho$ </sub>, . . . , Prover<sub> $\rho$ </sub> be a set of  $\rho$  provers. The Verifier has a message  $m \in \mathcal{M}$  from the message space  $\mathcal{M}$  which he redundantly encodes to  $M_1, \ldots, M_{\rho}$ .

- (i) In the keyed setting, the Verifier picks  $\rho$  different keys  $(K_1, \ldots, K_{\rho})$ , one for each of the corresponding
- (ii) The Verifier gives  $M_i$  to Prover<sub>i</sub>. In the case of a keyed scheme, Prover<sub>i</sub> may be also given an additional tag  $S_i$ , generated using the key  $K_i$ , and  $M_i$ .
- (iii) The Verifier stores some sort of information (say a *fingerprint* of the encoded message) which allows him to verify the responses made by the provers.
- (iv) On receiving the encoded message  $M_i$ , Prover generates a proving algorithm  $\mathcal{P}_i$ , which it uses to generate its responses during the auditing phase.
- (v) At any time, the Verifier picks an index i, where  $1 \le i \le \ell$ , and engages in a challenge-response protocol with  $Prover_i$ . In one execution of the challenge-response protocol, the Verifier picks a challenge c and gives it to Prover<sub>i</sub>, and the prover responds. The Verifier then verifies the correctness of the response (based on its fingerprint).
- (vi) The success probability  $\operatorname{succ}(\mathcal{P}_i)$  is the probability, computed over all the challenges, with which the Verifier accepts the response sent by Prover<sub>i</sub>.
- (vii) The Extractor is given a subset S of the proving algorithms  $\mathcal{P}_1, \ldots, \mathcal{P}_{\rho}$  (and in the case of a keyed scheme, the corresponding subset of the keys,  $\{K_i : i \in S\}$  and outputs a message  $\widehat{m}$ . The Extractor succeeds if  $\widehat{m} = m$ .

The above framework does not restrict any provers from interacting with other provers when they receive the encoded message. However, we assume that they do not interact after they have generated a proving algorithm. If we do not include this restriction, then it is not hard to see that one cannot have any meaningful protocol. For example, if provers can interact after they receive the encoded message, then it is possible that one prover stores the entire message and the other provers just relay the challenges to this specific prover and relay back its response to the verifier.

In contrast to a single-prover PoR scheme, there are two possible ways in which one can define the security of a multiple-prover PoR system. We define them next.

The first security definition corresponds to the "worst case" scenario and is the natural generalization of a single-server PoR system.

**Definition 2.1.** A  $\rho$ -prover MPoR scheme is  $(\eta, \nu, \tau, \rho)$ -threshold secure if there is an Extractor which, when given any  $\tau$  proving algorithms, say  $\mathcal{P}_{i_1}, \ldots, \mathcal{P}_{i_r}$ , succeeds with probability at least  $\nu$  whenever

$$\operatorname{succ}(\mathcal{P}_i) \geq \eta$$
 for all  $j \in I$ ,

where  $I = \{i_1, ..., i_{\tau}\}.$ 

We note that when  $\rho = \tau = 1$ , we get a standard single-server PoR system. Moreover, the definition captures the worst-case scenario in the sense that it only guarantees extraction if there exists a set of  $\tau$  proving algorithms, all of which succeed with high enough probability.

The above definition requires that all the  $\tau$  servers succeed with high enough probability. On the other hand, it might not be the case that all the proving algorithms of the servers picked by the Verifier succeed with the required probability. In fact, even verifying whether or not all the  $\tau$  proving algorithms have high enough success probability to allow successful extraction might be difficult (see, for example [15] for more details about this). However, it is possible that some of the proving algorithms succeed with high enough probability to compensate for the failure of the rest of the proving algorithms. For instance, since the provers are allowed to interact before they specify their proving algorithms, it might be the case that the colluding provers decide to store most of the message on a single prover. In this case, even a weaker guarantee that the average success probability is high enough might be sufficient to guarantee a successful extraction. In other words, it is possible to state (and as we show in this paper, achieve) a security guarantee with respect to the average case success probability over all the proving algorithms.

**Definition 2.2.** A  $\rho$ -prover MPoR scheme is  $(\eta, \nu, \rho)$ -average secure if the Extractor succeeds with probability at least  $\nu$  whenever

$$\frac{1}{\rho}\sum_{i=1}^{\rho}\operatorname{succ}(\mathcal{P}_i)\geq\eta.$$

Note that the average-case secure system reduces to the standard PoR scheme (with  $\tau = \rho$ ) when  $\rho = 1$ . The following example illustrates that average-case security is possible even when an MPoR system is not possible as per Definition 2.1.

**Example 2.3.** Suppose  $\eta = 0.7$ ,  $\nu = 0$  and  $\rho = 3$ . Further, suppose that  $succ(\mathcal{P}_1) = 0.9$ ,  $succ(\mathcal{P}_2) = 0.6$  and  $succ(\mathfrak{P}_3) = 0.6$ . Then the hypotheses of Definition 2.1 are not satisfied for  $\tau = 2$ . So even if the MPoR scheme is  $(\eta, \nu, \tau, \rho)$ -threshold secure, we cannot conclude that the Extractor will succeed. On the other hand, for the assumed success probabilities, the hypotheses of Definition 2.2 are satisfied. Therefore, if the MPoR scheme is  $(0.7, \nu, \tau)$ -average secure, the Extractor will succeed.

**Privacy guarantee.** We mentioned at the start of this section that PoR systems were introduced and studied to give assurance of the integrity of the data stored on remote storage. However, the confidentiality aspects of data have not been studied formally in the area of cloud-based PoR systems. There have been couple of ad hoc solutions that have been proposed in which the messages are encrypted and then stored on the cloud [9]. We believe that, in addition to the standard integrity requirement, privacy of the stored data when multiple provers are involved is also an important requirement. We model the privacy requirement as follows:

**Definition 2.4.** An MPoR system is called *t*-private if no set  $\mathcal{A}$  of adversarial provers of size at most *t* learns anything about the message stored by the Verifier.

Note that t = 0 corresponds to the case when the MPoR system does not provide any confidentiality to the message. The above definition captures the idea that, even if t provers collude, they do not learn anything about the message. We remark that we can achieve confidentiality without encrypting the message by using secret sharing techniques.

**Notation.** We fix the letter m for the original message,  $\mathcal{M}$  to denote the space from which the message m is picked and M to denote the encoded message. We fix v to denote the failure probability of the extractor and  $\eta$  to denote the success probability of a proving algorithm. In this paper, we are mainly interested in the case when v = 0 for both the worst-case and the average-case security. We use n to denote the number of message blocks, assuming the underlying PoR system breaks the message into blocks.

# 3 Primitives used in this paper

#### 3.1 Ramp schemes

In our construction, we use a primitive related to secret sharing schemes known as ramp schemes. A secret sharing scheme allows a trusted dealer to share a secret between *n* players so that certain subsets of players can reconstruct the secret from the shares they hold [6, 17].

It is well known that the size of each player's share in a secret sharing scheme must be at least the size of the secret. If the secret that is to be shared is large, then this constraint can be very restrictive. Schemes for which we can get a certain form of trade-off between share size and security are known as ramp schemes [7].

**Definition 3.1** (Ramp scheme). Let  $\tau_1$ ,  $\tau_2$  and n be positive integers such that  $\tau_1 < \tau_2 \le n$ . A  $(\tau_1, \tau_2, n)$ -ramp scheme is a pair of algorithms, say ShareGen and Reconstruct, such that, on input a secret S, ShareGen(S) generates *n* shares, one for each of the *n* players, such that the following two properties hold:

- (i) Reconstruction: Any subset of  $\tau_2$  or more players can pool together their shares and use Reconstruct to compute the secret S from the shares that they collectively hold.
- (ii) Secrecy: No subset of  $\tau_1$  or fewer players can determine any information about the secret S.

