

Proofs of Retrievability with Public Verifiability and Constant Communication Cost in Cloud

Jiawei Yuan
University of Arkansas at Little Rock
Little Rock, USA
jxyuan@ualr.edu

Shucheng Yu
University of Arkansas at Little Rock
Little Rock, USA
sxyu1@ualr.edu

ABSTRACT

For data storage outsourcing services, it is important to allow data owners to efficiently and securely verify that the storage server stores their data correctly. To address this issue, several proof-of-retrievability (POR) schemes have been proposed wherein a storage server must prove to a verifier that all of a client's data are stored correctly. While existing POR schemes offer decent solutions addressing various practical issues, they either have a non-trivial (linear or quadratic) communication complexity, or only support private verification - only the data owner can verify the remotely stored data. It remains open to design a POR scheme that achieves both public verifiability and constant communication cost simultaneously.

In this paper, we solve this open problem and propose the first POR scheme with public verifiability and constant communication cost. In our proposed scheme, the message exchanged between the prover and verifier is composed of a constant number of the underlying group elements. Different from existing private POR constructions, our scheme allows public verification and releases the data owners from burden of being staying online. Thorough analysis shows that our proposed scheme is efficient and practical. We prove the security of our scheme based on Computational Diffie-Hellman Problem, Strong Diffie-Hellman Assumption and Bilinear Strong Diffie-Hellman Assumption.

Categories and Subject Descriptors

H.3.2 [Information Storage and Retrieval]: Information Storage; D.4.6 [Security and Protection]: Cryptographic controls

General Terms

Storage, Integrity, Security, Algorithm

Keywords

Proofs of Retrievability, Cloud Storage, Public Verification,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

Integrity Check, Constant Communication, Polynomial Commitment

1. INTRODUCTION

Due to a number of unprecedented advantages: resource elasticity, on-demand self-service, location independent resource polling, etc[22], cloud storage is increasingly attracting organizations and individuals to outsource their local data to it. Currently, millions of users have been involved into cloud storage services including Amazon S3, Microsoft Skydrive, Google Cloud Storage and Dropbox. It is clear that cloud storage is a raising trend to become a pervasive service.

Despite the proliferation of cloud storage service, there are also increasing security concerns caused by the shortages of cloud servers. One significant concern is the data integrity, i.e., whether the cloud server indeed stores its clients' data correctly. Recently, a number of data loss events have been reported for those best-known storage providers[2, 1, 3, 4], including Amazon S3, Dropbox. To ensure the clients' confidence for the integrity of data previously stored on the cloud storage server, a reliable proof-of-retrievability (POR)[16] system is desirable. Specifically, in a POR system, the data storage server must prove that he actually stores the clients' data correctly. At the same time, the client can verify whether the proof generated by the server is valid. For the construction of a POR system, there are several aspects that mainly affect the performance of the system: 1) communication cost between server and client in a verification; 2) private verifiability or public verifiability; 3) storage overhead on server side and client side; 4) computation cost of server and client for each verification. Although several POR techniques[20, 11, 15] have been proposed, different limitations are existed in them, especially for the communication cost and the support of public verification.

By utilizing the idea of aggregating block integrity values, Shacham and Waters (SW)[20] proposed two fast POR schemes. In their POR constructions, one scheme supports the private POR verification and the other one achieves the public POR verification. However, both public and private SW schemes suffer from the communication cost, which is linear to the number of elements in the data block. For the integrity check of large data files, this kind of communication cost will obviously introduce high latency to the schemes and even makes them unaffordable in some bandwidth limited scenarios(e.g., mobile devices paid by limited data plan). To overcome the limitation in Ref.[20], Xu et al.[15] construct a private POR scheme with constant com-

munication cost by utilizing a recent proposed polynomial commitment technique. The main drawback of this scheme is that it only supports private verification. This kind of verification will cause heavy burden on both communication and computation to the data owner if he wants to share his data with other clients. Because the owner has to stay online and perform verification for every integrity check request from other clients. In Ref.[10], Dodis et al. enhance SW schemes by reducing the challenge message size, but no change is made to the response size, which is linear to the number of elements in the data block. To our best knowledge, there is no existing solution proposed for the public *POR* with constant communication cost. It is still an open problem as mentioned in Ref.[15].

In this paper, we solve this problem and design a secure and practical *POR* scheme with constant communication cost and public verifiability, called *PCPOR*. The main idea of our scheme can be summarized as follows: a data owner first breaks an ensure coded file into n blocks $\{m_i\}, 1 \leq i \leq n$ and generates a corresponding authentication tag $\{\sigma_i\}$ for each block. Then, all the data blocks and tags are outsourced to the cloud storage server. When a client wants to retrieve the data from the server and check whether the data are stored correctly, he can generate a challenge message and send it to the cloud server. By receiving the challenge message, the cloud server generates the proof for the correct data storage based on the message, the public key information and the previously stored tags, and then returns it to the client as response. After receiving the proof response, the client performs the verification algorithm to check data integrity using the public key information. In our *PCPOR* scheme, any client or auditing third-party can perform verification process without contacting the data owner, who can keep off-line after outsourcing his data. By tailoring a constant size polynomial commitment technique[17], our proposed scheme aggregates the proof information into a polynomial, which makes the size of proof information independent to the number of elements in data block and achieves better performance in storage cost than the existing *POR* schemes. Through the unique incorporation of techniques from Dotis et al.[10], our *PCPOR* reduces the challenge message complexity to constant level. Our thorough analysis shows that our proposed scheme is efficient and practical. Based on Computational Diffie-Hellman Problem (CDH), Strong Diffie-Hellman (SDH) Assumption and Bilinear Strong Diffie-Hellman (BSDH) Assumption, we prove that our scheme is secure.

Main contributions of this paper can be summarized as below.

- We construct the first secure and efficient *POR* scheme with constant communication cost and public verifiability.
- Our proposed scheme solves the open problem pointed out in Ref.[15] and omits the trade-off between communication cost and storage cost in the public SW scheme[20].
- We formally prove the security of our scheme. The advantages of our proposed scheme are validated via analysis.

