

# Interactive Message-Locked Encryption and Secure Deduplication

Mihir Bellare<sup>1</sup> and Sriram Keelveedhi

Dept. of Computer Science and Engineering, University of California San Diego, USA.

**Abstract.** This paper considers the problem of secure storage of out-sourced data in a way that permits deduplication. We are for the first time able to provide privacy for messages that are both correlated and dependent on the public system parameters. The new ingredient that makes this possible is interaction. We extend the message-locked encryption (MLE) primitive of prior work to interactive message-locked encryption (iMLE) where upload and download are protocols. Our scheme, providing security for messages that are not only correlated but allowed to depend on the public system parameters, is in the standard model. We explain that interaction is not an extra assumption in practice because full, existing deduplication systems are already interactive.

## 1 Introduction

THE SECURE DEDUPLICATION PROBLEM. Cloud storage providers such as Google, Dropbox and NetApp [31, 41, 51] derive significant cost savings from what is called *deduplication*. This means that if Alice and Bob upload the same data  $m$ , the service provider stores only one copy that is returned to Alice and Bob upon download.

Enter security, namely the desire of clients to keep their data private from the server. Certainly, Alice and Bob can conventionally encrypt their data under their passwords and upload the ciphertext rather than the plaintext. But then, even if they start from the same data  $m$ , they will end up with different ciphertexts  $C_A, C_B$ , foiling deduplication. The corresponding cost increase for the server would ultimately be passed to the clients in higher storage fees. It is thus in the interest of the parties to cooperate towards storage that is secure but deduplicatable.

Douceur et al. [30] provided the first solution, called convergent encryption (CE). The client encrypts its plaintext  $m$  with a *deterministic* symmetric encryption scheme under a  $k$  that is itself derived as a deterministic hash of the plaintext  $m$ . If Alice and Bob start with the same  $m$ , they will arrive at the same ciphertext, and thus deduplication is possible. Despite lacking an analysis until recently [12], CE has long been used in research and commercial systems [2, 4, 5, 18, 26, 27, 29, 35, 39, 46, 47, 52, 54], an indication of practitioners' interest in secure deduplication.

Scheme(s)	Type	Messages		STD/ROM
		Correlated	Param. dep.	
CE, HCE1, HCE2, RCE [12]	MLE	Yes	No	ROM
XtDPKE, XtESPKE, ... [12]	MLE	Yes	No	STD
BHK [10]	MLE	Yes	No	STD
ABMRS [1]	MLE	No	Yes	RO
FCHECK	iMLE	<b>Yes</b>	<b>Yes</b>	<b>STD</b>

**Fig. 1. Features of prior schemes (first four rows) and our scheme (last row).** We achieve security for the first time for messages that are *both* correlated *and* parameter dependent. Our scheme is in the standard model. The advance is made possible by exploiting interaction.

MLE. Bellare, Keelveedhi and Ristenpart (BKR) [12] initiated a theoretical treatment of secure deduplication aimed in particular at answering questions like, what security does CE provide and what can one prove about it? To this end they defined a primitive they called message-locked encryption (MLE). An MLE scheme specifies algorithms  $K, E, D, T$ . To encrypt  $m$ , let  $k \leftarrow K(p, m)$ , where  $p$  is a system-wide public parameter, and return ciphertext  $c \leftarrow E(k, m)$ . Decryption  $m \leftarrow D(k', c)$  recovers  $m$  as long as  $k' \leftarrow K(p, m)$  is any key derived from  $m$ . Tags, produced via  $t \leftarrow T(c)$ , are a way to test whether the plaintexts underlying two ciphertexts are the same or not, all encryptions of  $m$  having the same tag but it being hard to find differing plaintexts with matching tags.