**Example 3.2.** Suppose the dealer wishes to set up a (2, 4, n)-ramp scheme with the secret  $(a_0, a_1)$ . The dealer picks a finite field  $\mathbb{F}_q$  with q > n such that  $a_0, a_1 \in \mathbb{F}_q$ . The dealer picks random elements  $a_2, a_3$ independently from the field  $\mathbb{F}_q$  and constructs the following polynomial of degree 3 over the finite field  $\mathbb{F}_q$ :  $f(x) = a_0 + a_1x + a_2x^2 + a_3x^3$ . The share for any player  $\mathcal{P}_i$  is generated by computing  $s_i = f(i)$ . It is easy to see that if two or fewer players come together, they do not learn any information about the secret, and if at least four players come together, they can use Lagrange's interpolation formula to compute the function *f* as well as the secret. However, if three players pool together their shares, then they can learn some partial information about one of the other player's share. For concreteness, let q = 17. Then  $5a_1 \equiv 7s_3 + 9s_6 + s_{15} \mod 17$ ; therefore, players  $\mathcal{P}_3$ ,  $\mathcal{P}_6$  and  $\mathcal{P}_{15}$  can compute the value of  $a_1$ .

For completeness, we review some of the basic theory concerning the construction of ramp schemes. Linear codes have been used to construct ramp schemes for over thirty years since the work of McEliece and Sarwate [13]. We will consider a construction from an arbitrary code in this paper. The following relation between an arbitrary code (linear or non-linear) and a ramp scheme was shown by Paterson and Stinson [14].

**Theorem 3.3.** Let C be a code of length N, distance d and dual distance  $d^{\perp}$ . Let  $1 \le s < d^{\perp} - 2$ . Then there is  $a(\tau_1, \tau_2, N - s)$ -ramp scheme, where  $\tau_1 = d^{\perp} - s - 1$  and  $\tau_2 = N - d + 1$ .

Here s is the *rate* of the ramp scheme. If G is a generator matrix of a code C with dimension k, then  $|C| = q^k \ge q^{d^{\perp}-1}$ . In other words,  $k \ge d^{\perp} - 1$ .

**Construction 3.4.** The construction of a ramp scheme from a code is as follows. Let s and  $\rho$  be positive integers, and let  $(m_1, \ldots, m_5) \in \mathbb{F}^5$  be the message. Let C be a code of length  $n = \rho + s$  defined over a finite field  $\mathbb{F}$ . We also require that the first s entries of a codeword is the message to be encoded, i.e., the corresponding gen-

- On input  $\mathcal{P}$ , compute the vector  $R' = (r'_c : c \in \Gamma)$ , where  $r'_c = \mathcal{P}(c)$  for all  $c \in \Gamma$  (i.e., for every  $c, r'_c$  is the response computed by  $\mathcal{P}$  when it is given the challenge c).
- Find  $\widehat{M} \in \mathbb{M}^*$ , so that dist $(R', r^{\widehat{M}})$  is minimized.
- (iii) Output  $\widehat{m} = e^{-1}(\widehat{M})$ .

Figure 1: Extractor for Theorem 3.6.

erator matrix is in standard form. Select a random codeword ( $\mathbf{c}_1 = m_1, \dots, \mathbf{c}_s = m_s, \mathbf{c}_{s+1}, \dots, \mathbf{c}_{n+s}$ )  $\in C$ , and define the shares as  $(\mathbf{c}_{s+1}, \ldots, \mathbf{c}_{\rho+s})$ .

**Example 3.5.** One can use a Reed–Solomon code to construct a ramp scheme [13]. Let q be a prime and  $1 \le s < t \le n < q$ . It is well known that, for a prime q, there is an  $[N, k, N - \tau + 1]_q$  Reed–Solomon code with  $d^{\perp} = \tau + 1$ . This implies a  $(\tau - s, \tau, N)$ -ramp scheme over  $\mathbb{F}_q$ .

### 3.2 Single-prover PoR system

We start by fixing some notation for PoR schemes that we use throughout the paper. Let  $\Gamma$  be the *challenge* space, and let  $\Delta$  be the response space. We denote by  $y = |\Gamma|$  the size of a challenge space. Let  $\mathcal{M}^*$  be the space of all encoded messages. The response function  $\rho \colon \mathbb{M}^* \times \Gamma \to \Delta$  computes the response  $r = \rho(M, c)$  given the encoded message M and the challenge c.

For an encoded message  $M \in \mathcal{M}^*$ , we define the *response vector*  $r^M$  that contains all the responses to all possible challenges for the encoded message M. Finally, define the response code of the scheme to be

$$\mathcal{R} = \{r^M : M \in \mathcal{M}^*\}.$$

The codewords in  $\Re$  are just the response vectors that we defined above. Previously [15], we proved the following result for a single-prover PoR scheme.

**Theorem 3.6.** Suppose that  $\mathcal{P}$  is a proving algorithm for a PoR scheme with response code  $\mathcal{R}$ . If the success probability of the corresponding proving algorithm satisfies  $\operatorname{succ}(\mathfrak{P}) \geq 1 - \tilde{\mathsf{d}}/2\gamma$ , where  $\tilde{\mathsf{d}}$  is the Hamming distance of the code  $\Re$ , and y is the size of the challenge space, then the extractor described in Figure 1 always outputs the message m.

If we cast this in the security model defined in Section 1 (Definition 2.1 and Definition 2.2), then we have the following theorem.

**Theorem 3.7.** Suppose that  $\mathcal{P}$  is a proving algorithm for a single server PoR scheme with response code  $\mathcal{R}$ . Then there exists a (1 - d/2y, 0, 1, 1)-MPoR system, where d is the Hamming distance of the code  $\Re$ , and y is the size of the challenge space  $\Gamma$ .

Previously [15], we gave a modified version of the Shacham-Waters scheme which they showed is secure in the unconditional security setting. They argued that, in the setting of unconditionally security, any keyed PoR scheme should be considered to be secure when the success probability of the proving algorithm  $\mathcal{P}$ , denoted by succ(P), is defined as the average success probability of the prover over all possible keys (Theorem 3.8). The same reasoning extends to MPoR systems. Therefore, in what follows and in Section 6, when we say a scheme is an  $(\eta, \nu, \tau, \rho)$ -threshold-secure scheme, the term  $\eta$  is the average success probability where the average is computed over all possible keys. We denote the average success probability of a prover  $\mathcal{P}$  over all possible keys by  $succ_{avg}(\mathcal{P})$ . Previously [15], we showed the following:

**Theorem 3.8.** Let  $\mathbb{F}_q$  be the underlying field, and let  $\ell \geq 1$  be the Hamming weight of the challenges made by the Verifier. Let d be the Hamming distance of the space of the encoded message  $M^*$ . Suppose that

$$\operatorname{succ}_{\operatorname{avg}}(\mathcal{P}) \gtrsim 1 - \frac{\mathsf{d}^*(q-1)}{2\gamma q},$$

where  $y = q^n$  is the size of the challenge space and  $d^*$  is given by

$$d^* \approx \binom{n}{\ell} (q-1)^{\ell} - \binom{n-d}{\ell} (q-1)^{\ell} - \sum_{w \ge 1} \binom{d}{w} \binom{n-d}{\ell-w} \frac{(q-1)^{\ell}}{q}. \tag{3.1}$$

Then there exists an Extractor that always outputs  $\widehat{m} = m$ .

# 4 Worst-case MPoR based on ramp scheme

In this section, we give our first construction that achieves a worst-case security guarantee. The idea is to use a  $(\tau_1, \tau_2, \rho)$ -ramp scheme in conjunction with a single-server PoR system. The intuition behind the construction is that the underlying PoR system along with the ramp scheme provides the retrievability guarantee and the ramp scheme provides the confidentiality guarantee.