The rest of this paper is organized as follows: In section 2, We review and discuss the related works. Section 3 describes

the system model, security model and assumptions. We introduce the technique preliminaries of our work in Section 4, which is followed by the construction and security proof of our proposed scheme in Section 5. We evaluate the performance of our scheme in Section 6. Discussions about our paper are provided in Section 7 and we conclude the paper in Section 8

2. RELATED WORK

In Ref.[16], Juels et al. first defined the *POR* formally, which allows a storage server to convince a client that it can correctly retrieve a file previously stored at the server. In their proposed *POR* scheme, disguised blocks hidden among regular file blocks are utilized to detect data modified by the server. However, the communication cost in this scheme is inevitably large, i.e. $O(s\lambda^2)$, where λ is the security parameter and s is the number of elements in each erasure coded file block. This is because there are $O(\lambda)$ file blocks in the proof response, each of which has a size of $O(\lambda)$. What is more, the number of challenges supported by this scheme is a fixed priori and thus limits its application. With similar purpose of Ref.[16], Ateniese et al.[5] proposed an efficient but weaker provable data possession model using homomorphic authentication tag. However, an adversary was later introduced for this scheme by Shacham and Waters[20], which can answer a fraction of queries correctly with non-negligible probabilities.

To omit the limitation in Juels et al.'s *POR* scheme[16], Shacham and Waters (SW) [20] proposed two fast *POR* schemes based on the homomorphic linear authenticators[5], which enables the storage server to reduce the proof complexity by aggregating the authentication tags of individual file blocks. In their constructions, one proposed *POR* schemes supports private verification based on pseudorandom functions (PRFs) [12] and the other one achieves public verifiability by utilizing BLS signatures[7]. Compared to scheme in Ref.[16], the communication cost for proof response in Ref.[20] is reduced from $O(s\lambda^2)$ to $O(s\lambda)$ and can support unlimited number of challenges, where s is the number of elements in each erasure coded file block. At the same time, they first provide a security proof against arbitrary adversaries in the formal *POR* model. However, in SW schemes, the communication size for proof response is still linear to the value of s and the challenge cost is increased to $O(\lambda^2)$. When considering the integrity verification of large data files, the bandwidth consuming will introduce high delay, which can be worse and even unaffordable in some scenarios with restricted bandwidth, such as mobile devices paid by limited data plan.

Following SW schemes[20], several *POR* schemes are proposed recently to enhance it in terms of communication cost. In Ref.[10], by using a (γ, δ) -*hitter* introduced by Goldreich[18], Dodis et.al. reduce the size of challenge message in SW scheme from $O(\lambda^2)$ to $O(\lambda)$. Nevertheless, there is still no change made in this scheme to the response size, which is linear to the number of elements in a data block. To further improve *POR* scheme and overcome the limitations in previous ones, Xu et.al[15] proposed a private *POR* scheme with constant communication cost based on a recent proposed polynomial commitment technique[17]. By aggregating the proofs into a polynomial, the limitation of communication cost in SW's private *POR* scheme is omitted in Xu et.al.'s construction[15]. The main drawback of this

proposed scheme is that the data owner has to stay online and perform onerous works of verification for the each request of integrity check from other clients. These tasks will inevitably introduce heavy communication cost and computation cost to the data owner, especially when multiple clients submit integrity check at the same time. How to support *POR* scheme with public verifiability and constant communication cost is pointed out as an open problem.

3. MODEL AND ASSUMPTION

3.1 System Model

In this work, we follow the *POR* model that is consistent with most existing *POR* schemes[16, 20, 8, 15]. We consider a *POR* system that has three participating entities: *Data owner*, *Client*, and *Cloud server*. Data owner has a collection of data and stores them on cloud server after ensure coding together with the corresponding authentication tags. The client who shares the stored data with the owner can access it with integrity check, which can also be performed as an independent procedure. To check the integrity of data files, the client generates a challenge message and sends it to the cloud server. The cloud server then responses the computed proof for the selected file blocks to the client. After receiving the proof, the client can verify the integrity of data files through the verification algorithm. W.l.o.g., we define the 1-round version of our *POR* model, which contains four algorithms, *KeyGen*, *Setup*, *Prove* and *Verify*, as below:

- **KeyGen:** Given a selected security parameter λ , the randomized *KeyGen* algorithm outputs the system public key and private key as (PK, SK) .
- **Setup:** Given a data file $M \in \{0, 1\}^*$ and the public-private key pair (PK, SK) , the *Setup* algorithm generates the encoded file \hat{M} as well as the corresponding authentication tag σ , which will be stored on the server.
- **Prove:** Given the public key PK , encoded file \hat{M} , authentication tag σ and a challenge message *Chall*, the *Prove* algorithm produces a proof response *Prf*.
- **Verify:** Given the public key PK and the *Prf*, the *Verify* algorithm checks the data integrity and outputs result as either accept or reject.

3.2 Security Model

We consider the storage server as untrusted and potentially malicious, which is consistent with existing *POR* schemes [16, 20, 15]. In our construction, we would like our *POR* scheme to be correct and sound. For the correctness of our scheme, we require that our *Verify* algorithm accepts a valid proof generated from all key pairs (PK, SK) , all files $M \in \{0, 1\}^*$, all encoded files \hat{M} and authentication tags σ . For the soundness of our scheme, if any malicious cloud server can generate a proof and convince the *Verify* algorithm that it actually stores \hat{M} correctly, it has to yield up the right \hat{M} for the proof generation. By following Ref.[16, 20, 15], we define the security game for the soundness of our *POR* scheme as below.

DEFINITION 3.1. Let $\nabla = (KeyGen, Setup, Prove, Veiry)$ be a *POR* scheme and A be a probabilistic polynomial-time

adversary. Consider the following security game among a trust authority, a challenger and A .

- The trust authority runs $KeyGen(1^\lambda) \rightarrow (PK, SK)$ and gives PK to the adversary A .
- The adversary A chooses a data file M and sends it to the trust authority. The authority then runs $Setup(M, SK, PK) \rightarrow (\sigma, \hat{M})$ and responses the encoded data file \hat{M} together with the authentication tag σ back to A .
- With regard to the data file M chosen by adversary A , the challenger picks a random challenge *Chall* and sends it to A .
- According to the received challenge *Chall*, the adversary A produces a proof response *Prf* by running an arbitrary algorithm $Art(\hat{M}, \sigma, PK) \rightarrow Prf$ rather than the *Prove* algorithm. The proof response *Prf* is sent back to the challenger.
- The challenger verifies *Prf* by running *Verify* algorithm with *Prf* and PK . The output of $Verify(Prf, PK)$ is denoted as *Rst*.
- The adversary wins the game if and only if he can produce a *Prf* with data file \hat{M}' , $\hat{M}' \neq \hat{M}$ and make the challenger generate *Rst* as *accept* through the *Verify* algorithm.