Any MLE scheme enables deduplication. Alice, having  $m_A$ , computes and retains a key  $k_A \leftarrow K(p, m_A)$  and uploads  $c_A \leftarrow E(k, m_A)$ . The server stores  $c_A$ . Now Bob, having  $m_B$ , computes and retains a key  $k_B \leftarrow K(p, m_B)$  and uploads  $c_B \leftarrow E(k, m_B)$ . If the tags of  $c_A$  and  $c_B$  match, which means  $m_A = m_B$ , then the server deduplicates, storing only  $c_A$  and returning it to both Alice and Bob upon a download request. Both can decrypt to recover the common plaintext. CE is a particular MLE scheme in which key generation is done by hashing the plaintext.

MLE SECURITY. BKR [12] noted that MLE can only provide security for unpredictable data. (In particular, it cannot provide semantic security.) Within this range, two data dimensions emerge:

1. Correlation: Security holds even when messages being encrypted, although individually unpredictable, are related to each other.
2. Parameter-dependence: Security holds for messages that depend on the public parameters.

These dimensions are orthogonal, and the best would be security for correlated, parameter-dependent messages. This has not been achieved. What we have is schemes for correlated but parameter-independent messages [10, 12] and for non-correlated but parameter-dependent messages [1]. This past work is summarized in Fig. 1 and we now discuss it in a little more detail.

PRIOR SCHEMES. The definition of BKR [12], following [6], was security for correlated but parameter-independent messages. For this notion they proved security of CE in the ROM, gave new, secure ROM schemes, and made partial progress towards the challenging task of security without ROs. An efficient scheme in the standard model, also for correlated but parameter-independent messages, was provided in [10] assuming UCE-secure hash functions. (Specifically, against statistically unpredictable sources.)

Abadi, Boneh, Mironov, Raghunathan and Segev (ABMRS) [1] initiated treatment of security for parameter dependent messages, which they termed lock-dependent security. Achieving this is challenging. They gave a ROM solution that uses NIZK proofs to provide proofs of consistency. But to achieve security for parameter-dependent messages they were forced to sacrifice security for correlated messages. Their result assumes messages being encrypted are independently distributed.

QUESTIONS AND GOALS. The question we pose and address in this paper is, is it possible to achieve the best of both worlds, meaning security for messages that are both correlated *and* parameter dependent? This is important in practice. As indicated above, schemes for secure deduplication are currently deployed and in use in many systems [2,4,5,18,26,27,29,35,39,46,47,52,54]. In usage, messages are very likely to be correlated. For example, suppose Alice has uploaded a ciphertext  $c$  encrypting a paper  $m$  she is writing. She edits  $m$  to  $m'$ , and uploads the new version. The two plaintexts  $m, m'$  could be closely related, differing only in a few places. Also, even if messages of honest users are unlikely to depend on system parameters, attackers are not so constrained. Lack of security for parameter-dependent messages could lead to breaches. This is reflected for example in the BEAST attack on CBC in SSL/TLS [32]. We note that the question of achieving security for messages that are both correlated and parameter dependent is open both in the ROM and in the standard model.

CONTRIBUTIONS IN BRIEF. We answer the above questions by providing a deduplication scheme secure for messages that are both correlated and parameter dependent. Additionally, our scheme is standard-model, not ROM. The key new ingredient is interaction. In our solutions, upload and download are interactive protocols between the client and server. To specify and analyze these protocols, we define a new primitive, interactive MLE or iMLE. We provide a syntax and definitions of security, then specify and prove correct our protocols.

iMLE turns out to be interesting in its own right and yields some other benefits. We are able to provide the first secure deduplication scheme that permits incremental updates. This means that if a client's message changes only a little, for example due to an edit to a file, then, rather than create and upload an entirely new ciphertext, she can update the existing one with communication cost proportional only to the distance between the new and old plaintexts. This is beneficial because communication is a significant fraction of the operating expenditure in outsourced storage services. For example, transferring one gigabyte to the server costs as much storing one gigabyte for a month or longer in popular storage services [3, 40, 49]. In particular, backup systems, an important use case

for deduplication, are likely to benefit, as the operations here are incremental by nature. Incremental cryptography was introduced in [8, 9] and further studied in [14, 22, 34, 50].