We first present a schematic diagram of the working of an MPoR in Figure 2 and illustrate the scheme with the help of following example. We provide the details of the construction in Figure 3.

**Example 4.1.** Let  $\rho = 6$ . Suppose the Verifier and the provers use a PoR system  $\Pi$ . Let the message to be stored be (15, 3). The Verifier picks q = 17 and chooses two random elements 1,  $2 \in \mathbb{F}_{17}$  to construct a polynomial  $f(x) = 15 + 3x + x^2 + 2x^3$ . The Verifier picks an encoding function  $e(\cdot)$  and stores e(4) on Prover<sub>1</sub>, e(7) on Prover<sub>2</sub>, e(2) on Prover<sub>3</sub>, e(1) on Prover<sub>4</sub>, e(16) on Prover<sub>5</sub>, and e(8) on Prover<sub>6</sub>.

Let us suppose that the PoR scheme is such that, for a random challenge vector of dimension  $\rho$ , say (5, 2, 9, 13, 5, 6), where the *i*-th entry would be a challenge to Prover<sub>i</sub>, the corresponding responses of the provers form a vector (3, 14, 1, 13, 12, 14), where Resp<sub>i</sub> is the correct response of Prover<sub>i</sub>. In other words, on challenge 5 to Prover<sub>1</sub>, the correct response is 3, and so on.

During the audit phase, the Verifier picks any four provers and sends the challenges to the provers. Once all the provers that he chose reply, he verifies their response. For example, suppose the Verifier picks Prover<sub>1</sub>, Prover<sub>4</sub>, Prover<sub>4</sub> and Prover<sub>6</sub>. The Verifier then sends the challenge 5 to Prover<sub>1</sub>, 9 to Prover<sub>3</sub>, 13 to Prover<sub>4</sub> and 6 to Prover<sub>6</sub>. If it gets the responses 3, 1, 13 and 14 back, it accepts; otherwise, it rejects.

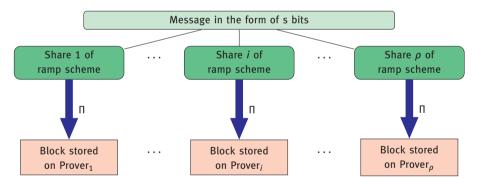


Figure 2: Schematic view of Ramp-MPoR system.

**Input.** The Verifier gets the message m as input. Let Prover<sub>1</sub>, . . . , Prover<sub> $\rho$ </sub> be the set of  $\rho$  provers. **Initialization Stage.** The Verifier performs the following steps for storing the message:

- (i) The Verifier chooses a single-server PoR system  $\Pi$  and a  $(\tau_1, \tau_2, \rho)$ -ramp scheme Ramp = (ShareGen, Reconstruct).
- (ii) The Verifier computes  $\rho$  shares of the message using the ramp scheme  $(m_1, \ldots, m_\rho) \leftarrow \text{ShareGen}(m)$ .
- (iii) The Verifier runs  $\rho$  independent copies of  $\Pi$  and generates the encoded share  $M_i = e(m_i) \in \mathbb{M}$  corresponding to each  $1 \le i \le \rho$ .
- (iv) The Verifier stores  $M_i$  on Prover<sub>i</sub>.

Challenge Phase. During the audit phase, the Verifier picks Prover; and runs the challenge-response protocol of  $\Pi$  with Prover;

Figure 3: Worst-case secure MPoR using a ramp scheme (Ramp-MPoR).

We note one of the possible practical deployments of the Ramp-MPoR stated in Figure 3. Let m be a message that consists of sk elements from  $\mathbb{F}_q$ . The Verifier breaks the message into k blocks of length s each. It then invokes a  $(\tau_1, \tau_2, n)$ -ramp scheme on each of these blocks to generate n shares of each of the k blocks. The Verifier then runs a PoR scheme  $\Pi$  to compute the encoded message to be stored on each of the servers by encoding its *k* shares, one corresponding to each of the *k* blocks.

We prove the following security result for the MPoR scheme presented in Figure 3.

**Theorem 4.2.** Let  $\Pi$  be an  $(\eta, 0, 1, 1)$ -threshold-secure MPoR with a response code of Hamming distance  $\tilde{d}$ and the size of challenge space y. Let Ramp = (ShareGen, Reconstruct) be a  $(\tau_1, \tau_2, \rho)$ -ramp scheme. Then Ramp-MPoR, defined in Figure 3, is an MPoR system with the following properties:

- (i) Privacy: Ramp-MPoR is  $\tau_1$ -private.
- (ii) Security: Ramp-MPoR is  $(\eta, 0, \tau_2, \rho)$ -threshold secure, where  $\eta = 1 \tilde{d}/2\gamma$ .

*Proof.* The privacy guarantee of Ramp-MPoR is straightforward from the privacy property of the underlying ramp scheme.

For the security guarantee, we need to demonstrate an Extractor that outputs a message  $\widehat{m} = m$  if at least t servers succeed with probability at least  $\eta = 1 - d/2y$ . The description of our Extractor is as follows:

- (i) The Extractor chooses  $\tau_2$  provers and runs the extraction algorithm of the underlying single-server PoR system on each of these provers. In the end, it outputs  $\widehat{M}_{i_i}$  for the corresponding provers Prover $_{i_i}$ . It defines  $S \leftarrow \{\widehat{M}_{i_1}, \ldots, \widehat{M}_{i_{\tau_0}}\}.$
- (ii) The Extractor invokes the Reconstruct algorithm of the underlying ramp scheme with the elements of S. It outputs whatever Reconstruct outputs.

Now note that the Verifier interacts with every Prover, independently. We know from the security of the underlying single-server PoR scheme (Theorem 3.6) that there is an extractor that always outputs the encoded message whenever  $\operatorname{succ}(\mathcal{P}_i) \geq \eta$ . Therefore, if all the  $\tau_2$  chosen proving algorithms succeed with probability at least  $\eta$ , then the set S will have  $\tau_2$  correct shares. From the correctness of the Reconstruct algorithm, we know that the message output in the end by the Extractor will be the message *m*.

As a special case of the above, we get a simple MPoR system which uses a replication code. A replication code has an encoding function

Enc: 
$$\Lambda \to \Lambda^{\rho}$$
 such that  $\operatorname{Enc}(x) = (\underbrace{x, x, \dots, x}_{\rho \text{ times}})$  for any  $x \in \Lambda$ .

This is the setting considered by Curtmola et al. [9].

We call a Ramp-MPoR scheme based on a replication code a Rep-MPoR. The schematic description of the scheme is presented in Figure 4, and the scheme is presented in Figure 5. Since a  $\rho$ -replication code is a  $(0, 1, \rho)$ -ramp scheme, a simple corollary to Theorem 4.2 is the following:

**Corollary 4.3.** Let  $\Pi$  be an  $(\eta, 0, 1, 1)$ -MPoR system with a response code of Hamming distance  $\tilde{d}$  and the size of challenge space y. Then Rep-MPoR, formed by instantiating Ramp-MPoR with the replication code based Ramp scheme, is an MPoR system with the following properties:

- (i) Privacy: It is 0-private.
- (ii) Security: It is  $(\eta, 0, 1, \rho)$ -threshold secure, where  $\eta = 1 \tilde{d}/2\gamma$ .

The issue with Rep-MPoR scheme is that there is no confidentiality of the file. We will come back to this issue later in Section 6.1.

# 5 Average-case secure MPoR system

In general, it is not possible to verify with certainty whether the success probability of a proving algorithm is above a certain threshold; therefore, in that case, it is unclear how the Extractor would know which proving algorithms to use for extraction as described in Section 4. In this section, we analyze the average-case security

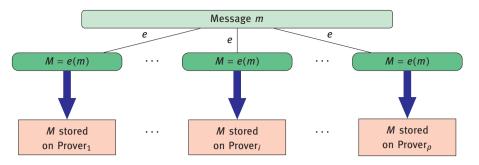


Figure 4: Schematic view of Rep-MPoR.