We say that ∇ is sound if any probabilistic polynomial-time adversary A can win the game with at most a negligible probability.

3.3 Assumption

DEFINITION 3.2. **Computational Diffie-Hellman (CDH) Problem**[9]

Let $a, b \xleftarrow{R} Z_p^*$. Given input as (g, g^a, g^b) , where g is a generator of a cyclic group G of order p . It is computationally intractable to compute the value g^{ab} .

DEFINITION 3.3. **t -Strong Diffie-Hellman (t -SDH) Assumption**[6]

Let $\alpha \xleftarrow{R} Z_p^*$. Given input as a $(t+1)$ -tuple $(g, g^\alpha, \dots, g^{\alpha^t}) \in G^{t+1}$, where G is a cyclic group of prime order p and g is the generator of G . For any probabilistic polynomial time adversary (PPT_{Adv}), the probability $Pr[PPT_{Adv}(g, g^\alpha, \dots, g^{\alpha^t}) = (c, g^{\frac{1}{\alpha+c}})]$ is negligible for any value of $c \in Z_p^* / -\alpha$.

DEFINITION 3.4. **t -Bilinear Strong Diffie-Hellman (t -BSDH) Assumption**[13]

Let $\alpha \xleftarrow{R} Z_p^*$. Given input as a $(t+1)$ -tuple $(g, g^\alpha, \dots, g^{\alpha^t}) \in G^{t+1}$, where G is a multiplicative cyclic group of prime order p and g is the generator of G . For any probabilistic polynomial time adversary (PPT_{Adv}), the probability $Pr[PPT_{Adv}(g, g^\alpha, \dots, g^{\alpha^t}) = (c, e(g, g)^{\frac{1}{\alpha+c}})]$ is negligible for any value of $c \in Z_p^* / -\alpha$.

4. TECHNIQUE PRELIMINARIES

4.1 Bilinear Map

A bilinear map[7] is a map $e : G \times G \rightarrow G_1$, where G and G_1 are two multiplicative cyclic groups of the same prime order p . A bilinear map has the following properties:

- **Bilinear:** For all $g_1, g_2 \in G$ and $a, b \xleftarrow{R} Z_p^*$, $e(g_1^a, g_2^b) = e(g_1, g_2)^{ab}$.
- **Computable:** There exists a computable algorithm that can compute e efficiently.
- **Non-degenerate:** For $g \in G$, $e(g, g) \neq 1$

4.2 Constant Size Polynomial Commitment

A secure polynomial commitment scheme enables a committer to commit to a polynomial with a short string, which can be used by a verifier to confirm claimed evaluations of the committed polynomial. By utilizing an algebraic property of polynomials $f(x) \in Z[x]$: $(x-r)$ perfectly divides the polynomial $f(x) - f(r)$, $r \xleftarrow{R} Z_p^*$, Kate et.al.[17] proposed a polynomial commitment scheme with constant communication size. In their scheme, a committer of polynomial $f(x)$ can generate a proof with constant size to verify the correctness of the polynomial evaluation $f(r)$, where $x = r$ is an index on the polynomial. Specifically, we summarize Kate et.al.'s scheme[17] as below.

- **Setup($1^\lambda, s$):** Given a security parameter λ and a fixed number s , a trust authority generates a public-private key pair (PK, SK) :

$$PK = (G, G_1, g, g^\alpha, \dots, g^{\alpha^{s-1}}) \quad SK = \alpha \xleftarrow{R} Z_p^*$$

where G and G_1 are two multiplicative cyclic groups with prime order p (λ bits security), g is the generator of G and $e : G \times G \rightarrow G_1$.

- **Commit($PK, f_{\vec{m}}(x)$):** For polynomial $f_{\vec{m}}(x) \in Z_p[x]$ with coefficient vector $\vec{m} = (m_0, m_1, \dots, m_{s-1}) \xleftarrow{R} Z_p^*$, a committer computes the commitment $C = g^{f_{\vec{m}}(\alpha)} \in G$ and publishes C .
- **CreateWitness($PK, f_{\vec{m}}(x), r$):** For any index $r \xleftarrow{R} Z_p^*$, the committer divides polynomial $f_{\vec{m}}(x) - f_{\vec{m}}(r)$ with $(x-r)$ and outputs \vec{w} as result, where $\vec{w} = (w_0, w_1, \dots, w_{s-1})$ and $f_{\vec{w}}(x) \equiv \frac{f_{\vec{m}}(x) - f_{\vec{m}}(r)}{(x-r)}$. Then, the witness ψ is computed as $\psi = \prod_{j=0}^{s-1} (g^{\alpha^j})^{w_j} = g^{f_{\vec{w}}(\alpha)}$ based on PK .
- **VerifyEval($PK, C, r, f_{\vec{m}}(r), \psi$):** A verifier verifies that $f_{\vec{m}}(r)$ is the evaluation at the index r of the polynomial committed to by C as:

$$e(C, g) \stackrel{?}{=} e(\psi, g^\alpha / g^r) \cdot e(g, g)^{f_{\vec{m}}(r)}$$

Due to the space limitation, please refer to Ref.[17] for the detail correctness and security proofs of this scheme.

5. CONSTRUCTION OF PCPOR

5.1 Scheme Description

Let $e : G \times G \rightarrow G_1$ and H be the one-way hash function[14], where G is a multiplicative cyclic group of prime order p and g be a generator of G . We define $f_{\vec{c}(x)}$ as a polynomial with coefficient vector $\vec{c} = (c_0, c_1, \dots, c_{s-1})$ and describe our PCPOR scheme as follows.

