

Secure Deduplication of Encrypted Data: Refined Model and New Constructions

Jian Liu^{*1}[0000-0001-6796-6828], Li Duan^{*2}[0000-0002-8383-0776],
Yong Li³[0000-0002-6920-0663], N. Asokan¹[0000-0002-5093-9871]

¹ Aalto University, Finland

`jian.liu@aalto.fi`, `asokan@acm.org`

² Paderborn University, Germany

`liduan@mail.upb.de`

³ Ruhr-University Bochum, Germany

`yong.li@rub.de`

Abstract

Cloud providers tend to save storage via cross-user deduplication, while users who care about privacy tend to encrypt their files on client-side. Secure deduplication of encrypted data (SDoE) which aims to reconcile this apparent contradiction is an active research topic. In this paper, we propose a formal security model for SDoE. We also propose two single-server SDoE protocols and prove their security in our model. We evaluate their deduplication effectiveness via simulations with realistic datasets.

1 Introduction

Cloud storage services are very popular. Providers of cloud storage services routinely use *cross-user deduplication* to save costs: if two or more users upload the same file, the storage provider stores only a single copy of the file. Users concerned about privacy of their data may prefer encrypting their files on client-side before uploading them to cloud storage. This thwarts deduplication since identical files are uploaded as completely different ciphertexts. Reconciling deduplication and encryption has been a very active research topic [3,10,4,16,20,14]. One proposed solution is *convergent encryption* (CE) [10] [4], which derives the file encryption key solely and deterministically from the file contents. As a result, identical files will always produce identical ciphertexts given identical public parameters. Unfortunately, a server compromised by the adversary can perform an *offline brute-force guessing attack* over the ciphertexts, due to the deterministic property of CE.

More recent solutions allow clients to encrypt their files using stronger encryption schemes while allowing the server to perform deduplication. They usually assume the presence of independent (trusted) third parties [3] [16] [20]. However, in a cloud storage setting, like in many other commercial client-server

* Equal contribution.

contexts, assuming the presence of an independent third party is unjustified in practice [14] since it is unclear who can bear the costs of such third parties. Moreover, such schemes *cannot* prevent *online brute-force guessing attacks* from a compromised active server.

Liu et al. proposed a single-server scheme for secure deduplication without the need for any third party [14]. Their scheme uses a per-file rate limiting strategy to prevent online brute-force guessing attacks by a compromised active server. However, their security model and proof only cover one round of the protocol (Section 9 in [15]). Consequently, their scheme is vulnerable to additional attacks when considering the long-term operation of the system which involves multiple rounds of the protocol.

In this paper, we make the following contributions:

- We propose a **formal security model** for the single-server “secure deduplication of encrypted data” (SDoE) (Section 2). We claim that a deduplication scheme proved secure in this model can guarantee that, for a certain file, (1) a compromised client cannot learn whether or not this file has already been uploaded by someone else (Section 2.1), and (2) the only way for a compromised server to uniquely determine this file is by doing an online brute-force attack (Section 2.2).
- We propose two new **single-server SDoE** schemes and prove their security in our model. (Section 4)
- We show that their deduplication effectiveness is reasonable via **simulations with realistic datasets**. (Section 5)

2 Syntax and Security Model

2.1 Syntax

We consider the generic setting for a cloud storage system where a set of clients (\mathcal{C} s) store their files on a single storage server (\mathcal{S}), and \mathcal{C} s and \mathcal{S} are always communicating through secure channels. The deduplication happens at server-side, i.e., the client always uploads encrypted files and the server knows whether to discard the uploaded file or not after the protocol execution. All these participants are generalized as *parties*. Each party has a party identifier pid ; a flag τ indicating whether it is corrupted or not. Each \mathcal{C} may have one or more sessions connecting to \mathcal{S} , where each session has a session identifier sid . The internal state $\Phi_{\mathcal{C},pid}$ of a \mathcal{C} is a list of tuples $\{(fid_i, k_i)\}$, which stores the identifier and the encryption key of each file owned by it. The internal state $\Phi_{\mathcal{S}}$ of \mathcal{S} contains a list $\mathbf{DB} = \{(b_i, fid_i, \varrho_i, \mathbf{LO}_i)\}$ and a list of current user identifiers \mathbf{PID} , where b_i is a bit indicating whether the file has been uploaded or not, fid_i is an identifier of an encrypted file ϱ_i and \mathbf{LO}_i is the list of owners of ϱ_i . Note that \mathbf{DB} contains all possible files.

Definition 1. *A single-server SDoE scheme $\Pi_{\text{dedup}} = (\text{Init}, \text{Reg}, \text{Upload}, \text{Download})$ is composed of an initialization algorithm Init and three sub-protocols Reg , Upload and Download . Each component is defined as follows.*

- $\text{Init}(1^\lambda, \text{aux}) \rightarrow \mathbf{PP}$. The Init algorithm takes the security parameter 1^λ and the auxiliary information aux as input and outputs the public parameter \mathbf{PP} , which includes the specification of the encryption scheme Π_{Enc} chosen for the files.
- $\text{Reg}(pid) \rightarrow \Phi_{pid}$. The protocol Reg register a new client \mathcal{C} with identifier pid and returns a new client internal state Φ_{pid} .
- $\text{Upload}(\mathbf{PP}, \Phi_c, \mathbf{DB}_s, \{\Phi_j\}, F) \rightarrow \Phi'_c, \mathbf{DB}'_s$. The protocol Upload involves the uploader \mathcal{C} , the server \mathcal{S} and a group of possible file owners $\{\mathcal{C}_j\}$. Taking as input Φ_c of \mathcal{C} , \mathbf{DB}_s of \mathcal{S} and $\{\Phi_j\}$ of $\{\mathcal{C}_j\}$, this protocol produces updated states Φ'_c and \mathbf{DB}'_s of \mathcal{C} and \mathcal{S} respectively. A file identifier fid is contained both in Φ'_c and \mathbf{DB}'_s and a pair (fid, k_f) is in Φ'_c with encryption/decryption key k_f conformant with Π_{Enc} .
- $\text{Download}(\Phi_c, \mathbf{DB}_s, fid_i) \rightarrow F$. The protocol Download involves the downloader \mathcal{C} and the server \mathcal{S} . Besides the internal state of \mathcal{C} and \mathcal{S} , this protocol takes as an extra input a file identifier fid_i and outputs a file F .

Definition 2. (Correctness) A Π_{dedup} is correct if

$$\begin{aligned} \forall (fid, k_f) \in \Phi'_c \leftarrow \text{Upload}(\mathbf{PP}, \Phi_c, \mathbf{DB}_s, \{\Phi_j\}, F) : \\ \Pr[\text{Download}(\Phi'_c, \mathbf{DB}_s, fid) = F] = 1 \end{aligned}$$

2.2 Game Setup

One way to model security in cryptography is by issuing security games played between an adversary (attacker) and a challenger. The challenger possesses some secret targeted by the adversary. As in the real world, the adversary can interact with the challenger by using different queries to the challenger. At the end of each game, the adversary outputs what it has learned about the secret and it wins if the output is correct. The restriction on queries are used to rule out the trivial cases of breaking the security of the scheme.

To initialize \mathbf{DB} for the security game, the challenger first generates a list of file owners with the corresponding identifier list \mathbf{PID} . After calling Init with given security parameter and the auxiliary information to generates the public parameter \mathbf{PP} , the challenger calls $\text{Reg}(id)$ for each id in \mathbf{PID} . Let \mathbf{PID}_h be the identifier set of honest file owners. Note that before interacting with adversaries, $\mathbf{PID} = \mathbf{PID}_h$, but \mathcal{A} can add new malicious identities to \mathbf{PID} by using $\text{RegisterCorrupt}(pid)$ queries described below. Then for each fid_i , the challenger chooses $b_i \stackrel{\$}{\leftarrow} \{0, 1\}$. If $b_i = 0$, the tuple would be $(0, fid_i, -, \emptyset, -)$. Otherwise, it chooses a uniformly random $\mathcal{C}_j \in \mathbf{PID}$ and calls $\text{Upload}(\mathbf{PP}, \Phi_j, \mathbf{DB}, \{\Phi_l\}, F_i)$ to upload F_i . Note that the updated state \mathbf{DB} now contains the ciphertext ϱ_i of F_i generated and \mathcal{C}_j is added to \mathbf{LO}_i . We denote as \mathbf{DB}_0 the content of \mathbf{DB} after initialization.