**Input.** The Verifier gets the message m as input. Let Prover<sub>1</sub>, ..., Prover<sub>p</sub> be the set of p provers.

Initialization stage. The Verifier performs the following steps for storing the message:

- (i) The Verifier chooses a single-server PoR system  $\Pi$ .
- (ii) Using the encoding scheme of  $\Pi$ , the Verifier generates the encoded message  $M = e(m) \in \mathcal{M}$  for  $1 \le i \le n$ .
- (iii) The Verifier stores the message M on all Prover<sub>i</sub> for  $1 \le i \le n$ .

Challenge phase. During the audit phase, the Verifier runs the challenge-response protocol of  $\Pi$  independently on each server.

Figure 5: Average-case secure MPoR (Rep-MPoR).

properties of the replication code based scheme, Rep-MPoR, described in the last section. This allows us an alternative guarantee that allows successful extraction where the extractor need not worry whether a certain proving algorithm succeeds with high enough probability or not.

Recall the scenario introduced in Example 2.3. Here we assumed  $succ(\mathcal{P}_1) = 0.9$ ,  $succ(\mathcal{P}_2) = 0.6$  and  $succ(\mathcal{P}_3) = 0.6$  for three provers. Suppose that successful extraction for a particular prover  $\mathcal{P}_i$  requires  $succ(\mathcal{P}_2) \geq 0.7$ . Then extraction would work on only one of these three provers. On the other hand, suppose we have an average-case secure MPoR in which extraction is successful if the average success probability of the three provers is at least 0.7. Then the success probabilities assumed above would be sufficient to guarantee successful extraction.

**Theorem 5.1.** Let  $\Pi$  be a single-server PoR system with a response code of Hamming distance  $\tilde{d}$  and the size of challenge space  $\gamma$ . Then Rep-MPoR, defined in Figure 5, is an MPoR system with the following properties:

- (i) Privacy: Rep-MPoR is 0-private.
- (ii) Security: Rep-MPoR is  $(1 \tilde{d}/2\gamma, 0, \rho)$ -average secure.

*Proof.* Since the message is stored in its entirety on each of the servers, there is no confidentiality.

For the security guarantee, we need to demonstrate an Extractor that outputs a message  $\widehat{m} = m$  if the average success probability of all the provers is at least  $\eta = 1 - \widetilde{d}/2\gamma$ . The description of our Extractor is as follows:

- (i) For all  $1 \le i \le n$ , use  $\mathcal{P}_i$  to compute the vector  $R_i = (r_c^{(i)} : c \in \Gamma)$ , where  $r_c^{(i)} = \mathcal{P}_i(c)$  for all  $c \in \Gamma$  (i.e., for every c,  $r_c^{(i)}$  is the response computed by  $\mathcal{P}_i$  when it is given the challenge c).
- (ii) Compute R as a concatenation of  $R_1, \ldots, R_\rho$  and find  $\widehat{M} := (\widehat{M}_1, \ldots, \widehat{M}_\rho)$  so that  $\operatorname{dist}(R, r^{\widehat{M}})$  is minimized.
- (iii) Compute  $m = e^{-1}(\widehat{M})$ .

Now note that the Verifier interacts with each Prover<sub>i</sub> independently and the Extractor uses the challengeresponse step with independent challenges. Let  $\eta_1, \ldots, \eta_\rho$  be the success probabilities of the  $\rho$  proving algorithms. Let  $\bar{\eta}$  be the average success probability over all the servers and challenges. Therefore,  $\bar{\eta} = \rho^{-1} \sum_{i=1}^{\rho} \eta_i$ .

First note that, in the case of Figure 5, the response code is of the form

$$\{(\underbrace{r,r,\ldots,r}_{o \text{ times}}):r\in\mathbb{R}\}.$$

It is easy to see that the distance of the response code is  $\rho \tilde{d}$  and the length of a challenge is  $\rho \gamma$ . From the definition of the extractor and Theorem 3.6, it follows that the extraction succeeds if

$$\frac{\eta_1 + \dots + \eta_\rho}{\rho} = \bar{\eta} \ge 1 - \frac{\tilde{\mathsf{d}}}{2y}.$$

### 5.1 Hypothesis testing for Rep-MPoR

For the purposes of auditing whether a file is being stored appropriately, it is necessary to have a mechanism for determining whether the success probability of a prover is sufficiently high. In the case of replication code based on MPoR with worst-case security, we are interested in the success probabilities of individual provers, and the analysis can be carried out as detailed in [15]. In the case of Rep-MPoR, however, we wish to determine whether the *average* success probability of the set of provers  $\{\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_o\}$  is at least  $\eta$ . This amounts to distinguishing the null hypothesis

$$H_0$$
: avg-succ( $\mathcal{P}_i$ ) <  $\eta$ 

from the alternative hypothesis

$$H_1$$
: avg-succ( $\mathcal{P}_i$ )  $\geq \eta$ .

Suppose we send *c* challenges to each server. If a given server  $\mathcal{P}_i$  has success probability succ( $\mathcal{P}_i$ ), then the number of correct responses received follows the binomial distribution  $B(c, succ(\mathcal{P}_i))$ . If the success probabilities  $\operatorname{succ}(\mathcal{P}_i)$  were the same for each server, then the sum of the number of successes over all the servers would also follow a binomial distribution. However, we are also interested in the case in which these success probabilities differ, in which case the total number of successes follows a Poisson binomial distribution, which is more complicated to work with. In order to establish a test that is conceptually and computationally easier to apply, we will instead rely on the observation that, in cases where the average success probability is high enough to permit extraction, the failure rates of the servers are relatively low.

For a given server  $\mathcal{P}_i$ , let  $f_i = 1 - \text{succ}(\mathcal{P}_i)$  denote the probability of failure. For r challenges, the number of failures follows the binomial distribution  $B(c, f_i)$ . Provided that r is sufficiently large and  $f_i$  is sufficiently low, then  $B(c, f_i)$  can be approximated by the *Poisson distribution*  $Pois(cf_i)$ . The Poisson distribution  $Pois(\lambda)$ is used to model the scenario where discrete events are occurring independently within a given time period with an expected rate of  $\lambda$  events during that period. The probability of observing k events within that period is given by

$$P(k) = \frac{e^{-\lambda} \lambda^k}{k!}.$$

Mean and variance of  $Pois(\lambda)$  are equal to  $\lambda$ . For our purposes, the advantage of using this approximation is that the sum of  $\rho$  independent variables following the Poisson distributions  $Pois(\lambda_1)$ ,  $Pois(\lambda_2)$ , ...,  $Pois(\lambda_0)$ is itself distributed according to the Poisson distribution  $Pois(\lambda_1 + \lambda_2 + \cdots + \lambda_\rho)$ , even when the  $\lambda_i$  all differ. In the case where the average failure probability is low, the distribution  $Pois(c(f_1 + f_2 + \cdots + f_\rho))$  should provide a reasonable approximation to the actual distribution of the total number of failed challenges.

**Example 5.2.** To demonstrate the appropriateness of the Poisson approximation for this application, suppose we have five servers, whose failure probabilities are expressed as  $\mathbf{f} = (f_1, f_2, \dots, f_5)$ . Let t be the number of trials per server and b the total number of observed failures out of the 5t trials. Table 1 gives both the exact cumulative probability  $Pr[B \le b]$  of observing up to b failures, and the Poisson approximation  $Pr_{Pois}[B \le b]$ of this cumulative probability, for a range of values for **f**.