- **KeyGen(1^λ) \rightarrow (PK, SK):**

Choose a random $(\lambda+1)$ bits safe prime p and generate a random signing keypair $((spk, ssk) \xleftarrow{R} SKg)$ using BLS signature[6]. Choose two random numbers $\alpha, \epsilon \xleftarrow{R} Z_p^*$ and compute $v \leftarrow g^\epsilon$, $\kappa \leftarrow g^{\alpha\epsilon}$ as well as $\{g^{\alpha^j}\}_{j=0}^{s-1}$. Then, the public and private keys are

$$PK = \{p, v, \kappa, spk, \{g^{\alpha^j}\}_{j=0}^{s-1}\}, \quad SK = \{\epsilon, ssk, \alpha\}$$

- **Setup(PK, SK, M) \rightarrow (M^*, σ, τ):**

Given a data file M , obtain M' by applying error correction code(e.g.Reed-Solomon code[19]). Split M' into n blocks, each of which is s elements long: $\{m_{ij}\}_{1 \leq i \leq n, 0 \leq j \leq s-1}$. Choose a random file name $name$ from some sufficiently large domain(e.g., Z_p^*). Let τ_0 be " $name||n$ "; the file tag τ is τ_0 together with a signature on τ_0 under ssk : $\tau \leftarrow \tau_0 || SSig_{ssk}(\tau_0)$. For each data block i , $1 \leq i \leq n$, an authentication tag is computed as:

$$\begin{aligned} \sigma_i &= (g^{H(name||i)} \cdot \prod_{j=0}^{s-1} g^{m_{ij}\alpha^j})^\epsilon \quad (1) \\ &= (g^{H(name||i)} \cdot g^{f_{\vec{\beta}_i}(\alpha)})^\epsilon \end{aligned}$$

where $\beta_{i,j} = m_{i,j}$ and $\vec{\beta}_i = \{\beta_{i,0}, \beta_{i,1}, \dots, \beta_{i,s-1}\}$. The processed file M^* together with the authentication tags σ_i are outsourced for storage, where M^* is $\{m_{i,j}\}, 1 \leq i \leq n, 0 \leq j \leq s-1$.

- **Verify(PK, τ) \rightarrow Chall:**

Stage 1: Verify the signature on τ : if the signature is not valid, reject and halt; otherwise, parse τ to recover $name$ and n . Now, choose a random l -elements subset I of the set $[1, n]$ and a random number $\rho \xleftarrow{R} Z_p^*$. Produce the challenge message as

$$Chall = \{r, \rho, I\}$$

where $r \xleftarrow{R} Z_p^*$. Challenge the data storage server with $Chall$.

- **Prove($PK, M^*, Chall, \tau$) \rightarrow (ψ, y, σ):**

Parse the challenge message $Chall$ as $\{r, \rho, I\}$ and generate a l -element set $Q = (i, v_i)$, where $i \in I, v_i = \rho^i \bmod p$. Based on the processed file $M^* = m_{i,j}, 1 \leq i \leq n, 0 \leq j \leq s-1$ and σ_i , compute

$$\sigma = \prod_{(i, v_i) \in Q} \sigma_i^{v_i} \quad (2)$$

We denote vector

$$\vec{A} = \left(\sum_{(i, v_i) \in Q} v_i m_{i,0}, \dots, \sum_{(i, v_i) \in Q} v_i m_{i,s-1} \right)$$

Compute

$$y = f_{\vec{A}}(r) \quad (3)$$

As we mentioned before, polynomials $f(x) \in Z[x]$ have the algebraic property that $(x-r)$ perfectly divides the polynomial $f(x) - f(r)$, $r \xleftarrow{R} Z_p^*$. Now, divide the polynomial $f_{\bar{A}}(x) - f_{\bar{A}}(r)$ with $(x-r)$ using polynomial long division, and denote the coefficients vector of the resulting quotient polynomial as $\bar{w} = (w_0, w_1, \dots, w_{s-1})$, that is, $f_{\bar{w}}(x) \equiv \frac{f_{\bar{A}}(x) - f_{\bar{A}}(r)}{x-r}$. Produce

$$\psi = \prod_{j=0}^{s-1} (g^{\alpha^j})^{w_j} = g^{f_{\bar{w}}(\alpha)} \quad (4)$$

Response $Prf = \{\psi, y, \sigma\}$.

- **Verify**(PK, Prf) $\rightarrow Rst$:

Stage 2: After receiving the proof response Prf , compute

$$\eta_i = v^{-H(name||i)v_i}, (i, v_i) \in Q \quad (5)$$

$$\eta = \prod_{(i, v_i) \in Q} \eta_i \quad (6)$$

where $v = g^\epsilon$ in PK . Parse Prf as $\{\psi, y, \sigma\}$ and check

$$e(\psi, \kappa \cdot v^{-r}) \stackrel{?}{=} e(\sigma \cdot \eta, g) \cdot e(g^{-y}, v) \quad (7)$$

where $\kappa = g^{\epsilon\alpha}$ in PK . If Eq.7 holds, then output Rst as accept; otherwise, output Rst as reject.

- **Correctness** For a storage server who honestly responds the challenge with a $Prf = \{\psi, y, \sigma\}$, we can analyze the correctness of Eq.7 from the left part and right part as:

Left Part:

$$\begin{aligned} & e(\psi, \kappa \cdot v^{-r}) \quad (8) \\ &= e(g^{f_{\bar{w}}(\alpha)}, g^{\epsilon(\alpha-r)}) \\ &= e(g, g)^{\frac{f_{\bar{A}}(\alpha) - f_{\bar{A}}(r)}{\alpha-r} \cdot \epsilon(\alpha-r)} \\ &= e(g, g)^{\epsilon(f_{\bar{A}}(\alpha) - f_{\bar{A}}(r))} \end{aligned}$$

Right Part:

$$\begin{aligned} & e(\sigma \cdot \eta, g) \cdot e(g^{-y}, v) \quad (9) \\ &= e(g^{\epsilon(\sum_{(i, v_i) \in Q} u_i + f_{\bar{A}}(\alpha)) + \epsilon \sum_{(i, v_i) \in Q} -u_i}, g) \cdot \\ & e(g^{-f_{\bar{A}}(r)}, g^\epsilon) \\ &= e(g, g)^{\epsilon f_{\bar{A}}(\alpha)} \cdot e(g, g)^{-\epsilon f_{\bar{A}}(r)} \\ &= e(g, g)^{\epsilon(f_{\bar{A}}(\alpha) - f_{\bar{A}}(r))} \end{aligned}$$

where $u_i = H(name||i)v_i$, $v = g^\epsilon$ and $\kappa = g^{\epsilon\alpha}$. From the above Eq.8 and Eq.9, it is easy to see that the scheme is correct if the storage sever generate the Prf honestly.