INTERACTION? One might question the introduction of interaction. Isn't a non-interactive solution preferable? Our answer is that we don't "introduce" interaction. It is already present. Upload and download in real systems is inherently and currently interactive, even in the absence of security. MLE is a cryptographic core, not a full deduplication system. If MLE is used for secure deduplication, the uploads and downloads will be interactive, even though MLE is not, due to extra flows that the full system requires. Interaction being already present, it is natural to exploit it for security. In doing so, we are taking advantage of an existing resource rather than introducing an entirely new one.

MLE considered a single client. But in a full deduplication system, there are multiple clients concurrently executing uploads and downloads. Our iMLE model captures this. iMLE is thus going further towards providing security of the full system rather than just a cryptographic core. We know from experience that systems can fail in practice even when a "proven-secure" scheme is used if the security model does not encompass the full range of attacker capabilities or security goals of the implementation. Modeling that penetrates deeper into the system, as with iMLE, increases assurance in practice.

We view iMLE as a natural extension of MLE. The latter abstracted out an elegant primitive at the heart of the secure deduplication problem that could be studied in isolation. We study the full deduplication system, leveraging MLE towards full solutions with added security features.

DUPLICATE FAKING. In a duplicate faking attack, the adversary concocts and uploads a perverse ciphertext  $c^*$  with the following property. When honest Alice uploads an encryption  $c$  of her message  $m$ , the server's test (wrongly) indicates that the plaintexts underlying  $c^*$ ,  $c$  are the same, so it discards  $c$ , returning  $c^*$  to Alice upon a download request. But when Alice decrypts  $c^*$ , she does not get back her original plaintext.

Beyond privacy, BKR [12] defined an integrity requirement for MLE called tag consistency whose presence provides security against duplicate faking attacks. The important tag consistency property is possessed by the prior MLE schemes of Fig. 1 and also by our new iMLE schemes.

Deterministic schemes provide tag consistency quite easily and naturally. But ABMRS [1] indicate that security for parameter-dependent messages requires randomization. Tag consistency now becomes challenging to achieve. Indeed, providing it accounts for the use of NIZKs and the corresponding cost and complexity of the ABMRS scheme [1].

In the interactive setting, we capture the requirement underlying tag consistency by a recovery condition that is part of our soundness definition and requirement. Soundness in particular precludes duplicate faking attacks in the interactive setting. Our scheme provides soundness, in addition to privacy for messages that are both correlated and parameter dependent. Our FCHECK solution uses composable point function obfuscation [17] and FHE [19–21, 28, 36, 38, 55].

CLOSER LOOK. We look in a little more detail at the main definitional and scheme contributions of our work.

Public parameters for an iMLE scheme are created by an `Init` algorithm. Subsequently, a client can register (`Reg`), upload (`Put`) and download (`Get`). Incremental schemes have an additional update (`Upd`). All these are interactive protocols between client and server. For soundness, we ask that deduplication happens as expected and that clients can recover their uploaded files even in the presence of an attacker which knows all the files being uploaded and also read the server’s storage at any moment. The latter condition protects against duplicate-faking attacks. Our security condition is modeled on that of BKR [12] and requires privacy for correlated but individually unpredictable messages that may depend on the public parameters.

Our `FCHECK` construction, described and analyzed in Section 4, achieves soundness as well as privacy for messages that are both correlated and parameter dependent, all in the standard model, meaning without recourse to random oracles. The construction builds on a new primitive we call MLE-Without-Comparison (MLEWC). As the name indicates, MLEWC schemes are similar to MLE schemes in syntax and functionality, except that they do not support comparison between ciphertexts. We show that MLEWC can be realized in the standard model, starting from point function obfuscation [17] or, alternatively, UCE-secure hash function families [10]. However, comparison is essential to enable deduplication. To enable comparison, `FCHECK` employs an interactive protocol using a fully homomorphic encryption (FHE) scheme [19–21, 28, 36, 38, 55], transforming the MLEWC scheme into an iMLE scheme.