2.3 Security against a compromised client

As noticed by Harnik [12], a client can use deduplication as a side-channel to obtain information about the contents of files of other clients. Here, we want to

model attacks from a compromised client. The intuition is that by interacting with the server, the client must not be able to learn whether a file already exists in the cloud storage. We allow the adversary \mathcal{A} that has compromised one or more clients to make the following types of oracle queries in the security experiments:

- RegisterCorrupt**(pid) The adversary \mathcal{A} can register a new corrupted client with identifier pid . If $pid \in \mathbf{PID}$, \mathcal{A} gets the state $\Phi_{C,pid}$ of pid , including all the file identifiers and the corresponding file key $\{(fid_i, k_i)\}$ owned by pid . Otherwise \mathcal{A} only gets an empty state. The challenger updates $\mathbf{PID} := \mathbf{PID} \cup \{pid\}$ and marks pid as corrupted. In both cases, \mathcal{A} can perfectly impersonate pid from this moment on.
- Send**(pid, sid, M) The corresponding oracle computes on the input message M following the SDoE protocol and returns the output message in the view of all corrupted parties to \mathcal{A} . This oracle models that an adversary can tamper with each single message in the SDoE protocol. For example, initiating Download with maliciously formed file identifier.
- Test**() \mathcal{A} signals the end of the security game to the challenger, ceases all the interaction with oracles, and outputs a pair (F^*, b^*) . The adversary can only query this oracle once. Note that this query is not an abstraction of any attack but serves as a measure of the adversarial success.

Let λ be the security parameter. Given the queries described above, we define the security experiment $\mathbf{Exp}_{C,\Pi}^{\text{SDoE}}(\lambda)$ for a SDoE protocol Π against compromised clients as follows: $\mathbf{Exp}_{C,\Pi}^{\text{SDoE}}(\lambda) = 1$ if \mathcal{A} replies to **Test**() with (F^*, b^*) and either of the following cases happens:

- If $b^* = 0$ and $\nexists(fid_j, k_j) \in \bigcup_{pid \in PID_h} \Phi_{C,pid}$, $(b_i, fid_i, \varrho_i, \mathbf{LO}_i) \in \mathbf{DB}_0$ s.t., $E(k_j, F^*) = \varrho_i$. (i.e., F^* has not been uploaded before).
- If $b^* = 1$ and $\exists(fid_j, k_j) \in \bigcup_{pid \in PID_h} \Phi_{C,pid}$, $(b_i, fid_i, \varrho_i, \mathbf{LO}_i) \in \mathbf{DB}_0$ s.t., $E(k_j, F^*) = \varrho_i$. (i.e., F^* has been uploaded before).

But *none* of the following events happens before \mathcal{A} outputs (F^*, b^*) :

- \mathcal{A} has issued **RegisterCorrupt**(pid) with $pid \in \mathbf{PID}_h$ (i.e., \mathcal{A} cannot directly read internal information of honest file owners.)
- \mathcal{A} has issued **Send**(pid, sid, M) with $pid \in \mathbf{PID}_h$. (i.e., \mathcal{A} cannot force an honest owner to send/receive any messages. However, \mathcal{A} can use the **Send** queries to fully control the behavior of compromised clients.)

Definition 3. We define the advantage of an adversary \mathcal{A} in the experiment $\mathbf{Exp}_{C,\Pi}^{\text{SDoE}}(\lambda)$ as

$$\mathbf{Adv}_{C,\Pi}^{\text{SDoE}}(\lambda) = \Pr[\mathbf{Exp}_{C,\Pi}^{\text{SDoE}}(\lambda) = 1] - \frac{1}{2}$$

2.4 Security against a compromised server

The intuition behind the security definition is that a SDoE scheme is secure against a compromised server cannot be uniquely determined by the compromised server. The queries (adversary’s ability) captures the essence of concrete attacks, such as registering malicious clients, uploading and tampering some messages. Those attacks may come from a malicious server colluding with some external clients. For simplicity, we now assume that all files have been uploaded into \mathbf{DB}_0 . We allow the adversary \mathcal{A} that has compromised the server to make the following types of queries in the security experiments:

RegisterCorrupt(pid) The same as that of compromised clients.

Send(pid, sid, M) The same as that of compromised clients.

AccessDB() The adversary \mathcal{A} gets all the ϱ_i and the owner list of each ϱ_i in Φ_S with $b_i = 1$. If this is the t_a -th query made by \mathcal{A} , then for all the t -th queries with $t > t_a$, \mathcal{A} also gets the updated Φ_s items with $b_i = 1$ in addition to the response for other queries.

Execute(pid, P, F) As the initiator, \mathcal{A} invokes a complete (sub-)protocol $P \in \{\text{Upload, Download}\}$ with party pid on the input file F and obtains all the messages exchanged, following the description of P .

Test() \mathcal{A} outputs two files F_0, F_1 with equal length. Upon receiving F_0, F_1 , the challenger chooses $b \xleftarrow{\$} \{0, 1\}$ and replies with a ciphertext $\varrho_b = \text{Enc}(k_{f_b}, F_b)$. \mathcal{A} performs the above queries and then outputs a bit b' .

We define the security experiment $\mathbf{Exp}_{S,II}^{\text{SDoE}}(\lambda)$ for a SDoE protocol II against partially compromised server as follows: $\mathbf{Exp}_{S,II}^{\text{SDoE}}(\lambda) = 1$ if \mathcal{A} replies to **Test**() with $b' = b$, but *none* of the following events happens before \mathcal{A} outputs the bit b' :

- \mathcal{A} has issued **RegisterCorrupt**(pid) with $pid \in \mathbf{PID}_h$.
- \mathcal{A} has issued **Execute**(pid, P, F), $F \in \{F_0, F_1\}$. (\mathcal{A} have not included F_0 or F_1 in its online brute-force attacks.)
- \mathcal{A} has issued **Send**(pid, sid, M) with $pid \in \mathbf{PID}_h$.

Definition 4. We define the advantage of an adversary \mathcal{A} in the experiment $\mathbf{Exp}_{S,II}^{\text{SDoE}}(\lambda)$ as

$$\mathbf{Adv}_{S,II}^{\text{SDoE}}(\lambda) = \Pr[\mathbf{Exp}_{S,II}^{\text{SDoE}}(\lambda) = 1] - \frac{1}{2}$$

Remark. We exclude trivial wins by preventing the adversary from corrupting honest clients, but this is not overly restrictive because the adversary can still steer honest clients by issuing **Execute** queries.

3 PAKE based Deduplication

Bellovin and Merritt [7] proposed a *password authenticated key exchange* (PAKE) protocol to resist offline brute-force attacks even through users choose low-entropy passwords. PAKE enables two parties to set up a session key iff they

hold the same secret (“password”). Otherwise, neither party can learn anything about the key output by the other party.

Bellare et al. provided a game-based definition for the security of PAKE [5]. A random bit b is chosen at the beginning of the game. They assume that there is an adversary \mathcal{A} that has complete control over the environment (mainly, the network), and is allowed to query the following oracles:

$\text{Send}(U_i, M)$: causes message M to be sent to instance U_i , which computes following the protocol and gives the result to \mathcal{A} . If this query causes U_i to accept or terminate, this will also be shown to \mathcal{A} .

$\text{Execute}(A_i, B_j)$: causes the protocol to be executed to completion between A_i and B_j , and outputs the transcript of the execution.

$\text{Reveal}(U_i)$: output k_{U_i} , which is the session key held by U_i .

$\text{Test}()$: if $b = 1$, output the session key k_{U_i} ; otherwise, output a string drawn uniformly from the space of session keys. Note that Test is queried only once.

$\text{Corrupt}(U_i)$: output U_i 's password.

Let $\text{Succ}_{\mathcal{A}}^{\text{PAKE}}(\lambda)$ be the event that \mathcal{A} outputs a bit $b' = b$ but none of the following events happens:

1. a $\text{Reveal}(U_i)$ query occurs;
2. a $\text{Reveal}(U_j)$ query occurs where U_j is the partner of U_i ;
3. a $\text{Corrupt}(U_i)$ query occurs before U_i defined its key k_{U_i} and a $\text{Send}(U_i, M)$ query occurred.