5.2 Security Proof

5.2.1 The underlying authenticator is unforgeable

THEOREM 5.1. *If t -SDH Assumption holds and $g^{f_{\bar{c}}(x)}$ can be forged by an existed probabilistic polynomial time adversary A , we can construct an algorithm B that uses A to efficiently compute the solution to t -SDH problem.*

PROOF. Suppose there exists a probabilistic polynomial time adversary A can forge $f_{\bar{c}}^1(\alpha)$ such that $g^{f_{\bar{c}}^1(\alpha)} = g^{f_{\bar{c}}(\alpha)}$, where \bar{c} is the coefficient vector, he/she obtains $g^{f_{\bar{c}}^2(\alpha)} = g^{f_{\bar{c}}(\alpha)} / g^{f_{\bar{c}}^1(\alpha)} = g^{f_{\bar{c}}(\alpha) - f_{\bar{c}}^1(\alpha)} \in Z_p[x]$. Since $f_{\bar{c}}^1(\alpha) = f_{\bar{c}}(\alpha)$ and $f_{\bar{c}}^2(\alpha) = 0$, i.e., α is a root of polynomial $f_{\bar{c}}^2(x)$ and by factoring $f_{\bar{c}}^2(x)$ [21], B can easily find $SK = \alpha$ and solve the instance of the t -SDH problem given by the system parameters. \square

THEOREM 5.2. *If the signature scheme used for file tags is existentially unforgeable, CDH Problem is hard, t -SDH Assumption and t -BSDH Assumption hold. In our proposed scheme, the prover's response (y, ψ, σ) is unforgeable.*

PROOF. Suppose a probabilistic polynomial time adversary can generate a response (y', ψ', σ') to forge (y, ψ, σ) after receiving a challenge message from the verifier, $(y', \psi', \sigma') \neq (y, \psi, \sigma)$. Since both the forged and the true responses can be accepted by the verification algorithm, we can get the following two equations:

$$e(\psi, \kappa \cdot v^{-r}) = e(\sigma \cdot \eta, g) \cdot e(g^{-y}, v) \quad (10)$$

$$e(\psi', \kappa \cdot v^{-r}) = e(\sigma' \cdot \eta, g) \cdot e(g^{-y'}, v) \quad (11)$$

Dividing Eq.10 with Eq.11, we obtain:

$$\frac{e(\psi, g)^{\epsilon(\alpha-r)}}{e(\psi', g)^{\epsilon(\alpha-r)}} = \frac{e(g, g)^{\epsilon E - \sum_{(i, v_i) \in Q} H(name||i)v_i - y}}{e(g, g)^{\epsilon E' - \sum_{(i, v_i) \in Q} H(name||i)v_i - y'}}$$

$$\left(\frac{e(\psi, g)}{e(\psi', g)} \right)^{\epsilon(\alpha-r)} = e(g, g)^{\epsilon(E-E') + y' - y} \quad (12)$$

where we denote σ as $g^{E\epsilon}$ and σ' as $g^{E'\epsilon}$ for simplicity.

Now we do a case analysis on whether $\sigma = \sigma'$.

Case 1: $\sigma \neq \sigma'$. Since $g^{E\epsilon} = \sigma$ and $g^{E'\epsilon} = \sigma'$, we get $E \neq E'$. As $\left(\frac{e(\psi, g)}{e(\psi', g)} \right)^{\epsilon(\alpha-r)}$, $e(g, g)^{y' - y}$ and $e(g, g)^{E\epsilon} = e(g, \sigma)$ are known to the adversary, we rewrite Eq.12 as

$$\Upsilon = e(g, g)^{\epsilon E'} \quad (13)$$

where we denote $\Upsilon = e(g, g)^{E\epsilon + y' - y} / \left(\frac{e(\psi, g)}{e(\psi', g)} \right)^{\epsilon(\alpha-r)}$ as a known value to the adversary.

Recall that in this proof, the CDH Problem is hard. If the any probabilistic polynomial time adversary A can find E' with non-negligible probability and make the Eq.13 hold, we can construct an algorithm B that uses A to solve the instance of CDH Problem. Specifically, given $E' \neq E$ found by A , which makes Eq.13 hold, B can get $\Upsilon = e(g, g)^{\epsilon E'}$ as solution for CDH Problem of $e(g, g)^\epsilon$ and $e(g, g)^{E'}$. Therefore, no probabilistic polynomial time adversary can find a valid forged response $(y, \psi, \sigma) \neq (y', \psi', \sigma')$ and $\sigma \neq \sigma'$ with non-negligible probability.

Case 2: $\sigma = \sigma'$. In this case, we can rewrite Eq.12 as:

$$\left(\frac{e(\psi, g)}{e(\psi', g)} \right)^{\epsilon(\alpha-r)} = e(g, g)^{y' - y} \quad (14)$$

We further do a case analysis on whether $y = y'$.

Case 2.1: $y = y'$. As $(y, \psi, \sigma) \neq (y', \psi', \sigma')$, $\sigma = \sigma'$ and $y = y'$, we can obtain that $\psi \neq \psi'$. In this case, since $y = y'$,

we further rewrite the Eq.14 as:

$$\left(\frac{e(\psi, g)}{e(\psi', g)} \right)^{\epsilon(\alpha-r)} = 1 \quad (15)$$

We know that $\psi \neq \psi'$, i.e., $\frac{e(\psi, g)}{e(\psi', g)} \neq 1$, and $\epsilon \neq 0$, we can infer $\alpha = r$ from Eq.15. In our scheme, r is known to the adversary (i.e., the adversary can find $SK = \alpha$). As we mentioned in Theorem 5.1, by finding $SK = \alpha$, we can solve the instance of the t-SDH problem by given the system parameters. Therefore, no valid forged response $(y, \psi, \sigma) \neq (y', \psi', \sigma')$ and $y = y'$ can be found by probabilistic polynomial time adversary A with non-negligible probability.