As an example of using the given formula to calculate a confidence interval, suppose we do 200 trials on each of five servers (so there are 1000 trials in total), and we observe 50 failures in total. Then the resulting confidence interval is [0, 63.29). Suppose we wish to know whether the success probability is at least  $\eta = 0.9$ . We have  $(1-0.9) \times 1000 = 100$ . This is outside of that interval, and hence we conclude there is enough evidence to reject  $H_0$  at the 95 % significance level. However, to test whether the success probability was greater than 0.95 we see that  $(1 - 0.95) \times 1000 = 50$ . Since 50 lies within the interval, we conclude there is insufficient evidence to reject  $H_0$  at the 95 % significance level.

t	b	$Pr[B \leq b]$	$Pr_{Pois}[B \leq b]$	t	b	F
<b>f</b> = (0	0.1, 0.1	,0.1,0.1,0.1)		f = (0	.1,0.	1,
200	5	$2.556545692 \times 10^{-38}$	$3.261456422 \times 10^{-36}$	20	0	(
200	10	$1.450898832 \times 10^{-32}$	$1.137687971 \times 10^{-30}$	20	5	(
200	50	$5.995167631 \times 10^{-9}$	$2.401592276 \times 10^{-8}$	20	10	(
200	100	0.5265990813	0.5265622074	20	15	(
100	0	$1.322070819 \times 10^{-23}$	$1.928749864 \times 10^{-22}$	20	20	(
100	5	$6.272915577 \times 10^{-17}$	$5.567756307 \times 10^{-16}$	40	0	7
100	10	$1.135691814 \times 10^{-12}$	$6.450152972 \times 10^{-12}$	40	5	(
100	15	$1.662665039 \times 10^{-9}$	$6.357982164 \times 10^{-9}$	40	10	(
100	20	$4.557480806 \times 10^{-7}$	0.000001235187232	40	15	(
200	0	$1.747871252 \times 10^{-46}$	$3.720076039 \times 10^{-44}$	40	20	(
200	5	$2.556545692 \times 10^{-38}$	$3.261456422 \times 10^{-36}$	100	20	4
200	10	$1.450898832 \times 10^{-32}$	$1.137687971 \times 10^{-30}$	100	25	(
200	15	$6.757345217 \times 10^{-28}$	$3.340076418 \times 10^{-26}$	100	30	(
200	20	$5.962487876 \times 10^{-24}$	$1.905558774 \times 10^{-22}$	100	35	(
500	20	$1.240463044 \times 10^{-84}$	$1.084188102 \times 10^{-79}$	100	40	(
500	25	$3.140367419 \times 10^{-79}$	$1.697380630 \times 10^{-74}$	f _ (0	0.01, 0	٠ ،
500	30	$2.935666094 \times 10^{-74}$	$9.912214279 \times 10^{-70}$	20	0.01, 0	).U )
500	35	$1.193158517 \times 10^{-69}$	$2.542280876 \times 10^{-65}$	20	5	(
500	40	$2.369596756 \times 10^{-65}$	$3.218593843 \times 10^{-61}$	20	10	(
<b>f</b> = ((	0.01.0.	01, 0.01, 0.01, 0.01)		20	15	
200	5	0.06613951161	0.06708596299	20	20	
200	10	0.5830408032	0.5830397512	40	0	(
200	20	0.9985035184	0.9984117410	40	5	(
200	50	≈ 1	≈ 1	40	10	(
<b>f</b> – (C	200	1, 0.02, 0.03, 0.04)		40	15	(
1 – (c 200	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	$9.651421837 \times 10^{-22}$	$6.180223643 \times 10^{-20}$	40	20	(
200	10	$5.539867010 \times 10^{-17}$	$1.744235672 \times 10^{-15}$	100	20	(
200	20	0.09020056729	0.1076778797	100	25	(
200	50	0.9999999198	0.9999991415	100	30	(
			0.7777771413	100	35	(
•		01, 0.03, 0.04, 0.05)	7	100	40	(
200	5	$8.312224722 \times 10^{-8}$	$1.196952269 \times 10^{-7}$	<b>f</b> = (0	.02,0	0.0
200	10	0.00006809921297	0.00008550688580	20	0	(
200	20	0.06901537242	0.07274102693	20	5	(
200	50	0.9999582547	0.9999397284	20	10	(
				20	15	(
						•
				20	20	:

t	b	$Pr[B \leq b]$	$Pr_{Pois}[B \leq b]$					
f = (0.1, 0.1, 0.1, 0.1, 0.1)								
20	0	0.00002656139888	0.00004539992984					
20	5	0.05757688648	0.06708596299					
20	10	0.5831555123	0.5830397512					
20	15	0.9601094730	0.9512595983					
20	20	0.9991924263	0.9984117410					
40	0	$7.055079108 \times 10^{-10}$	$2.061153629 \times 10^{-9}$					
40	5	0.00003871193246	0.00007190884076					
40	10	0.008071249954	0.01081171886					
40	15	0.1430754340	0.1565131351					
40	20	0.5591747822	0.5590925860					
100	20	$4.557480806 \times 10^{-7}$	0.000001235187232					
100	25	0.00003540113222	0.00007160717427					
100	30	0.001002549708	0.001594027332					
100	35	0.01231948910	0.01621388016					
100	40	0.07508928967	0.08607000083					
f = (0	0.01,0	0.01, 0.01, 0.01, 0.01)						
20	0	0.3660323413	0.3678794412					
20	5	0.9994654657	0.9994058153					
20	10	0.9999999999	0.9999999900					
20	15	1.000000000	1.000000000					
20	20	1.000000000	1.000000000					
40	0	0.1339796748	0.1353352833					
40	5	0.9839770930	0.9834363920					
40	10	0.9999931182	0.9999916922					
40	15	0.9999999996	1.000000000					
40	20	0.999999999	1.000000000					
100	20	0.9999999367	0.9999999198					
100	25	0.999999999	1.000000001					
100	30	0.999999999	1.000000001					
100	35	0.999999999	1.000000001					
100	40	0.999999999	1.00000001					
<b>f</b> = (0	f = (0.02, 0.0075, 0.0075, 0.0075, 0.0075)							
20	0		0.09536916225					
20	5		0.9672561739					
20		0.9999843669	0.9999642885					
20	15	0.9999999995	0.9999999958					
20	20	1.000000000	1.000000000					
40	0	0.007986825382	0.009095277109					
40	5	0.6699740391	0.6684384858					
40	10	0.9927425867	0.9909776597					
40	15	0.9999835852	0.9999661876					
40	20	0.9999999935	0.9999999715					
100	20	0.9999999935	0.9999999715					
100	25	0.9999999998	1.000000001					
100	30	0.9999999998	1.000000001					
100	35	0.9999999998	1.000000001					
100	40	0.9999999998	1.000000001					
100	40	0.777777777	1.00000001					

 $\textbf{Table 1:} \ Comparison \ between \ exact \ cumulative \ probability \ and \ approximation \ by \ Poisson \ distribution.$ 

Let b denote the number of incorrect responses we have received from the  $c\rho$  challenges given to the provers. Suppose that  $H_0$  is true, so that the expected number of failures is at least  $\eta \rho c$ . Based on our approximation, the probability that the number of failures is at most b is at most

$$\sum_{i=0}^{b} \frac{e^{-\eta \rho c} (\eta \rho c)^{i}}{i!}.$$

If this probability is less than 0.05, we reject  $H_0$  and accept the alternative hypothesis. However, if the probability is greater than 0.05, then there is not enough evidence to reject  $H_0$  at the 5 % significance level, and so we continue to suspect that the file is not being stored appropriately.