The advantage of \mathcal{A} attacking a PAKE protocol is defined to be

$$\text{Adv}_{\mathcal{A}}^{\text{PAKE}}(\lambda) \stackrel{\text{def}}{=} 2\text{Pr}[\text{Succ}_{\mathcal{A}}^{\text{PAKE}}(\lambda)] - 1.$$

The PAKE protocol is considered secure if passwords are uniformly and independently drawn from a dictionary of size n :

$$\text{Adv}_{\mathcal{A}}^{\text{PAKE}}(\lambda) \leq \frac{n_{se}}{n} + \text{negl}(\lambda),$$

where n_{se} is the number of Send queries (to distinct instances U_i). The intuition behind this definition is that only online brute-force attacks are allowed in a secure PAKE protocol.

PAKE-based SDoE. Liu et al. presented a PAKE-based SDoE scheme that does not depend on any additional independent servers [14]. Their scheme allows an uploader to securely obtain the decryption key of another user who has previously uploaded the same file. Specifically, the uploader \mathcal{C} first sends a short hash of its file (10-20 bits long) to \mathcal{S} . \mathcal{S} finds other clients who may hold the same files based on the short hash, and lets them run a single round PAKE protocol (routed through \mathcal{S}) with the long hashes of their files as inputs. At the end of the protocol, the uploader gets the key of another \mathcal{C} if and only if they indeed hold the same file. Otherwise, it gets a random key. The PAKE-based SDoE scheme Π_{PAKE} is shown in Figure 1. Notice that Π_{PAKE} uses additively homomorphic encryption for the key transformation. Namely, after PAKE, the

uploader sends $Enc(pk, k_{iR} + r)$ to the server. The goal of using additively homomorphic encryption $Enc()$ is to: (1) guarantee the privacy of k_{iR} ; and (2) allow the server to compute $Enc(pk, (k'_{iR} + k_{F_j}) - (k_{iR} + r))$, such that the uploader can get k_{F_j} iff $k'_{iR} = k_{iR}$. We noticed that the use of additively homomorphic encryption [15] is an overkill, since both of these two goals can be achieved more efficiently by using a one-time-pad $k_{iR} \oplus r$: (1) the server cannot learn anything about k_{iR} since r is randomly chosen by the uploader; (2) the server can compute $(k_{iR} \oplus r) \oplus (k'_{iR} \oplus k_{F_j})$, such that the uploader can also get k_{F_j} iff $k'_{iR} = k_{iR}$.

If there is no match on the short hash, \mathcal{S} lets the uploader run PAKE with dummy checkers to hide the fact that the file has not been uploaded before. In addition, \mathcal{C} s protect themselves against online brute-force attacks by limiting the number of PAKE instances they will participate in for each file.

Security against compromised clients. As pointed out by Liu et al. themselves in [15], additional attacks are possible when considering the long-term operation of the system. For example, a malicious client can upload a file and then pretend to be offline. Later it uploads the same file using another identity. If it gets the the same key as the one it got before, it knows that the file has been uploaded by someone else. Another attack is also targeting the PAKE phase. The adversary uploads a file F with the identity of \mathcal{C}_1 in the first protocol run. It then uses a different identity \mathcal{C}_2 to upload F again. By observing whether \mathcal{C}_1 is involved in the PAKE phase with \mathcal{C}_2 for F , the adversary knows if there are other owners of F .

In the next section, we will introduce two protocols that are immune to those attacks and prove their security under our new model.

4 New SDoE Schemes

Recall that in Π_{PAKE} , there are two possible cases when an uploader uploads a file: it either gets the key of a previous uploader of the same file or gets a new random key. As described in the previous section, a malicious \mathcal{C} can distinguish between these two cases.

In this section, we address the issue in Π_{PAKE} by having \mathcal{C} s always get random keys when they upload their files. We propose two schemes. The first scheme ($\Pi_{\text{PAKE, re-enc}}$) borrows the idea from proxy re-encryption [1]. \mathcal{S} only keeps a single copy of duplicated files. When \mathcal{C} wants to download its file, \mathcal{S} re-encrypts the file so that \mathcal{C} will download the same ciphertext as the one it uploaded. However, this scheme requires public-key operations on the entire file, which is not efficient for large files. So we propose a second scheme ($\Pi_{\text{PAKE, popular}}$) that only deduplicates popular files and only protects the privacy of unpopular files. For unpopular files, \mathcal{C} s get random keys and download the same ciphertexts as they uploaded. If those files become popular later, \mathcal{S} deletes all duplicated copies and provides a way to help \mathcal{C} s to transform their keys to the right key.

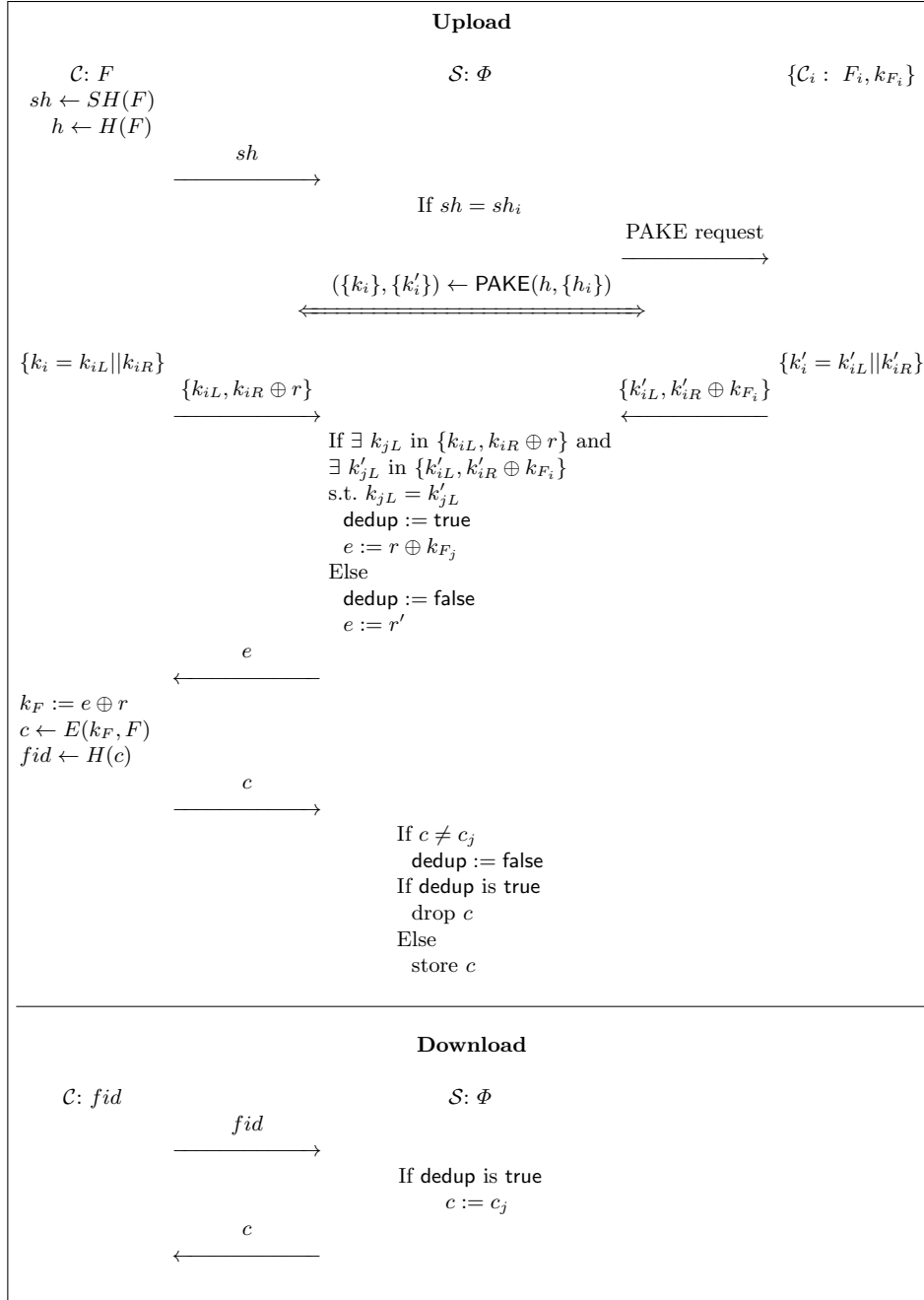


Fig. 1: PAKE-based deduplication scheme [14].