Case 2.2: $y \neq y'$. From Eq.14 and $y \neq y'$, we can imply that $\alpha \neq r$. In this case, we show how to construct an algorithm B , using the existed adversary, that can break the t-BSDH Assumption with a valid solution $(-r, \left(\frac{e(\psi, g)}{e(\psi', g)}\right)^{\frac{1}{y'-y}})$.

We denote ψ as g^θ and ψ' as $g^{\theta'}$, and then we can rewrite Eq.14 as :

$$\begin{aligned} \left(\frac{e(\psi, g)}{e(\psi', g)} \right)^{\epsilon(\alpha-r)} &= \frac{e(g, g)^{-y}}{e(g, g)^{-y'}} \\ \theta\epsilon(\alpha-r) + y &= \theta'\epsilon(\alpha-r) + y' \\ \frac{\epsilon(\theta - \theta')}{y' - y} &= \frac{1}{\alpha - r} \end{aligned} \quad (16)$$

Therefore, algorithm B can compute

$$\left(\frac{e(\psi, v)}{e(\psi', v)} \right)^{\frac{1}{y'-y}} = e(g, g)^{\frac{\epsilon(\theta-\theta')}{y'-y}} = e(g, g)^{\frac{1}{\alpha-r}} \quad (17)$$

and return $(-r, e(g, g)^{\frac{1}{\alpha-r}})$ as a solution for t-BSDH instance. It is easy to see that the success probability of solving the instance is the same as the success probability of the adversary, and the time required is a small constant larger than the time required by the adversary.

Therefore, *Theorem 5.2* is proved. \square

6. PERFORMANCE EVALUATION

In this section, we numerically evaluate the performance our proposed *PCPOR* scheme in terms of communication cost, computation cost and storage cost. We compare our *PCPOR* scheme with the existing *POR* techniques[20, 10, 15] and summarize the result in Table 1. For simplicity, in the following part of this paper, we denote the complexity of one multiplication operation on Group G as MUL and that of one exponentiation operation on Group G as EXP¹. Furthermore, we use ZADD and ZMUL to represent the addition and multiplication operations on Z_p^* respectively. PRF is used to denote pseudorandom function and $\|G\|$ denotes the size(in number of bits) of a group element on G .

6.1 Communication

In our proposed *PCPOR* scheme, the communication cost comes from the challenge message *Chall* and the proof response *Prf* in each verification request. For the challenge

¹When the operation is on the elliptic curve, EXP means scalar multiplication operation and MUL means one point addition operation.

message, it consists of a l -elements subset I and two random elements $\rho, r \xleftarrow{R} Z_p^*$. By utilizing Goldreich[12]'s (γ, δ) -*hitter*², we can represent the subset I with $\log^{|F|} + 3\log^{(1/\delta)}$ bits, where $|F|$ is the size of the encoded data file, and δ is the error probability. In particular, given a data file less than 1024 TB (i.e. $|F| = 2^{43}$ bits, which is enough for most practical scenarios) and set $\delta = 2^{-80}$ same as in Ref.[15], the subset I can be represented with 283 bits. As a result, the total cost of *Chall* in our *PCPOR* scheme is $2\lambda + 283$ bits when the data file size is smaller than 1024 TB. For the proof response, by aggregating proof information into 3 elements ψ, σ and y , the total size of the proof response is $2\|G\| + \lambda$ bits, where ψ and σ are two group elements and y the result of a polynomial. Therefore, the total complexity of communication cost in our *PCPOR* scheme is $O(\|G\|)$. When setting reasonable parameters for our *PCPOR* scheme (e.g. $\lambda = 160$, $\|G\| = 1024$ bits and the data file is less than 1024 TB), the total communication cost becomes constant size.

Now, we compare the existing *POR* schemes [20, 10, 15] with our *PCPOR* scheme and summarize the result in Table 1. In Ref.[20], the complexity of challenge message and proof response are $O(\lambda^2)$ and $O(s\lambda)$ respectively, where s is the number of elements in each encoded block. Compared to our *PCPOR* scheme, which has constant communication cost in for both challenge and response processes, this scheme have a communication size of proof response linear to s and a challenge message cost λ times to ours. In the *POR* scheme proposed by Dodis et al.[10], the size of challenge message is reduced to $O(\lambda)$, which is the same as in our *PCPOR* scheme. However, the cost of proof response in their scheme is still $O(s\lambda)$. Considering only private verification, Xu et al.[15]'s *POR* scheme reduces communication cost to constant size in each single verification. Nevertheless, it is worthy of notice that their scheme only supports private verification, which centralizes all the verification tasks to the data owner and thus introduces onerous burden in communication to it. Specifically, if the data owner wants to share his outsourced data with other individuals/organizations, it has to stay online and processes all verifications by himself. As a result, the communication cost for the owner increases linearly to the number of verification requests at the same time. Differently, with the public verifiability of our *PCPOR* scheme, each challenger is able to conduct the verification independently with constant communication cost and the data owner can go off-line after outsourcing his data.

6.2 Computation

As mentioned in Section 5.1, our *PCPOR* scheme is composed of 4 algorithms: *KeyGen*, *Setup Prove* and *Verify*. Among these algorithms, *KeyGen* and *Setup* are performed by the data owner as data preparation process. To generate the public key PK as well as the private key SK for the system, the data owner performs $(s+3)$ EXP operations using the *KeyGen* algorithm. In the *Setup* procedure, to process an encoded data file with n blocks, each of which has s elements, $(2s+2)n$ EXP and $(s+1)$ MUL operations are needed to generate the authentication tags. Note that the data preparation process in our *PCPOR* scheme is one-time cost for the data owner, it can go off-line after finishing this process. For the computation cost in each verification request,

²In Goldreich[12]'s (γ, δ) -*hitter*, it is guaranteed that any subset $Sub \subset [1, n]$ with size $|Sub| \geq |F|(1 - \text{gamma})$, $Pr[I \cap Sub \neq \emptyset]$. For more details, please refer to Ref.[12].