We then move on to the problem of incremental updates. Supporting incremental updates over MLE schemes turns out to be challenging: deterministic MLE schemes cannot support incremental updates, as we show in the full version [11], while randomized MLE schemes seem to need complex machinery such as NIZK proofs of consistency [1] to support incremental updates while retaining the same level of security as deterministic schemes, which makes them unfit for practical usage. We show how interaction can be exploited to solve this problem. We describe an efficient ROM scheme `IRCE` that supports incremental updates. The scheme, in its simplest form, works like the randomized convergent encryption (RCE) scheme [12], where the message is encrypted with a random key using a blockcipher in counter (`CTR`) mode, and the random key is encrypted with a key derived by hashing the message. We show that this indirection enables incremental updates. However, RCE does not support strong tag consistency and hence cannot offer strong security against duplicate faking attacks. We overcome this in `IRCE` by including a simple response from the server as part of the upload process. We remark that `IRCE` is based off a core MLE (non-interactive) scheme permitting incremental updates, interaction being used only for tag consistency.

<pre> Run(<math>1^\lambda, P, \text{inp}</math>) <math>n \leftarrow 1; i \leftarrow 1; M \leftarrow \epsilon; \mathbf{a}[1, 1] \leftarrow \text{inp}[1]; \mathbf{a}[2, 1] \leftarrow \text{inp}[2]</math> While <math>T[n] = \text{False}</math>   (<math>\mathbf{a}[n, i + 1], M, T[n]</math>) <math>\leftarrow_s P[n, i](1^\lambda, \mathbf{a}[n, i], M)</math>   If <math>n = 2</math> then <math>n \leftarrow 1; i \leftarrow i + 1</math> else <math>n \leftarrow 2</math> Ret <math>\text{last}(\mathbf{a}[1]), \text{last}(\mathbf{a}[2])</math> </pre>
<pre> Msgs(<math>1^\lambda, P, \text{inp}, r</math>) <math>n \leftarrow 1; i \leftarrow 1; j \leftarrow 1; \mathbf{a}[1, 1] \leftarrow \text{inp}[1]; \mathbf{a}[2, 1] \leftarrow \text{inp}[2]; M \leftarrow \epsilon</math> While <math>T[n] = \text{False}</math>   (<math>\mathbf{a}[n, i + 1], M, T[n]</math>) <math>\leftarrow_s P[n, i](1^\lambda, \mathbf{a}[n, i], M; r[n, i]); M[j] \leftarrow M; j \leftarrow j + 1</math>   If <math>n = 2</math> then <math>n \leftarrow 1; i \leftarrow i + 1</math> else <math>n \leftarrow 2</math> Ret <math>M</math> </pre>

**Fig. 2. Top:** Running a 2-player protocol  $P$ . **Bottom:** The  $\text{Msgs}$  procedure returns the messages exchanged during the protocol when invoked with specified inputs and coins.

## 2 Preliminaries

We let  $\lambda \in \mathbb{N}$  and  $1^\lambda$  denote the security parameter and its unary representation. The empty string is denoted by  $\epsilon$ . We let  $|S|$  denote the size of a finite set  $S$  and let  $s \leftarrow_s S$  denote sampling an element from  $S$  at random and assigning it to  $s$ . If  $a, b \in \mathbb{N}$  and  $a < b$ , then  $[a]$  denotes the set  $\{1, \dots, a\}$  and  $[a, b]$  denotes the set  $\{a, \dots, b\}$ . For a tuple  $\mathbf{x}$ , we let  $|\mathbf{x}|$  denote the number of components in  $\mathbf{x}$ , and  $\mathbf{x}[i]$  denote the  $i$ -th component, and  $\text{last}(\mathbf{x}) = \mathbf{x}[|\mathbf{x}|]$ , and  $\mathbf{x}[i, j] = \mathbf{x}[i] \dots \mathbf{x}[j]$  for  $1 \leq i \leq j \leq |\mathbf{x}|$ . A binary string  $s$  is identified with a tuple over  $\{0, 1\}$ . The guessing probability of a random variable  $X$ , denoted by  $\mathbf{GP}(X)$ , is defined as  $\mathbf{GP}(X) = \max_x \Pr[X = x]$ . The conditional guessing probability  $\mathbf{GP}(X | Y)$  of a random variable  $X$  given a random variable  $Y$  are defined via  $\mathbf{GP}(X | Y) = \sum_y \Pr[Y = y] \cdot \max_x \Pr[X = x | Y = y]$ .