4.1 PAKE-based Deduplication with Re-encryption

The first scheme $\Pi_{\text{PAKE, re-enc}}$ is shown in Figure 2. It is similar to Π_{PAKE} . In the following description, the details of client authentication and file ownership authentication are omitted. We assume that the owners of each file are stored in an ordered list with respect to the upload time points. In the case that there are more than one owner of a candidate file, the newest checker is chosen by \mathcal{S} for the PAKE phase. After PAKE, instead of masking k_{F_i} with k'_{iR} , \mathcal{C}_i generates a random number r_i , and masks both k_{F_i} and k'_{iR} with r_i . \mathcal{C} only sends k_{iL} to \mathcal{S} . If there is an index j s.t. $k_{jL} = k'_{jL}$, \mathcal{S} knows that \mathcal{C} is uploading the same file with \mathcal{C}_j . Then, it keeps $(r_j - k_{F_j})$ and sends $(k'_{jR} + r_j)$ to \mathcal{C} . Otherwise, it sends a random number r' . \mathcal{C} calculates its file key as $k_F := e - k_{jR}$ and then encrypts its file as $F \cdot g^{k_F}$. Notice that if F is detected to be duplicated, k_F is just the randomness r_j generated by \mathcal{C}_j . \mathcal{S} can just drop this ciphertext if deduplication happens and stores the $fid = H(C)$ as an alias of the file. Later, when \mathcal{C} wants to download F , \mathcal{C} re-encrypts c_j to k_F : $c := c_j \cdot g^{r_j - k_{F_j}} = F \cdot g^{k_{F_j}} g^{r_j - k_{F_j}} = F \cdot g^{r_j} = F \cdot g^{k_F}$. Notice that c_j may be deduplicated already. In this case, \mathcal{S} need to calculate $(r_0 - k_{F_0}) + (r_1 - k_{F_1}) + \dots + (r_j - k_{F_j}) = (r_j - k_{F_0})$, and then transfer c_0 to \mathcal{C} 's ciphertext. We follow the same dummy checker and rate limiting strategy as Π_{PAKE} . The correctness of $\Pi_{\text{PAKE, re-enc}}$ is trivial.

Security against compromised clients. Security of SDoE schemes cannot be directly reduced to the semantic security of PAKE schemes in [5]. This technical impossibility in the proof lies in the fact that the password (the hash of the file) is always known to the adversary in SDoE prior to any other interactions in the PAKE protocol. To overcome this difficulty, we expand the original definition of the model in [5] in the following way, which we call the constrained PAKE security game. Let $sk = sk_L || sk_R$ be the session key computed in the $\text{Test}()$ session, where $|sk_L| = |sk_R| = \frac{1}{2}|sk|$.

- The setup of this game is the same as in the original PAKE game except that each party now holds an additional secret $s_u \in \mathcal{K}$. A public function $f : \mathcal{K} \times \{0, 1\}^* \rightarrow \{0, 1\}^{\frac{1}{2}|sk|}$ can be queried by the adversary as $f(p_i, \cdot)$.
- The $\text{Test}()$ -query now returns $tk = tk_L || tk_R$, where $|tk_L| = |tk_R| = \frac{1}{2}|sk|$. The first half of tk is always the same as the the first half of the real session key, i.e., $tk_L = sk_L$. If $b = 1$, $tk_R = sk_R \oplus f(s_i, T_i)$, where $T_{i,s}$ is the transcript of this session. Otherwise $tk_R = sk_R \oplus r$, where $r \xleftarrow{\$} \{0, 1\}^{\frac{1}{2}|sk|}$. The adversary wins if she outputs $b' = b$.
- $\text{Corrupt}(u)$ only returns the password PW_u but not the additional secret s_u .
- A session involving π_i^s and π_j^t is fresh if both the following condition holds
 - no $\text{Reveal}(s, i)$ or $\text{Reveal}(t, j)$ is made before $\text{Test}()$.
 - no $f(p_i, T_{i,s})$ is made before $\text{Test}()$.

The winning condition and the advantage of an adversary in a constrained-PAKE game is defined in the same way as in the PAKE game.⁴

⁴ We also assume that the implicit authentication property is preserved in the PAKE protocol as in the ideal functionality $\mathcal{F}_{\text{same-input-pake}}$ in [15]. The extension of the constrained-PAKE with implicit authentication is straight forward.

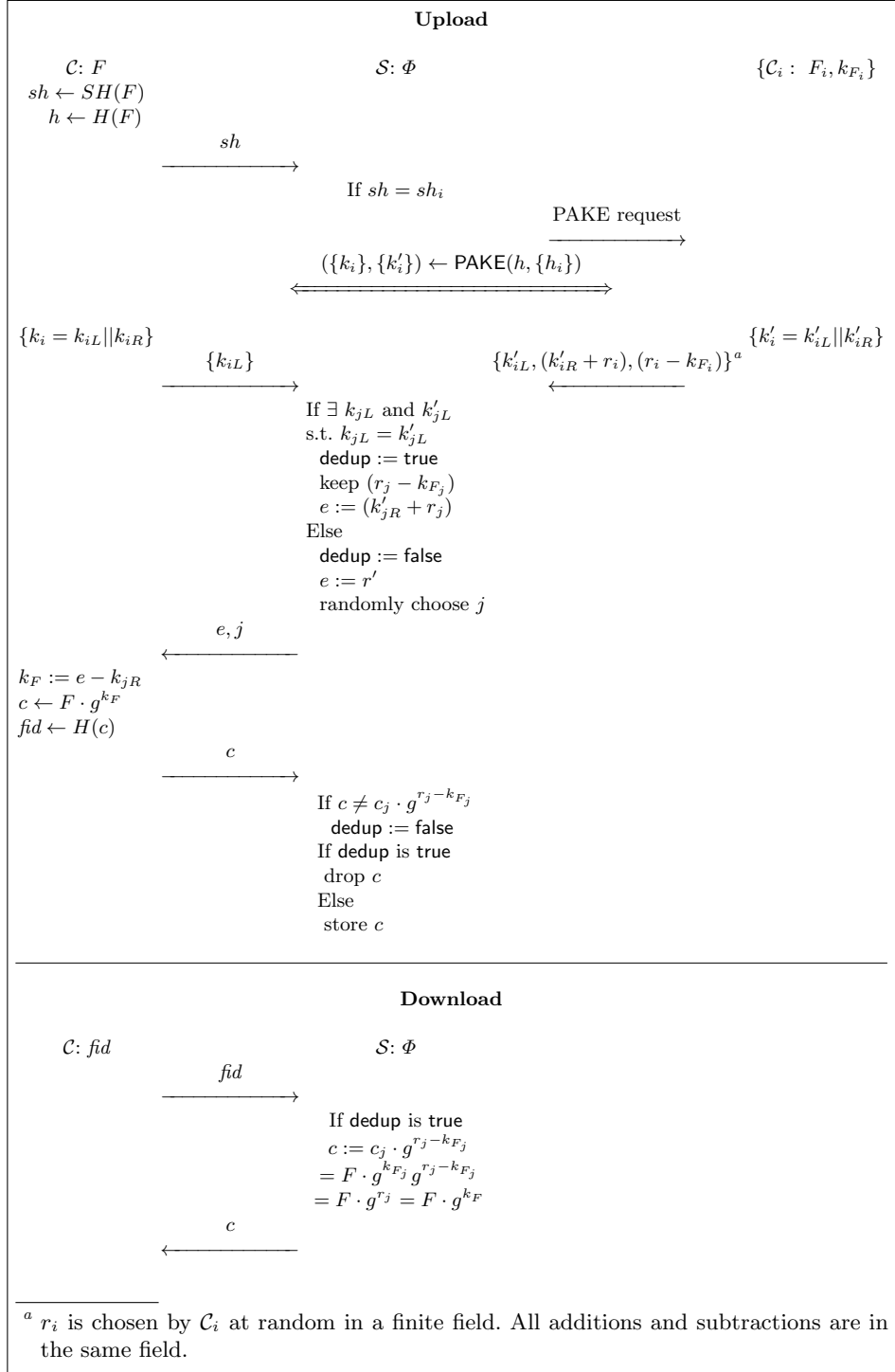


Fig. 2: PAKE-based deduplication via ciphertext transformation.