| Scheme | Public verifiability | Comm. complexity (Challenge) (bits) | Comm. complexity (Response) (bits) | Comp. Cost (Server) | Comp. Cost (Challenger) | Storage Cost (Server) (bits) | Storage Cost (Challenger) (bits) | Data Preparation |
|--------|----------------------|-------------------------------------|------------------------------------|--|--|------------------------------|----------------------------------|--|
| [20] | Yes | $O(\lambda^2)$ | $O(s\lambda)$ | $l(\text{MUL}+\text{EXP})+s(\text{ZADD}+\text{ZMUL})$ | $(s+l)\text{MUL}+(s+l)\text{EXP}+2\text{Pairing}$ | $(1+\frac{1}{s}) F $ | $\lambda+ G $ | $(s+1)\text{MUL}+(s+2)n\text{EXP}$ |
| [10] | No | $O(\lambda)$ | $O(s\lambda)$ | $(sl+l)(\text{ZMUL}+\text{ZADD})$ | $(s+l)(\text{ZMUL}+\text{ZADD})+l\text{PRF}$ | $(1+\frac{1}{s}) F $ | 2λ | $l\text{PRF}+n(\text{ZMUL}+\text{ZADD})$ |
| [15] | No | $O(\lambda)$ | $O(\lambda)$ | $(s-1)(\text{EXP}+\text{MUL})+(s+l+sl)(\text{ZMUL}+\text{ZADD})$ | $2\text{EXP}+l\text{PRF}+l(\text{ZMUL}+\text{ZADD})$ | $(1+\frac{1}{s}) F $ | $3\lambda+80$ | $l\text{PRF}+n(\text{ZMUL}+\text{ZADD})$ |
| PCPOR | Yes | $O(\lambda)$ | $O(G)$ | $(l+s-1)\text{MUL}+(2l+s-1)\text{EXP}+sl\text{ZADD}+s(l+2)\text{ZMUL}$ | $l\text{MUL}+(2l+2)\text{EXP}+3\text{Pairing}$ | $(1+\frac{1}{s}) F $ | $\lambda+4 G $ | $(s+1)\text{MUL}+(2s+2)n\text{EXP}$ |

Table 1: Complexity Summary: in this table, MUL is one multiplication operation on Group G , EXP is one exponentiations operation on Group G , PRF denotes the pseudorandom function, ZADD and ZMUL represent the addition and multiplication operations on Z_p^* respectively; λ is the security parameter chosen for the system, $||G||$ is the size(in number of bits) of a group element on G , $|F|$ is size of data file, n is number of encoded blocks for the data file, s is the number of elements in each block and l is number of blocks selected for verification. Note: when the security parameter λ is fixed, the communication complexities in our PCPOR scheme and Ref.[15] become constant.

the server needs to perform $(l+s-1)\text{MUL}$, $(2l+s-1)\text{EXP}$, $sl\text{ZADD}$ and $s(l+2)\text{ZMUL}$ operations to generate the proof response, where l is number of blocks chosen in each verification(Our discussion in Section.7 shows that the value of l can be small). After receiving the proof information, the challenger first conducts $l\text{MUL}$ and $2l\text{EXP}$ operations to generate η . Then, 2 more EXP and 3 Pairing operations are needed for the final verification. Therefore, as shown in Table 1, the total computation cost for the challenger in one verification is $l\text{MUL}+(2l+2)\text{EXP}+3\text{Pairing}$.

We now compare our PCPOR scheme with the existing POR schemes[20, 10, 15] and show the result in Table 1. Beginning with the public POR scheme in Ref.[20], although our PCPOR will introduce $s-1$ more MUL, $l+s-1$ more EXP and $2s$ more ZMUL operations to the server side, it achieves the same computation cost level for the challenger side(client). As the cloud sever side is always much powerful than the challenger side(e.g. Amazon EC2 vs Mobile devices), the additional computations brought to server side in our PCPOR can be easily handled in practical scenarios and will have little influence on the scheme’s performance. When considering the two private POR schemes[10, 15] as shown in Table 1, which have almost the same computation cost for challenger side except for 2 more EXP operations in Ref.[15], our PCPOR scheme will cause relative higher computation cost by replacing the operations on Z_p^* to operations on G in each verification. However, it is notable that these private POR schemes[10, 15] have to centralize all the verification tasks to the data owner(i.e., the computation cost on the data owner is linear to the number of simultaneous verification requests), which requires the data owner to keep online and causes heavy computation burden on it. On the contrary, in our PCPOR, the computation tasks for each verification is distributed to the challengers without contacting the data owner. This public verifiability makes our scheme easy to scale and have practical computation cost on each challenger. For the computation cost for data preparation, our PCPOR scheme introduces sn more EXP operations compared to the public POR scheme in Ref.[20]. Since the data preparation process is one-time cost in our scheme and will not influence the realtime verification process, the additional operations in this process can be accepted in practical scenarios. Compared to the private POR schemes in Ref.[10, 15], as shown in Table 1, our PCPOR scheme and scheme in Ref.[20] need relative

higher computation cost for the preparation due to the requirement of group operations for the public verifiability. Nevertheless, as we mentioned above in this section, these additional one-time computation cost in data preparation process is reasonable and acceptable.

6.3 Storage

In this section, we first analyze the storage cost of our PCPOR scheme on both challenger side and server side, and then compare it with the existing POR schemes[20, 10, 15]. The results of our analysis are summarized in Table 1. For the challenger side, our PCPOR scheme only requires the challenger to store partial public key $PK : \{p, v, k, spk, g\}$ to generate the challenge message *Chall* and perform verification algorithm *Verify*. Thus, the size of storage cost for each challenger is $4||G|| + \lambda$ bits. Compared to the existing POR schemes[20, 10, 15], which require $||G|| + \lambda$ bits, 2λ bits and $3\lambda + 80$ bits respectively for the challenger side storage cost, our PCPOR scheme achieves the same storage cost level as demonstrated in Table 1. For storage overhead on the server side, it mainly comes from the authentication tags for the encoded data blocks. In our PCPOR scheme, each authentication tag is a group element with λ bits, thus the total size for tags is $n\lambda$ bits. As the total encoded data file size $|F| = ns\lambda$ bits, we represent the total storage overhead on server side as $(1+\frac{1}{s})|F|$ bits, which equals to Ref.[20, 10, 15]’s storage cost on server sides when the values of s and λ are same.

Note that the communication cost in our PCPOR scheme is independent to s as we mentioned in Section 6.1, which enables our scheme to reduce the storage cost on server side by adjusting the value of s . However, in Ref.[20, 10], as the communication cost for proof response is linear to the value of s , reducing the storage cost will lead to the sacrifice of communication performance.