The Hamming distance between  $s_1, s_2 \in \{0, 1\}^\ell$  is given by  $\text{HAMM}(s_1, s_2) = \sum_{i=1}^\ell (s_1[i] \oplus s_2[i])$ . We let  $\text{patch}_{\text{HAMM}}(s_1, \delta)$  be the string  $s$  such that  $s[i] = s_1[i]$  if  $i \notin \delta$  and  $s[i] = \neg s_1[i]$  if  $i \in \delta$  and  $\text{diff}_{\text{HAMM}}(s_1, s_2) = \{i : s_1[i] \neq s_2[i]\}$ .

Algorithms are randomized and run in polynomial time (denoted by PT) unless otherwise indicated. We let  $y \leftarrow A(a_1, \dots; r)$  denote running algorithm  $A$  on  $a_1, \dots$  with coins  $r$  and assigning the output to  $y$ , and let  $y \leftarrow_s A(a_1, \dots)$  denote the same operation with random coins. We let  $[A(a_1, \dots)]$  denote the set of all  $y$  that have non-zero probability of being output by  $A$  on inputs  $a_1, \dots$ . Adversaries are either algorithms or tuples of algorithms. A negligible function  $f$  approaches zero faster than the polynomial reciprocal; for every polynomial  $p$ , there exists  $n_p \in \mathbb{N}$  such that  $f(n) \leq 1/p(n)$  for all  $n \geq n_p$ .

We use the code-based game playing framework of [16] along with extensions of [53] and [12] when specifying security notions and proofs.

A two player  $q$ -round protocol  $P$  is represented through a  $2 \times q$ -tuple  $(P[i, j])_{i \in [2], j \in [q]}$  of algorithms where  $P[i, j]$  represents the action of the  $i$ -th player invoked for the  $j$ -th time. We let  $P[1]$  denote the player who initiates

the protocol, and  $P[2]$  denote the other player. Each algorithm is invoked with  $1^\lambda$ , an input  $\mathbf{a}$ , and a message  $M \in \{0, 1\}^*$ , and returns a 3-tuple consisting of an output  $\mathbf{a}'$ , an outgoing message  $M' \in \{0, 1\}^*$ , and a boolean  $T$  to indicate termination. The `Run` algorithm (Fig. 2) captures the execution of  $P$ , and `Msgs` (Fig. 2) returns the messages exchanged in an instance of  $P$ , when invoked with specified inputs and coins.

**ADVERSARIAL MODEL.** A secure deduplication system (built from an iMLE scheme) will operate in a setting with a server and several clients. Some clients will be controlled by an attacker, while others will be legitimate, belonging to honest users and following the protocol specifications. A resourceful attacker, apart from controlling clients, could gain access to server storage, and interfere with communications. Our adversarial model captures an iMLE scheme running in the presence of an attacker with such capabilities.

We now walk through an abstract game  $G$ , and explain how this is achieved. The games in the rest of the paper, for soundness, security, and other properties of iMLE largely follow this structure. The game  $G$  sets up and controls a server instance. The adversary  $A$  is invoked with access to a set of procedures. Usually, the objective of the game involves  $A$  violating some property guaranteed to legitimate clients like  $L$ , such as ability to recover stored files, or privacy of data.

The `MSG` procedure can send arbitrary messages to the server and can be used to create multiple clients, and run multiple instances of protocols, which could deviate from specifications.

The `INIT` and `STEP` procedures control a single legitimate client  $L$ . The `INIT` procedure starts protocol instances on behalf of  $L$ , using inputs of  $A$ 's choice. The `STEP` procedure advances a protocol instance by running the next algorithm. Together, these procedures let  $A$  run several legitimate and corrupted protocol instances concurrently.