We can express this test neatly using a *confidence interval*. We define a 95 % upper confidence bound by

$$\lambda_{\mathrm{U}} = \inf \left\{ \lambda \mid \sum_{i=0}^{b} \frac{e^{-\lambda} \lambda^{i}}{i!} < 0.05 \right\}.$$

This represents the smallest parameter choice for the Poisson distribution for which the probability of obtaining b or fewer incorrect responses is less than 0.05. Then  $[0, \lambda_{\rm II}]$  is a 95 % confidence interval for the mean number of failures, so we reject  $H_0$  whenever  $\eta nr$  lies outside this interval. The value of  $\lambda_U$  can be determined easily by exploiting a connection with the chi-squared distribution [18]. We have

$$\sum_{i=0}^{b} \frac{e^{-\lambda} \lambda^{i}}{i!} = \Pr(\chi_{2b+2}^{2} > 2\lambda),$$

and so the appropriate value of  $\lambda_{\rm II}$  can readily be obtained from a table for the chi-squared distribution.

We give a comparison between exact cumulative probability and approximation by Poisson distribution in Table 1.

# Optimization using the keyed Shacham-Waters scheme

In the last three sections, we gave constructions of MPoR scheme using ramp schemes, linear secret-sharing schemes, replication codes and a single-prover PoR system. In this section, we show a specific instantiation of our scheme using the keyed scheme of Shacham and Waters [15, 16] for a single-server PoR system.

### 6.1 Extension of the keyed Shacham-Waters scheme to MPoR

If we instantiate the Rep-MPoR scheme (described in Section 4) with the modified Shacham-Waters scheme of [15], then we need one key that consists of n + 1 values in  $\mathbb{F}_q$ . However, in this case, we do not have any privacy. In particular, we have the following extension of Corollary 4.3.

**Corollary 6.1.** Let  $\Pi$  be an (n,0,0,1)-PoR system of Shacham and Waters [16] with a response code of Hamming distance d and the size of challenge space y, where d is given by equation (3.1). Then Rep-MPoR instantiated with the Shacham-Waters scheme is an MPoR system with the following properties:

- (i) *Privacy: It is* 0-*private.*
- (ii) Security: It is  $(\eta, 0, 1, \rho)$ -threshold secure, where  $\eta = 1 \frac{\tilde{d}(q-1)}{2\gamma q}$ .
- (iii) Storage Overhead: The Verifier needs to store n + 1 field elements, and every Prover<sub>i</sub> needs to store 2n field elements.

*Proof.* The results follow by combining Theorem 3.8 with Corollary 4.3.

The issue with the Rep-MPoR scheme is that there is no confidentiality of the file. In what follows, we improve the privacy guarantee of the MPoR scheme described above. Our starting point would be an instantiation of the Ramp-MPoR scheme, defined in Figure 3, with the Shacham-Waters scheme. We then reduce the storage on the Verifier through two steps.

### 6.2 Optimized version of the multi-server Shacham-Waters scheme

We follow two steps to get an MPoR scheme based on the Shacham-Waters scheme with a reduced storage requirement for the Verifier, while improving the confidentiality guarantee.

- (i) In the first step, stated in Theorem 6.2, we improve the privacy guarantee of the MPoR scheme to get a  $\tau_1$ -private MPoR scheme (where  $\tau_1 < \rho$  is an integer). The Verifier in this scheme has to store  $\rho(n+1)$ field elements. When the underlying field is  $\mathbb{F}_q$ , the verifier has to store  $\rho(n+1)\log q$  bits.
- (ii) In the second step, stated in Theorem 6.3, we reduce the storage requirement of the Verifier from  $\rho(n+1)$ to  $\tau_1(n+1)$  field elements for some integer  $\tau_1 < \rho$  without affecting the privacy guarantee. When the underlying field is  $\mathbb{F}_q$ , the verifier has to store  $\tau_1(n+1)\log q$  bits.

**Step 1.** To improve the privacy guarantee of Corollary 6.1 to say,  $\tau_1$ -private (as per Definition 2.4), we use a Ramp-MPoR scheme and  $\rho$  different keys, where each key consists of n+1 values in  $\mathbb{F}_q$ . The Verifier generates  $\rho$  shares of every message block using a ramp scheme, then encodes the shares, and finally computes the tag for each of these encoded shares.

We follow with more details. Let  $m = (m[1], \dots, m[k])$  be the message. The Verifier computes the shares of every message block  $(m[1], \ldots, m[k])$  using a  $(\tau_1, \tau_2, \rho)$ -ramp scheme. It then encodes all the shares using the encoding scheme of the PoR scheme. Let the resulting encoded shares be  $M_i[1], \ldots, M_i[n]$  for  $1 \le i \le \rho$ . In other words, the result of the above two steps are  $\rho$  encoded shares, each of which is an *n*-tuple in  $(\mathbb{F}_a)^n$ . The Verifier now picks random values  $a^{(i)}, b_1^{(i)}, \ldots, b_n^{(i)} \in \mathbb{F}_q$  for  $1 \le i \le \rho$  and computes the tags as follows:

$$S_i[j] = b_i^{(i)} + a^{(i)}M_i[j]$$
 for  $1 \le i \le \rho$ ,  $1 \le j \le n$ .

The verifier gives  $Prover_i$  the tuple of encoded messages  $(M_i[1], \ldots, M_i[n])$  and the corresponding tags  $(S_i[1], \ldots, S_i[n])$ . We call this scheme the Basic-MPoR scheme. The following is straightforward from Theorem 4.2.

**Theorem 6.2.** Let  $\Pi$  be an  $(\eta, 0, 0, 1)$ -PoR scheme of Shacham and Waters [16] with a response code of Hamming distance  $\tilde{d}$  and the size of challenge space  $y = q^n$ , where  $\tilde{d}$  is given by equation (3.1). Let Ramp be  $a\left(\tau_{1},\tau_{2},\rho\right)$ -ramp scheme. Then Basic-MPoR defined above is an MPoR scheme with the following properties:

- (i) Privacy: Basic-MPoR is  $\tau_1$ -private.
- (ii) Security: Basic-MPoR is  $(\eta, 0, \tau_2, \rho)$ -threshold secure, where  $\eta = 1 \frac{d(q-1)}{2\gamma q}$ .
- (iii) Storage Overhead: The Verifier needs to store  $\rho(n+1)$  field elements and every Prover; needs to store 2nfield elements.

In the construction mentioned above, the Verifier needs to store  $\rho(n+1)$  elements of  $\mathbb{F}_q$ , which is almost the same as the total storage requirements of all the provers. In [15], we encountered the same issue, where the Verifier has to store as much secret information as the size of the message. This seems to be the general drawback in the unconditionally secure setting. However, in the case of MPoR, we can improve the storage requirement of the Verifier as shown in the next step.

**Step 2.** In this step, we improve the above-described MPoR scheme to achieve considerable reduction on the storage requirement of the Verifier. The resulting scheme also provides unbounded audit capability against computationally unbounded adversarial provers, and it also ensures  $\tau_1$ -privacy.

The main observation that results in the reduction in the storage requirements of the Verifier is the fact that we can partially derandomize the keys generated by the Verifier. We use one of the most common techniques in derandomization. The keys in this scheme are generated by  $\tau_1$ -wise independent functions. Our construction works as follows: We pick n + 1 random polynomials,  $f_1(x), \ldots, f_n(x), g(x) \in \mathbb{F}_q[x]$ , each of degree at most  $\tau_1$  – 1. Then the Verifier computes the secret key by evaluating the polynomials  $f_i(x)$  and g(x) on  $\rho$  different values, say

$$b_j^{(i)} = f_j(i)$$
 and  $a_i = g(i)$ 

<sup>1</sup> A function is a  $\tau_1$ -wise independent function if every subset of  $\tau_1$  outputs is independent and equally likely. It should be noted that this does not imply that all the outputs of the function are mutually independent.