Theorem 1 *If there exists a ppt adversary \mathcal{C} in $\mathbf{Exp}_{\mathcal{C}, \Pi_{\text{PAKE, re-enc}}}^{\text{SDoE}}(\lambda)$ with advantage $\epsilon_{\mathcal{C}}$, then there also exists a ppt adversary \mathcal{A} with advantage $\epsilon_{\mathcal{A}}$ in the underlying constrained-PAKE game against Π in the random oracle model such that*

$$\epsilon_{\mathcal{C}} \leq \frac{q_H}{2^{l_h}} + q_H \cdot \epsilon_{\mathcal{A}}$$

where Π is the PAKE oracle, l_h the length of the long hash and q_H is the number of distinct files \mathcal{C} has queried for short hash, hash or uploaded.

Proof. We use the sequence of games technique introduced in [19]. We assume that the hash function is simulated by the challenger and all files are of equal length.

Game 0. This is the original game $\mathbf{Exp}_{\mathcal{C}, \Pi_{\text{PAKE, re-enc}}}^{\text{SDoE}}(\lambda)$.

$$\epsilon_{\mathcal{C}} = \mathbf{Adv}_{\mathcal{C}}^{\text{Game } 0}(\lambda) \tag{1}$$

Game 1. Let $\mathbf{F}' = \{F'_1, F'_2, \dots, F'_{q_H}\}$ be the set of distinct files that \mathcal{C} has issued H -queries or used for $\text{Send}()$ -queries before \mathcal{C} queries $\text{Test}()$. Let (F^*, b^*) be the output of \mathcal{C} . If $\exists F'_i \in \mathbf{F}' : H(F'_i) = H(F^*) \wedge F'_i \neq F^*$, abort the game. Then,

$$\mathbf{Adv}_{\mathcal{C}}^{\text{Game } 0}(\lambda) \leq \frac{q_H}{2^{l_h}} + \mathbf{Adv}_{\mathcal{C}}^{\text{Game } 1}(\lambda) \tag{2}$$

This rule makes sure that no hash collision happens.

Game 2. The challenger makes a guess of an index $i \in \{1, \dots, q_H\}$ and if $F^* \neq F'_i$, the challenger aborts the game⁵. Thus, we have

$$\mathbf{Adv}_{\mathcal{C}}^{\text{Game } 1}(\lambda) \leq q_H \cdot \mathbf{Adv}_{\mathcal{C}}^{\text{Game } 2}(\lambda) \tag{3}$$

Game 3. In this game, the random oracle in $\Pi_{\text{PAKE, re-enc}}$ is replaced (implicitly) by the random process for password generation. More specifically, we define $H : \{0, 1\}^* \rightarrow PW$, where PW is the password space and all the public parameters of Π (for example, group order and generator as in EKE2 [6]) are included in the public parameters of $\Pi_{\text{PAKE, re-enc}}$. This replacement has no impact on \mathcal{C} 's view since all passwords and parameters in Π are also sampled uniformly at random as required. Thus

$$\mathbf{Adv}_{\mathcal{C}}^{\text{Game } 2}(\lambda) = \mathbf{Adv}_{\mathcal{C}}^{\text{Game } 3}(\lambda) \tag{4}$$

We now construct an adversary \mathcal{A} using \mathcal{C} against the underlying PAKE scheme Π . Let d be the number of distinct passwords used by parties initialized by \mathcal{A} 's PAKE challenger, where $d \geq |\mathbf{F}|^6$. \mathcal{A} maps all files in \mathbf{F} to PAKE parties $P_1, \dots, P_{|\mathbf{F}|}$ with different passwords.⁷ \mathcal{A} then sets up the list \mathbf{PID}_h and \mathbf{PID} as described in the model and binds the identifiers in \mathbf{PID} to each file randomly. Each file identifier fid_i and encryption key k_{F_i} are chosen according to the protocol definition and then used to build each $\Phi_{\mathcal{C}, pid}$. Finally, \mathcal{A} stores all the $(b_i, fid_i, F_i, k_{F_i}, \mathbf{LO}_i)$ tuples as \mathbf{DB}_0 . To answer the Hash queries $H(F'_j)$,

⁵ The rationale of this abort rule can be found in Game 3.

⁶ We assume a polynomial sized file space.

⁷ This mapping ensures that all hash queries are answerable.

1. \mathcal{A} searches for $(F'_j, \{P'_j\})$.
2. If found, \mathcal{A} issues $\text{Corrupt}(P'_j)$ and let the output of $\text{Corrupt}(P'_j)$ be $PW_{P'_j}$.
3. \mathcal{A} returns $PW_{P'_j}$ to \mathcal{C} .

To answer $\text{RegisterCorrupt}(pid)$ for $pid \notin \mathbf{PID}_h$, \mathcal{A} simply enrolls this pid in \mathbf{PID} . If $pid \in \mathbf{PID}_h$, \mathcal{C} fails automatically. To answer $\text{send}(pid, M)$, for $pid \notin \mathbf{PID}_h$, \mathcal{A} answers the send-queries exactly as in the SDoE protocol. The hash-value or the PAKE messages to be returned are obtained in the same way as when answering hash queries.

During every PAKE for uploading F_i , \mathcal{A} uses Send to involve an oracle π_s^i run by party P_i . Afterwards, except for one session involving F^* , \mathcal{A} uses Reveal on each accepted process $P_{i,s}$ of party P_i to get session key $k_{i,s}$.

Denote as P^* the PAKE party bound with F^* . Recall that $\text{Test}()$ can only be queried once in the constrained-PAKE game. Since the abort rule in **Game 2** is *not* triggered, now \mathcal{A} can successfully bind his unique test session T^* into the answers to \mathcal{C} when \mathcal{C} uploads F^* . \mathcal{A} uses the PAKE Test queries to get a challenge session key tk^* . Then \mathcal{A} queries for $f(P^*, T^*)$ and computes $k_R^* = tk_R^* \oplus f(P^*, T^*)$. Finally \mathcal{A} chooses an r^* and uses k_R^* as defined in SDoE protocol $\Pi_{\text{PAKE, re-enc}}$.

Let i be the original index of F^* in \mathbf{DB}_0 . \mathcal{A} outputs 1 if $b^* = b_i$ and 0 otherwise. Note that if $b = 0$ in the constrained PAKE experiment, the right half of tk^* is random a bit-string and so is $k_R^* = tk_R^* \oplus f(P^*, T^*)$. As a consequence, e is also random. Therefore in this case, \mathcal{C} also has no advantage. On the other hand, if $b = 1$, k_R^* is correctly distributed as in $\Pi_{\text{PAKE, re-enc}}$. The probability \mathcal{C} outputs the correct b^* is the same as \mathcal{A} outputs the correct b . Thus we have

$$\epsilon_{\mathcal{A}} = \text{Adv}_{\mathcal{C}}^{\text{Game 3}}(\lambda) \quad (5)$$

By combining (1) to (5), we have proved Theorem 1.

Security against compromised server. Next, we prove the security of $\Pi_{\text{PAKE, re-enc}}$ against a compromised server, which leads to the following theorem.

Theorem 2 *If there exists a ppt adversary \mathcal{S} in $\text{Exp}_{\mathcal{S}, \Pi_{\text{PAKE, re-enc}}}^{\text{SDoE}}(\lambda)$ with advantage $\epsilon_{\mathcal{S}}$ when $sh(F_0) = sh(F_1)$, then there also exist a ppt adversary \mathcal{A} with advantage $\epsilon_{\mathcal{A}}$ in the underlying IND-KPA game against Π_{enc} in the random oracle model and a ppt but passive adversary \mathcal{B} against the PAKE-protocol Π with advantage $\epsilon_{\mathcal{B}}$ such that*

$$\epsilon_{\mathcal{S}} \leq \frac{2C \cdot N_e}{|\mathbf{K}|} + N_e \cdot T \cdot \epsilon_{\mathcal{B}} + \left(\frac{2^{l_{sh}}}{q_H} \right)^2 (2|\mathbf{F}|^2 \cdot \epsilon_{\mathcal{A}})$$

where C is the maximal number of owners of each file, \mathbf{K} the key space of Π_{enc} , N_e the number of Execute queries, T the maximal number of PAKE sessions in each Execute query, \mathbf{F} the file space, l_{sh} the length of the short hash, q_H the number of distinct files that \mathcal{S} has queried for its hash or short hash and Π_{enc} is the encryption scheme for files.

Proof. First we consider two different cases for \mathcal{S} to win.