The `STATE` procedure returns the server's state, which includes stored ciphertexts, public parameters, etc.. In some games, it also returns the state and parameters of  $L$ . `STATE` provides only read access to the server's storage. This restriction is necessary. If  $A$  is allowed to modify the storage of the server, then it can always tamper with the data stored by the clients, making secure deduplication impossible.

We assume that  $A$  can read, delay and drop messages between the server and legitimate clients. However,  $A$  cannot tamper with message contents, reorder messages within a protocol, or redirect messages from one protocol instance to another. This assumption helps us simplify the protocol descriptions and proofs. Standard, efficient techniques can be used to transform the protocols from this setting to be secure in the presence of an attacker that can tamper and reorder messages [7].

### 3 Interactive message-locked encryption

**DEFINITION.** An interactive message-locked encryption scheme iMLE consists of an initialization algorithm `Init` and three protocols `Reg`, `Put`, `Get`. Initialization

<u>MAIN</u> ( $1^\lambda$ ) // <u>REG</u> <sub>IMLE</sub> <sup>A</sup> ( $1^\lambda$ )	<u>STEP</u> ( $j$ ) // Advance by one step.
$\text{win} \leftarrow \text{False}; \sigma_S \leftarrow \text{\$ Init}(1^\lambda)$	$\mathbf{P} \leftarrow \mathbf{PS}[j]; n \leftarrow \mathbf{N}[j]; i \leftarrow \mathbf{rd}[j]$
$\mathbf{A}^{\text{REG,INIT,STEP,MSG,STATE}}(1^\lambda); \text{Ret win}$	If $\mathbf{T}[j, n]$ then return $\perp$
<u>REG</u> // Set up the legitimate client $L$ .	If $n = 2$ then $\text{inp} \leftarrow \sigma_S$ else $\text{inp} \leftarrow \mathbf{a}[j, i]$
$(\sigma_C, \sigma_S) \leftarrow \text{\$ Run}(\text{Reg}, \epsilon, \sigma_S)$	$(\text{outp}, \mathbf{M}[j], \mathbf{T}[j, n]) \leftarrow \text{\$ P}[n, i](1^\lambda, \text{inp}, \mathbf{M}[j])$
<u>INIT</u> ( $\mathbf{P}, \text{inp}$ ) // Start a protocol with $L$ .	If $n = 2$ then
If $\mathbf{P} \notin \{\text{Put}, \text{Get}\}$ then ret $\perp$	$\sigma_S \leftarrow \text{outp}; \mathbf{N}[j] \leftarrow 1; \mathbf{rd}[j] \leftarrow \mathbf{rd}[j] + 1$
$\mathbf{p} \leftarrow \mathbf{p} + 1; j \leftarrow \mathbf{p}; \mathbf{PS}[j] = \mathbf{P}$	Else $\mathbf{a}[j, i + 1] \leftarrow \text{outp}; \mathbf{N}[j] \leftarrow 2$
$\mathbf{a}[j, 1] \leftarrow \text{inp}; \mathbf{N}[j] \leftarrow 1; \mathbf{M}[j] \leftarrow \epsilon; \text{Ret } j$	If $\mathbf{T}[j, 1] \wedge \mathbf{T}[j, 2]$ then <u>WINCHECK</u> ( $j$ )
	Ret $\mathbf{M}[j]$
<u>MSG</u> ( $\mathbf{P}, i, \mathbf{M}$ ) // Send message to server.	<u>WINCHECK</u> ( $j$ ) // Check if $\mathbf{A}$ has won.
If $\mathbf{P} \notin \{\text{Reg}, \text{Put}, \text{Get}, \text{Upd}\}$ then ret $\perp$	If $\mathbf{PS}[j] = \text{Put}$ then
$(\sigma_S, \mathbf{M}, \mathbf{N}, \mathbf{T}) \leftarrow \text{\$ P}[2, i](1^\lambda, \sigma_S, \mathbf{M}); \text{Ret } \mathbf{M}$	$(\sigma_C, m) \leftarrow \mathbf{a}[j, 1]$
	$f \leftarrow \text{last}(\mathbf{a}[j]); T[f] \leftarrow m$
	If $\mathbf{PS}[j] = \text{Get}$ then
	$(\sigma_C, f) \leftarrow \mathbf{a}[j, 1]; m' \leftarrow \text{last}(\mathbf{a}[j])$
	$\text{win} \leftarrow \text{win} \vee (m' \neq T[f])$