7. DISCUSSION

In this section, we discuss about the error detection probability of our PCPOR scheme. As we mentioned in the *Setup* algorithm of PCPOR scheme, Reed-Solomon code with rate ℓ is adopted for the data file encoding. For a ℓ Reed-Solomon encoded data file, any ℓ fraction of encoded data blocks can recover the original file. If a data file encoded with ℓ Reed-Solomon code cannot be recovered from the erasure decoding, the probability of accessing a uncorrupted encoded

data block will be less than ℓ . In this case, when we randomly choose l independent encoded data blocks and all these blocks are uncorrupted, the probability should be less than ℓ^l .

In our *PCPOR* scheme, we can set $\ell = 0.98$ as previous *POR* scheme[15] does. In this case, by checking 200 encoded data blocks, the challenger can have at least 98.24% confidence that the stored data on the server is not corrupted if the *Verify* algorithm outputs result as accept. 99.99% confidence can be guaranteed if 1000 encoded data blocks are checked by the challenger.

8. CONCLUSIONS

Proofs of Retrievability(*POR*) technique enables individuals/organizations to ensure the integrity of their outsourced data on a untrusted server(e.g.public cloud storage platform). Nevertheless, the existing *POR* schemes either have limitation on communication cost linear to the number of elements in a data block, or only consider private verification. These limitations cause a severe scalability issue in data file size or user number for practical use. In this work, we proposed the first public *POR* scheme with constant communication cost: *PCPOR*. By uniquely incorporating the polynomial commitment technique and other cryptographic techniques, our *PCPOR* scheme achieves constant communication size, efficient computation performance as well as low storage overhead. What is more, based on the simultaneous realized public verifiability in our *PCPOR* scheme, we release the data owner from onerous verification tasks and make each client check integrity independently, which need to be centralized to the data owner in previous private *POR* scheme with constant communication size. Our security proof based on CDH Problem, SDH Assumption and BSDH Assumption shows that our *PCPOR* scheme is secure. Our thorough analysis demonstrates the efficiency and scalability of our scheme.

9. REFERENCES

- [1] Amazon forum. major outage for amazon s3 and ec2, <https://forums.aws.amazon.com/thread.jspa?threadID=19714&start=15&tstart=0>.
- [2] Amazon web service. summary of the amazon ec2 and amazon rds service disruption in the us east region, <http://aws.amazon.com/message/65648/>.
- [3] Business insider. amazon's cloud crash disaster permanently destroyed many customers' data, <http://www.businessinsider.com/amazon-lost-data-2011-4>.
- [4] Dropbox. dropbox forums on data loss topic, <http://forums.dropbox.com/tags.php?tag=data-loss>.
- [5] G. Ateniese, R. Burns, R. Curtmola, J. Herring, L. Kissner, Z. Peterson, and D. Song. Provable data possession at untrusted stores. In *Proceedings of the 14th ACM conference on Computer and communications security*, CCS '07, pages 598–609, New York, NY, USA, 2007. ACM.
- [6] D. Boneh and X. Boyen. Short signatures without random oracles. pages 56–73. Springer-Verlag, 2004.
- [7] D. Boneh, B. Lynn, and H. Shacham. Short signatures from the weil pairing. In *Proceedings of the 7th International Conference on the Theory and Application of Cryptology and Information Security: Advances in Cryptology*, ASIACRYPT '01, pages 514–532, London, UK, UK, 2001. Springer-Verlag.
- [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*, CCSW '09, pages 43–54, New York, NY, USA, 2009. ACM.
- [9] W. Diffie and M. Hellman. New directions in cryptography. *IEEE Trans. Inf. Theor.*, 22(6):644–654, Sept. 1976.
- [10] Y. Dodis, S. Vadhan, and D. Wichs. Proofs of retrievability via hardness amplification. In *Proceedings of the 6th Theory of Cryptography Conference on Theory of Cryptography*, TCC '09, pages 109–127, Berlin, Heidelberg, 2009.
- [11] C. Erway, A. Küpçü, C. Papamanthou, and R. Tamassia. Dynamic provable data possession. In *Proceedings of the 16th ACM conference on Computer and communications security*, CCS '09, pages 213–222, New York, NY, USA, 2009. ACM.
- [12] O. Goldreich, S. Goldwasser, and S. Micali. How to construct random functions. *J. ACM*, 33(4):792–807, Aug. 1986.
- [13] V. Goyal. Reducing trust in the pkg in identity based cryptosystems. In *Proceedings of the 27th annual international cryptology conference on Advances in cryptology*, CRYPTO'07, pages 430–447, Berlin, Heidelberg, 2007. Springer-Verlag.
- [14] P. Hawkes, M. Paddon, and G. G. Rose. On corrective patterns for the sha-2 family, 2004.
- [15] X. Jia and C. Ee-Chien. Towards efficient provable data possession. In *Proceedings of the 7th ACM Symposium on Information, Computer and Communications Security*, ASIACCS '12, Seoul, Korea, 2012.
- [16] 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*, CCS '07, pages 584–597, New York, NY, USA, 2007. ACM.
- [17] A. Kate, G. M. Zaverucha, and I. Goldberg. Constant-size commitments to polynomials and their applications. In *ASIACRYPT*, pages 177–194, 2010.
- [18] G. Oded. A sample of samplers - a computational perspective on sampling (survey). *Electronic Colloquium on Computational Complexity (ECCC)*, 4(20), 1997.
- [19] I. S. Reed and G. Solomon. Polynomial Codes Over Certain Finite Fields. *Journal of the Society for Industrial and Applied Mathematics*, 8(2):300–304, 1960.
- [20] H. Shacham and B. Waters. Compact proofs of retrievability. In *Proceedings of the 14th International Conference on the Theory and Application of Cryptology and Information Security: Advances in Cryptology*, ASIACRYPT '08, pages 90–107, Berlin, Heidelberg, May 2008. Springer-Verlag.
- [21] V. Shoup. *A computational introduction to number theory and algebra*. Cambridge University Press, New York, NY, USA, 2005.
- [22] G. Timothy and M. M. Peter. The nist definition of cloud computing. NIST SP - 800-145, September 2011.