1. \mathcal{S} has issued $\text{Execute}(pid, P, F)$ and seen at least one file key collides into any of the file keys of the equivalent ciphertexts of F_0 or F_1 .
2. \mathcal{S} has not seen any colliding keys by issuing $\text{Execute}(pid, P, F)$.

In the first case, each $\text{Execute}(pid, P, F)$ reveals at most one real file key $k \in \mathbf{K}$. There are at most C owners of each file, and each of whom has an equivalent file key. Thus, seeing one key increases the probability of \mathcal{S} by at most $\frac{C}{|\mathbf{K}|}$ to decrypt each F_b correctly.

Let the advantage of \mathcal{S} in the second case be ϵ'_S . With the union bound we have

$$\epsilon_S \leq \frac{2C \cdot \mathbf{N}_e}{|\mathbf{F}|} + \epsilon'_S \quad (6)$$

To further analyze ϵ'_S , two types of adversaries are considered.

1. Adversaries recovered at least one complete session key generated by honest clients in the PAKE. We call these adversaries as type 1 adversaries.
2. Adversaries did not recover any complete session keys generated by honest clients. We call these adversaries as type 2 adversaries.

With a simple probability argument, it can be deduced that

$$\epsilon'_S \leq \epsilon_1 + \epsilon_2 \quad (7)$$

where ϵ_1 is the advantage of type 1 adversary and ϵ_2 is the advantage of type 2 adversary. Furthermore, we assume that the hash function is simulated by the challenger and all files are of equal length.

Next, we prove Theorem 2 by proving the following 2 lemmas.

Lemma 1 (*Bounding of the advantage of the type 1 adversary*) *If there exists any type 1 adversary \mathcal{A}_1 with advantage ϵ_1 and running time t_1 , then there also exists a constrained PAKE adversary \mathcal{B} with advantage ϵ_B and running time $t_B \approx t_1$ such that $\epsilon_1 \leq \mathbf{N}_e \cdot T \cdot \epsilon_B$.*

Proof. (lemma 1, sketch) \mathcal{B} can answer \mathcal{A}_1 's queries with the PAKE oracle, including transcripts of the Test-session and obtain at least one session key $sk_{\mathcal{A}_1}$ recovered by \mathcal{A}_1 . If $sk_{\mathcal{A}_1}$ has the same session-id owned by the Test-session, \mathcal{B} then outputs $(sk_{\mathcal{A}_1} = k_b)$, where k_b is the reply of PAKE Test(). Since there are at most $\mathbf{N}_e \cdot T$ sessions, \mathcal{B} wins with advantage $\epsilon_B \geq \frac{1}{\mathbf{N}_e \cdot T} \cdot \epsilon_1$.

Note that in our protocol, if the session key is leaked to \mathcal{A}_1 , then the encryption key k_F is also leaked to \mathcal{A}_1 and vice versa. The confidentiality of k_F is the basis of the remaining proof.

Lemma 2 (*Bounding of the advantage of the type 2 adversary*) *If there exists any type 2 adversary \mathcal{A}_2 with advantage ϵ_2 and running time t_2 , then there also exists a IND-KPA adversary \mathcal{A} with advantage ϵ_A and running time $t_A \approx t_2$ such that $\epsilon_2 \leq \left(\frac{2^{l_{sh}}}{q_H}\right)^2 (2|\mathbf{F}|^2 \cdot \epsilon_A)$.*

Proof. (lemma 2) **Game 0**. This is the original game $\mathbf{Exp}_{\mathcal{S}, \Pi_{\text{PAKE, re-enc}}}^{\text{SDoE}}(\lambda)$.

$$\epsilon_{\mathcal{A}_2} = \mathbf{Adv}_{\mathcal{S}}^{\text{Game } 0}(\lambda) \quad (8)$$

Game 1. If either F_0 or F_1 chosen by \mathcal{S} has a unique short hash value, abort the game. We add this rule since the short hash is also stored as part of the ciphertext. If any $sh(F_j)$ is unique, \mathcal{S} can simply learn F_j by computing and comparing the short hash values of all file candidates. Fix F_j and let **CollSH** be the event that $sh(F_j)$ does equal to any other $sh(F_i), F_i \in \mathbf{F}$. Then $\Pr[\mathbf{CollSH}] = \frac{q_H}{2^{l_{sh}}}$. Thus we have

$$\Pr[\exists F_i, F_k \in \mathbf{F}, F_i \neq F_0 \wedge F_k \neq F_1 : \\ sh(F_i) = sh(F_0) \wedge sh(F_k) = sh(F_1)] \geq \left(\frac{q_H}{2^{l_{sh}}}\right)^2 \quad (9)$$

Therefore we have

$$\mathbf{Adv}_{\mathcal{S}}^{\text{Game } 0}(\lambda) \leq \left(\frac{2^{l_{sh}}}{q_H}\right)^2 \mathbf{Adv}_{\mathcal{S}}^{\text{Game } 1}(\lambda) \quad (10)$$

Note that l_{sh} is sub-polynomial in λ so the loss factor is not exponential.

Game 2. The challenger guesses two files F_j and F_k . If $\{F_j, F_k\} \neq \{F_0, F_1\}$, abort the game. Thus

$$\mathbf{Adv}_{\mathcal{S}}^{\text{Game } 1}(\lambda) \leq (|\mathbf{F}|)^2 \mathbf{Adv}_{\mathcal{S}}^{\text{Game } 2}(\lambda) \quad (11)$$

Now we show how to construct \mathcal{A} against Π_{Enc} from \mathcal{S} . \mathcal{A} can guess $\{F_0, F_1\}$ since **Game 2** does not abort. In the setup phase, \mathcal{A} includes the public parameters of Π_{Enc} in the public parameters of $\Pi_{\text{PAKE, re-enc}}$, and queries for its own challenge ciphertext ϱ_b with $\{m_0 = F_0, m_1 = F_1\}$ in its IND-KPA game. \mathcal{A} fixes this ϱ_b as the ciphertext of F_0 and use other random keys (conforming to security parameter) to encrypt all other files as described in the model and this protocol⁸.

The **Send** and **RegisterCorrupt** queries can be answered as in the proof for security against compromised clients. For **AccessDB**, \mathcal{A} simply gives \mathcal{S} all the ciphertexts and owner lists at that time. Whenever a query from \mathcal{S} results in an observable database change (i.e., new ciphertexts are added or new owners are added to files), \mathcal{A} updates the server state and gives the ciphertext and/or the changed owner lists to \mathcal{S} . For **Execute** with $F \notin \{F_0, F_1\}$, \mathcal{A} can use the homomorphic property of Π_{Enc} to correctly generate all the transcript. Since \mathcal{A} knows all other keys and ciphertexts, \mathcal{A} can answer all the queries from \mathcal{S} .

If \mathcal{S} queries **Test**(), \mathcal{A} replies with ϱ_b of her own and outputs whatever \mathcal{S} outputs. Since the probability that \mathcal{A} correctly simulates the SDoE-game for \mathcal{S} is exactly $\frac{1}{2}$, we have

$$\mathbf{Adv}_{\mathcal{S}}^{\text{Game } 2}(\lambda) \leq 2\epsilon_{\mathcal{A}} \quad (12)$$

By combining (8) and (12), we have proved Lemma 2.

⁸ Recall that in the security game for compromised server, we assume that every file has been uploaded into \mathbf{DB}_0 , so are F_0 and F_1 .

By combining (6) and (7) and the two lemmas, we have proved Theorem 2.

4.2 PAKE-based deduplication on popular files

Our second scheme ($\Pi_{\text{PAKE,popular}}$) is shown in Figure 3. It tries to avoid using public-key operations to encrypt the entire file. The penalty is that it only deduplicates popular files. The idea is the same as $\Pi_{\text{PAKE,re-enc}}$, except that instead of deleting the duplicated files directly, \mathcal{S} keeps them until they become popular.

Note that, for unpopular files, the views of both \mathcal{S} and \mathcal{C} are similar to those in $\Pi_{\text{PAKE,re-enc}}$, except that XOR is used to replace addition and subtraction and a symmetric-key encryption scheme $E()$ is used to replace $F \cdot g^{k_F}$. So the security argument for $\Pi_{\text{PAKE,re-enc}}$ still holds for unpopular files here.