**Input.** The Verifier gets a message  $m = (m[1], \ldots, m[n])$  as input. Let Prover<sub>1</sub>, ..., Prover<sub>p</sub> be the set of  $\rho$  provers. Let q be a prime number greater than  $\rho$ .

Initialization Stage. The Verifier performs the following steps:

- The Verifier chooses n+1 random polynomials of degree at most  $\tau_1-1, f_1(x), \ldots, f_n(x), g(x) \in \mathbb{F}_q[x]$  and a  $(\tau_1, \tau_2, \rho)$ -ramp scheme Ramp = (ShareGen, Reconstruct).
- (ii) For every server i, the Verifier does the following:
  - (a) Compute  $\rho$  shares of every message block using the share generation algorithm of Ramp as follows:  $(m_1[j], \ldots, m_p[j]) \leftarrow \mathsf{ShareGen}(m[j]) \text{ for } 1 \leq j \leq n.$
  - (b) The Verifier encodes the message as  $e(m_i[j]) = M_i[j]$  for  $1 \le j \le n$ ,  $1 \le i \le \rho$ .
  - (c) Compute  $b_1^{(i)} = f_1(i), \ldots, b_n^{(i)} = f_n(i), a^{(i)} = g(i).$
  - (d) Compute the tag  $S_i[j] = b_i^{(i)} + a^{(i)}M_i[j]$  for  $1 \le j \le n$ .
- (iii) The Verifier gives  $\{(M_i[j], S_i[j])\}_{1 \le j \le n}$  to Prover<sub>i</sub>.

Challenge Phase. During the audit phase, the Verifier picks a prover, Proveri, and runs the challenge-response algorithm of a single-server Shacham-Waters scheme. It computes the corresponding keys by computing the random polynomials chosen during the set up phase.

Figure 6: MPoR using optimized Shacham-Waters scheme (SW-MPoR).

for  $1 \le i \le n$  and  $1 \le i \le p$ . The Verifier then computes the encoded shares and their corresponding tags as in Basic-MPoR, i.e.,

$$S_i[j] = b_i^{(i)} + a^{(i)}M_i[j]$$
 for  $1 \le i \le \rho$ ,  $1 \le j \le n$ .

Figure 6 is the formal description of this scheme. For the scheme described in Figure 6, we prove the following result.

**Theorem 6.3.** Let Ramp = (ShareGen, Reconstruct) be a  $(\tau_1, \tau_2, \rho)$ -ramp scheme. Let  $\Pi$  be a single-prover Shacham-Waters scheme [16] with a response code of Hamming distance d and the size of challenge space y. Then SW-MPoR, defined in Figure 6, is an MPoR system with the following properties:

- (i) Privacy: SW-MPoR is  $\tau_1$ -private.
- (ii) Security: SW-MPoR is  $(\eta, 0, \tau_2, \rho)$ -threshold secure, where  $\eta = 1 \frac{\bar{d}(q-1)}{2\gamma q}$ . (iii) Storage Overhead: The Verifier needs to store  $\tau_1(n+1)$  field elements, and every Prover<sub>i</sub> (for  $1 \le i \le \rho$ ) needs to store 2n field elements.

*Proof.* The privacy guarantee of SW-MPoR is straightforward from the secrecy property of the underlying ramp scheme.

For the security guarantee, we have to show an explicit construction of the Extractor that, on input proving algorithms  $\mathcal{P}_1, \ldots, \mathcal{P}_{\rho}$ , outputs m if  $\operatorname{succ}(\mathcal{P}_i) > \eta$  for at least  $\tau_2$  proving algorithms. However, there is a subtle issue that we have to deal with before using the proof of Theorem 4.2, because of the relation between every message and tag pair. It was previously noted by us [15] that if the adversarial prover learns the secret key, then it can break the PoR scheme. We first argue that a set of  $\tau_1$  colluding provers cannot have an undue advantage from exploiting the linear structure of the message-tag pairs.

We now prove that any set of  $\tau_1$  provers do not learn anything about the keys generated using n+1polynomials of degree at most  $\tau_1$  – 1. The idea is very similar to the single-prover case. Previously [15], we noted that in the single prover case, for an n-tuple encoded message, the key is a tuple of n + 1 uniformly random elements  $(a, b_1, \ldots, b_n)$  in  $\mathbb{F}_q$ . Further, from the point of view of a prover, there are q possible keys – the value of a determines the n-tuple  $(b_1, \ldots, b_n)$  uniquely, but a is completely undetermined. In the MPoR case, we have  $\rho$  keys. Each prover in a given set of  $\tau_1$  provers has q possible keys, as discussed above. However, it is conceivable that they can use their collective knowledge to learn something about the keys. In what follows, we show that they cannot determine any additional information about their keys by combining the information they hold collectively.

Let  $I = \{i_1, \ldots, i_{\tau_1}\}$  be the indices of any arbitrary set of  $\tau_1$  provers. Let  $S_i$  denote the set of possible keys for Prover<sub>i</sub>, for  $i \in I$ . Consider any list of  $\tau_1$  keys  $(K_{i_1}, K_{i_2}, \ldots, K_{i_{\tau_1}})$ . Recall that  $K_i$  (for  $i \in I$ ) has the form  $(a^{(i)}, b_1^{(i)}, \ldots, b_n^{(i)})$ , where  $a^{(i)}$  and  $b_i^{(i)}$  (for  $1 \le j \le n$ ) are generated by random polynomials of degree  $\tau_1$ . We first consider  $a^{(i)}$  (for  $i \in I$ ). Note that the vector  $(b_1^{(i)}, \ldots, b_n^{(i)})$  is defined uniquely by  $a^{(i)}$  and the set of all encoded message-tag pairs. We have already shown that any set of  $\tau_1$  provers cannot learn anything about the random polynomial g(x) used to generate the  $a^{(i)}$  for all  $i \in I$ . We use the following well-known fact to show that any set of  $\tau_1$  provers does not learn any additional information about the keys.

**Fact 6.4.** Let t > 0 be an integer, let q be a prime number, and let  $\mathbb{F}_q$  be a finite field. Let  $h_0, h_1, \ldots, h_{t-1} \in \mathbb{F}_q$  be random elements picked uniformly at random. Define  $h(x) = \sum_{i=0}^{t-1} h_i x^i$  for all  $\alpha \in \mathbb{F}_q$ . Then,

$$\Pr[h(x_1) = y_1 \wedge \dots \wedge h(x_\tau) = y_t] = \prod_{i=1}^t \Pr[h(x_{\alpha_i}) = y_i].$$
 (6.1)

Since h(x) is uniformly distributed in  $\mathbb{F}_q$ , the probability computed in equation (6.1) is actually equal to  $q^{-t}$ .

By construction, g(x) is a random polynomial of degree at most  $\tau_1 - 1$ . Fact 6.4 then implies that any combination of  $\{a^{(i)}\}_{i \in I}$  is equally likely. A similar argument, with the  $a^{(i)}$ 's replaced by the  $b_j^{(i)}$ 's (for all  $i \in I$  and  $1 \le j \le n$ ) and the polynomial g(x) replaced by  $f_j(x)$  (for  $1 \le j \le n$ ), gives that all sets of  $\tau_1$  keys are equally likely. In other words, the set of provers in the set I cannot determine any additional information about their keys by combining the information they hold collectively.

We now complete the security proof by describing an Extractor that outputs the file if  $\tau_2$  provers succeed with high enough probability. The description of the Extractor and its analysis is the same as that of Theorem 4.2. We give it for the sake of completeness.