Deduplication effectiveness will be negatively affected if only popular files are deduplicated. In the next section, we show that this affection is small via simulations with realistic datasets.

5 Simulation

The authors in [14] did a realistic simulation to measure the deduplication effectiveness of Π_{PAKE} . They used a dataset comprising of Android application popularity data to represent the predominance of media files. We follow their simulation but with two improvements.

First, we expanded the data set to a more reasonable size since their dataset is relatively small (7 396 235 “upload requests” in total, of which 178 396 are for distinct files). In order to measure how the system behaves as the number of unique files increases, a larger dataset is needed. Since such data was not available, we used the Synthetic Minority Over-sampling (SMOTE) Technique [8] to generate extra samples. Given a set of input samples and the amount of required over-samplings, SMOTE performs the following for each input sample:

1. Compute x nearest neighbors for the input sample.
2. Randomly choose a neighbor and a point on the line segment joining the input sample to the selected neighbor. This point is a new, generated sample.
3. Repeat step 2 until the requested amount of over-sampling has been reached.

For example, if the amount of needed over-sampling is 200%, it will be repeated twice.

We used the (file size, popularity) pairs of the original dataset as the input samples in the SMOTE algorithm. The amount of over-sampling was 500% and for each input sample five nearest neighbors were considered when the new samples were computed. The hashes for the synthetic samples were chosen randomly. These new samples were combined with the samples from the original dataset into a *expanded* dataset. The expanded dataset contains 110 942 571 files of which 2 675 917 are unique. See Figure 4a for the file popularities of the original dataset and the expanded dataset.

Second, we adjust the distribution of upload request to better reflect the real world cases. In [14], they map the dataset to a stream of upload requests by

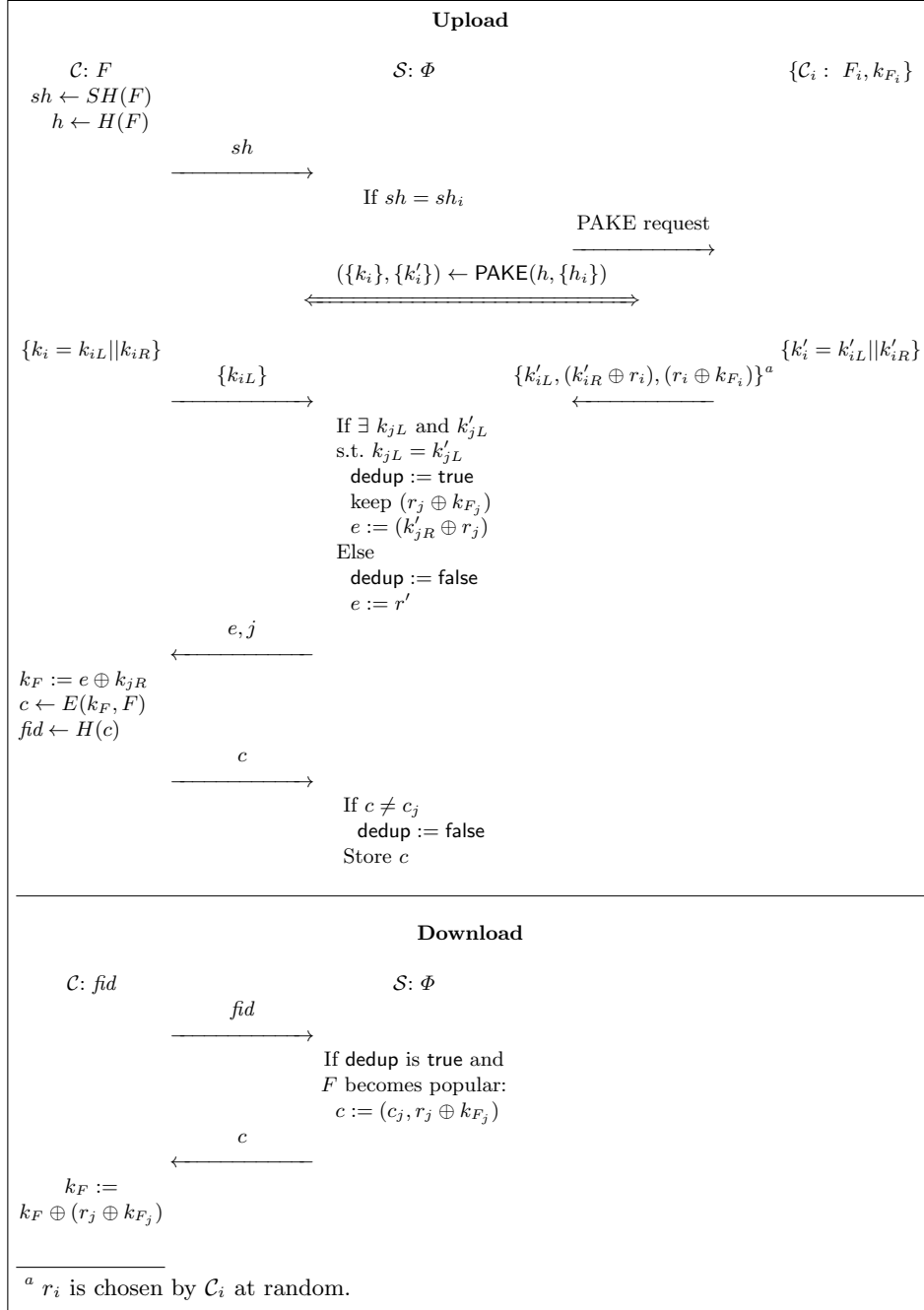


Fig. 3: PAKE-based deduplication on popular files.

generating the request in random order, i.e., a file that has x copies generates x upload requests that are uniformly distributed during the simulation. We argue that this cannot precisely capture the upload stream in real world: a file usually has less upload requests when it was generated, and becomes increasingly popular (more and more people hold it). To capture this case, we assume the upload requests of a single file follows normal distribution $\mathcal{N}(\mu, \sigma^2)$ where μ and σ are chosen randomly. Specifically, for a file F_i that has x_i total copies, the number of copies of F_i uploaded at time point t is $y_i = \frac{1}{\sigma_i \sqrt{2\pi}} e^{-\frac{(t-\mu_i)^2}{2\sigma_i^2}} x_i$. Then the total number of files uploaded at time point t is $\sum y_i$ and we assume that they are uploaded in random order. We do this for all time points and measure the final deduplication percentage.

Parameters. We follow [14], setting the number of possible files as 825 000, $l_{sh} = 13$ and $(n_{\text{RL}} + n'_{\text{RL}}) = 100$ (i.e., a \mathcal{C} will run PAKE at most 100 times for a certain file as both uploader and checker). We use these parameters in our simulations and measure deduplication effectiveness using the *deduplication percentage* ρ :

$$\rho = \left(1 - \frac{\text{Number of all files in storage}}{\text{Total number of upload requests}}\right) \cdot 100\% \quad (13)$$

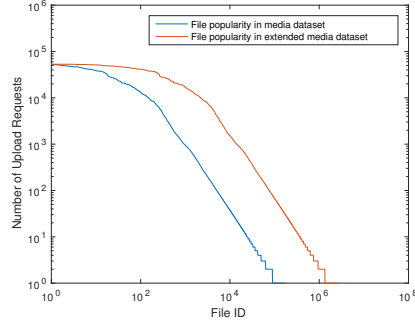
Rate limiting. We first assume that all \mathcal{C} s are online during the simulation and all files will be deduplicated (not limited to popular files). We run simulations with different combinations of RL_u and RL_c that satisfies $RL_u + RL_c = 100$, to see how selecting specific values for rate limits affects the deduplication effectiveness. Figure 4b shows that setting $RL_u = RL_c = 50$ maximises ρ to be 94.85%, which is close to the perfect deduplication percentage of 97.59%.

Offline rate. Note that \mathcal{C} s cannot participate in the deduplication protocol if they are offline, which may negatively affect deduplication effectiveness. To estimate this impact, we assign an *offline rate* to each \mathcal{C} as its probability to be offline during one run of the deduplication protocol. We set rate limits $RL_u = 50$ and $RL_c = 50$, and measured ρ by varying the offline rate. Figure 4c shows that ρ is still reasonably high even for relatively high offline rates of up to 70%, but drops quickly beyond that.