- (i) The Extractor chooses  $\tau_2$  provers and runs the extraction algorithm of the underlying single-server PoR system on each of these provers. In the end, it outputs  $\widehat{M}_{i_j}$  for the corresponding provers Prover<sub> $i_j$ </sub>. It defines  $\mathcal{S} \leftarrow \{\widehat{m}_{i_1}, \ldots, \widehat{m}_{i_{\tau_2}}\}$ . Note that the Extractor of the underlying PoR scheme has already computed  $e^{-1}$  on the set  $\{\widehat{M}_{i_1}, \ldots, \widehat{M}_{i_{\tau_2}}\}$ .
- (ii) The Extractor invokes the Reconstruct algorithm of the underlying ramp scheme with the elements of  $\tilde{S}$  to compute m'.

Now note that the Verifier interacts with every  $\operatorname{Prover}_i$  independently. We know from the security of the underlying  $\operatorname{PoR}$  scheme of Shacham–Waters that there is an extractor that always outputs the encoded message whenever  $\operatorname{succ}_{\operatorname{avg}}(\mathcal{P}_i) \geq \eta$ . Therefore, if all the  $\tau_2$  chosen proving algorithms succeed with probability at least  $\eta$  over all possible keys, then the set  $\mathbb{S}$  will have  $\tau_2$  correct shares. From the correctness of the Reconstruct algorithm and  $e^{-1}(\cdot)$ , we know that the message output in the end by the Extractor will be the message m.

For the storage requirement, the Verifier has to store the coefficients of all the random polynomials  $f_1(x), \ldots, f_n(x), g(x)$ , which amounts to a total of  $\tau_1(n+1) = \tau_1 n + n$  field elements.

### 7 Conclusion and future works

In this paper, we studied PoR systems when multiple provers are involved (MPoR). We motivated and defined the security of MPoR in the worst-case (Definition 2.1) and the average-case (Definition 2.2) settings, and extended the hypothesis testing techniques used in the single-server setting [15] to the multi-server setting. We also motivated the study of confidentiality of the outsourced message. We gave MPoR schemes which are secure under both these security definitions and provide reasonable confidentiality guarantees even when there is no restriction on the computational power of the servers. At the end of this paper, we looked at an optimized version of MPoR system when instantiated with the unconditionally secure version of the Shacham–Waters scheme [16]. We exhibited that, in the multi-server setting with computationally unbounded provers, one can overcome the limitation that the verifier needs to store as much secret information as the provers.

Our paper leaves several open problems. We list two of them below:

(i) Our approach works in the privately verifiable setting, i.e., the entity that wishes to verify the validity of stored data is the same entity that stored the data. It would be interesting to see if our schemes can be extended to publicly verifiable setting.

(ii) We assume that the provers do not interact with each other after they receive the encoded files. There is a vast literature on mitigating collusion. It is an interesting direction to see if our schemes can be combined with the recent advances in secure scheme against colluding players in the distributed setting to remove our assumption.

# Notation used in this paper

challenge С

 $C^{\perp}$ dual of a code C

۴h distance of the response code

d distance of a codeword dΤ dual distance of a code

dist Hamming distance between two vectors

G generator matrix of a code

k length of a message K key (in a keyed scheme) l number of message-blocks

m message

*i*-th message block m[i]

message outputted by the Extractor m

 $\mathfrak{M}$ message space M encoded message M[i]*i*-th encoded message

 $M_i[i]$ *i*-th encoded message on Prover<sub>i</sub>

 $\mathcal{M}^*$ encoded message space number of provers n N codeword length

 $\mathcal{P}_i$ proving algorithm of *i*-th Prover qorder of underlying finite field

response

 $r^{M}$ response vector for encoded message M

S tag (in a keyed scheme)

succ(𝒫) success probability of proving algorithm

 $\mathcal{R}^*$ response code Γ challenge space

number of possible challenges y

Δ response space number of users ρ privacy threshold τ

Acknowledgment: Thanks to Andris Abakuks and Simon Skene for some helpful discussions of statistics.

### References

- [1] G. Ateniese, R. C. Burns, R. Curtmola, J. Herring, O. Khan, L. Kissner, Z. N. J. Peterson and D. Song, Remote data checking using provable data possession, ACM Trans. Inform. Sys. Security 14 (2011), Paper No. 12.
- G. Ateniese, R. C. Burns, R. Curtmola, J. Herring, L. Kissner, Z. N. J. Peterson and D. X. Song, Provable data possession at untrusted stores, in: Proceedings of the 14th ACM Conference on Computer and Communications Security, ACM, New York (2007), 598-609.

- [3] G. Ateniese, Ö. Dagdelen, I. Damgård and D. Venturi, Entangled cloud storage, IACR Cryptology ePrint Archive (2012), https://eprint.iacr.org/2012/511.pdf.
- [4] G. Ateniese, R. Di Pietro, L. V. Mancini and G. Tsudik, Scalable and efficient provable data possession, in: Proceedings of the 4th International Conference on Security and Privacy in Communication Networks, ACM, New York (2008), 1-9.
- [5] G. Ateniese, S. Kamara and J. Katz, Proofs of storage from homomorphic identification protocols, in: Advances in Cryptology—ASIACRYPT 2009, Springer, Berlin (2009), 319-333.
- [6] G. R. Blakley, Safeguarding cryptographic keys, in: Proceedings of the National Computer Conference, AFIPS, New York (1979), 313-317.
- [7] G. R. Blakley and C. Meadows, Security of ramp schemes, in: Advances in Cryptology—CRYPTO 1985, Springer, Berlin (1985), 242-268.
- [8] K. D. Bowers, A. Juels and A. Oprea, Proofs of retrievability: Theory and implementation, in: Proceedings of the 2009 ACM Workshop on Cloud Computing Security, ACM, New York (2009), 43-54.
- [9] R. Curtmola, O. Khan, R. C. Burns and G. Ateniese, MR-PDP: Multiple-replica provable data possession, in: The 28th International Conference on Distributed Computing Systems, IEEE Press, Piscataway (2008), 411–420.
- [10] Y. Dodis, S. P. Vadhan and D. Wichs, Proofs of retrievability via hardness amplification, in: Theory of Cryptography, Springer, Berlin (2009), 109-127.
- [11] A. Juels and B. S. Kaliski, Jr., PORs: Proofs of retrievability for large files, in: Proceedings of the 14th ACM Conference on Computer and Communications Security, ACM, New York (2007), 584-597.
- [12] S. Kamara and K. Lauter, Cryptographic cloud storage, in: Financial Cryptography and Data Security, Springer, Berlin (2010), 136-149.
- [13] R. J. McEliece and D. V. Sarwate, On sharing secrets and Reed-Solomon codes, Comm. ACM 24 (1981), 583-584.
- [14] M. B. Paterson and D. R. Stinson, A simple combinatorial treatment of constructions and threshold gaps of ramp schemes, Cryptogr. Commun. 5 (2013), 229-240.
- [15] M. B. Paterson, D. R. Stinson and J. Upadhyay, A coding theory foundation for the analysis of general unconditionally secure proof-of-retrievability schemes for cloud storage, J. Math. Cryptol. 7 (2013), 183-216.
- [16] H. Shacham and B. Waters, Compact Proofs of Retrievability, in: Advances in Cryptology—ASIACRYPT 2008, Springer, Berlin (2009), 90-107.
- [17] A. Shamir, How to share a secret, *Comm. ACM* **22** (1979), 612–613.
- [18] K. Ulm, Simple method to calculate the confidence interval of a standardized mortality ratio (SMR), Amer. J. Epidemiology **131** (1990), 373-375.
- [19] C. Wang, Q. Wang, K. Ren and W. Lou, Privacy-preserving public auditing for data storage security in cloud computing, in: IEEE Proceedings INFOCOM 2010, IEEE Press, Piscataway (2010), 1-9.
- [20] Summary of the Amazon S3 Service Disruption in the Northern Virginia (US-EAST-1) Region, https://aws.amazon.com/message/41926/.
- [21] Why is decentralized and distributed file storage critical for a better web?, https://coincenter.org/entry/why-is-decentralized-and-distributed-file-storage-critical-for-a-better-web.