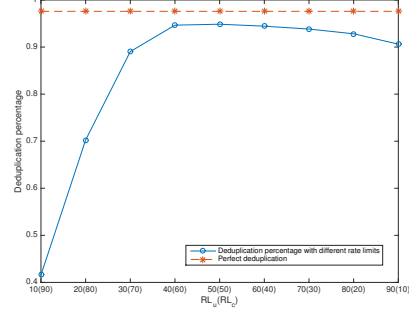
Popularity threshold. By far, all of the simulation results are for $\Pi_{\text{PAKE, re-enc}}$. Recall that $\Pi_{\text{PAKE, popular}}$ only deduplicates popular files which have a number of copies that are larger than a threshold, called *popularity threshold*. To investigate how this strategy affects deduplication effectiveness, we set rate limits $RL_u = 50$ ($RL_c = 50$), offline rate as 0.5, and run the simulation with different popularity thresholds. Figure 4d shows that ρ drops quickly if the popularity thresholds is larger than 32.

6 Related Work

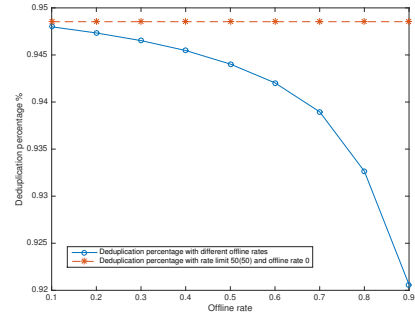
The first SDoE scheme is *convergent encryption* (CE) [10], which uses $H(F)$ as a key to encrypt F . In this way, different copies of F result in the same ciphertext.



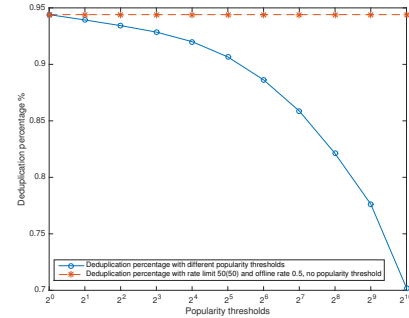
(a) file popularity.



(b) dedup. percentage vs. rate limits.



(c) dedup. percentage vs. offline rates.



(d) dedup. percentage vs. popularity thresholds.

Fig. 4: Simulation results.

However, a compromised passive \mathcal{S} can easily perform an offline brute-force attack over a predictable file. Bellare et al. recently formalized CE and proposed *message-locked encryption* (MLE) and its interactive version (iMLE) [2], which uses a semantically secure encryption scheme but produces a deterministic tag [4]. So it still suffers from the same attack. More recent work has attempted to improve MLE in several respects. Qin et al. [17] and Lei et al. [13] made MLE support *Rekeying* to protect key compromise and enable dynamic access control in the cloud storage. Zhao and Chow [21] proposed updatable MLE so that an encrypted file F can be efficiently updated with $O(\log|F|)$ computational cost. None of these improvements make MLE secure against offline brute-force attack.

DupLESS is a SDoE scheme that improves the security of CE against offline brute-force attacks [3]. In the key generation phase of CE, they introduce another secret which is provided by a third party and identical for all C_s . It adopts oblivious PRF to protect C_s ' files and the third party's secret. Duan [11] and Shin [18] later used decentralized architectures to distributed the trust of the

third party in DupLESS. *Cloudedup* is a SDoE scheme that introduces a third party for encryption and decryption [16]. Stanek et al. propose another SDoE scheme that only deduplicates popular files [20].

7 Conclusions

In this paper, we revisited the problem of secure deduplication of encrypted data (SDoE). We proposed a formal security model for this problem. We also proposed two single-server SDoE protocols and proved their security in our model. We showed that both of them can achieve reasonable deduplication effectiveness via simulations with realistic datasets.

Acknowledgments

This work was supported in part by TEKES - the Finnish Funding Agency for Innovation (CloSer project, 3881/31/2016) and by Intel (Intel Collaborative Research Institute for Secure Computing, ICRI-SC).

References

1. G. Ateniese, K. Fu, M. Green, and S. Hohenberger. Improved proxy re-encryption schemes with applications to secure distributed storage. *ACM Trans. Inf. Syst. Secur.*, 9(1):1–30, Feb. 2006.
2. M. Bellare and S. Keelveedhi. Interactive message-locked encryption and secure deduplication. In *IACR International Workshop on Public Key Cryptography*, pages 516–538. Springer, 2015.
3. M. Bellare, S. Keelveedhi, and T. Ristenpart. DupLESS: Server-aided encryption for deduplicated storage. In *USENIX Security*, pages 179–194. USENIX Association, 2013.
4. M. Bellare, S. Keelveedhi, and T. Ristenpart. Message-locked encryption and secure deduplication. In *EUROCRYPT*, volume 7881 of *LNCS*, pages 296–312. Springer, 2013.
5. M. Bellare, D. Pointcheval, and P. Rogaway. Authenticated key exchange secure against dictionary attacks. In B. Preneel, editor, *Advances in Cryptology — EUROCRYPT 2000: International Conference on the Theory and Application of Cryptographic Techniques Bruges, Belgium, May 14–18, 2000 Proceedings*, pages 139–155, Berlin, Heidelberg, 2000. Springer Berlin Heidelberg.
6. M. Bellare, D. Pointcheval, and P. Rogaway. Authenticated key exchange secure against dictionary attacks. In *EUROCRYPT*, pages 139–155, 2000.
7. S. M. Bellare and M. Merritt. Encrypted key exchange: password-based protocols secure against dictionary attacks. In *Research in Security and Privacy, 1992. Proceedings., 1992 IEEE Computer Society Symposium on*, pages 72–84, May 1992.
8. N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. Smote: Synthetic minority over-sampling technique. *J. Artif. Int. Res.*, 16(1):321–357, June 2002.
9. H. Dang and E. C. Chang. Privacy-preserving data deduplication on trusted processors. In *2017 IEEE 10th International Conference on Cloud Computing (CLOUD)*, pages 66–73, June 2017.

10. J. Douceur, A. Adya, W. Bolosky, P. Simon, and M. Theimer. Reclaiming space from duplicate files in a serverless distributed file system. In *Distributed Computing Systems, 2002. Proceedings. 22nd International Conference on*, pages 617–624, 2002.
11. Y. Duan. Distributed key generation for encrypted deduplication: Achieving the strongest privacy. In *CCSW*, pages 57–68. ACM, 2014.
12. D. Harnik, B. Pinkas, and A. Shulman-Peleg. Side channels in cloud services: Deduplication in cloud storage. *IEEE Security and Privacy*, 8(6):40–47, Nov. 2010.
13. L. Lei, Q. Cai, B. Chen, and J. Lin. *Towards Efficient Re-encryption for Secure Client-Side Deduplication in Public Clouds*, pages 71–84. Springer International Publishing, Cham, 2016.
14. J. Liu, N. Asokan, and B. Pinkas. Secure deduplication of encrypted data without additional independent servers. In *Proceedings of the 22Nd ACM SIGSAC Conference on Computer and Communications Security, CCS '15*, pages 874–885, New York, NY, USA, 2015. ACM.
15. J. Liu, N. Asokan, and B. Pinkas. Secure deduplication of encrypted data without additional independent servers. Cryptology ePrint Archive, Report 2015/455, 2015. <http://eprint.iacr.org/2015/455>.
16. P. Puzio, R. Molva, M. Onen, and S. Loureiro. Cloudedup: Secure deduplication with encrypted data for cloud storage. In *CloudCom*, pages 363–370. IEEE Computer Society, 2013.
17. C. Qin, J. Li, and P. P. C. Lee. The design and implementation of a rekeying-aware encrypted deduplication storage system. *Trans. Storage*, 13(1):9:1–9:30, Feb. 2017.
18. Y. Shin, D. Koo, J. Yun, and J. Hur. Decentralized server-aided encryption for secure deduplication in cloud storage. *IEEE Transactions on Services Computing*, PP(99):1–1, 2017.
19. V. Shoup. Sequences of games: a tool for taming complexity in security proofs. *IACR Cryptology ePrint Archive*, 2004:332, 2004.
20. J. Stanek, A. Sorniotti, E. Androulaki, and L. Kencl. A secure data deduplication scheme for cloud storage. In *FC*, pages 99–118, 2014.
21. Y. Zhao and S. S. Chow. Updatable block-level message-locked encryption. In *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, ASIA CCS '17*, pages 449–460, New York, NY, USA, 2017. ACM.