**Fig. 3.** The REC game. The STATE procedure returns  $\sigma_S, \sigma_C$ .

Init sets up server-side state:  $\sigma_S \leftarrow \text{\$ Init}(1^\lambda)$ . Each protocol  $\mathbf{P}$  consists of two players - a client  $\mathbf{P}[1]$  (meaning that the client always initiates), and a server  $\mathbf{P}[2]$ . All server-side algorithms  $\mathbf{P}[2, \cdot]$  take server-side state  $\sigma_S$  as input, and produce an updated state  $\sigma'_S$  as output. The Reg protocol registers new users; here,  $\text{Reg}[1]$  takes no input and returns client parameters  $\sigma_C \in \{0, 1\}^*$ . The Put protocol stores files on the server; here,  $\text{Put}[1]$  takes plaintext  $m \in \{0, 1\}^*$  and  $\sigma_C$  as inputs, and outputs an identifier  $f \in \{0, 1\}^*$ . The Get protocol retrieves files from the server; here,  $\text{Get}[1]$  takes identifier  $f$  and  $\sigma_C$  as inputs, and outputs plaintext  $m \in \{0, 1\}^*$ .

**SOUNDNESS.** We require two conditions. First is deduplication, meaning that if a client puts a ciphertext of a file already on the server, then the storage should not grow by the size of the file. A small increase towards book-keeping information, that is independent of the size of the file, is permissible. More precisely, there exists a bound  $\ell : \mathbb{N} \rightarrow \mathbb{N}$  such that for all server-side states  $\sigma_S \in \{0, 1\}^*$ , for all valid client parameters (derived through Reg with fresh coins)  $\sigma_C, \sigma'_C$ , for all  $m \in \{0, 1\}^*$ , the expected increase in size of  $\sigma''_S$  over  $\sigma'_S$  when  $(f', \sigma'_S) \leftarrow \text{\$ Run}(\text{Put}, (\sigma_C, m), \sigma_S)$  and  $(f', \sigma''_S) \leftarrow \text{\$ Run}(\text{Put}, (\sigma'_C, m), \sigma'_S)$  is bounded by  $\ell(\lambda)$ .

The second condition is correct recovery of files: if a legitimate client puts a file on the server, it should be able to get the file later. We formalize this requirement by the REC game of Fig. 3, played with an adversary  $\mathbf{A}$ , which gets



access to procedures REG, INIT, STEP, MSG, STATE. We provide an overview of these procedures here.

The REG procedure sets up a legitimate client  $L$  by running  $\text{Run}(1^\lambda, \text{Reg}, (\epsilon, \sigma_S))$ . The INIT procedure lets  $A$  run protocols on behalf of  $L$ . It takes input  $\text{inp}$  and  $P$ , where  $P$  has to be one of Put, Get, and  $\text{inp}$  should be the a valid input for  $P[1, 1]$ . A new instance of  $P$  is set up, and  $A$  is returned  $j \in \mathbb{N}$ , an index to the instance. The STEP procedure takes input  $j$ , advances the instance by one algorithm unless the current instance has terminated. The outgoing message is returned to  $A$ . The inputs and outputs of the protocol steps are all stored in an array  $\mathbf{a}$ . The STATE procedure returns  $\sigma_S, \sigma_C$ .

If an instance  $j$  has terminated, then STEP runs WINCHECK, which maintains a table  $T$ . If  $j$  is an instance of Put, then  $m$  and identifier  $f$  are recovered from  $\mathbf{a}[j]$  and  $T[f]$  gets  $m$ . If  $j$  is an instance of Get, then WINCHECK obtains  $f$  and the recovered plaintext  $m'$ , and checks if  $T[f] = m'$ . If this fails, either because  $T[f]$  is some value different from  $m'$ , or is undefined, then WINCHECK sets the win flag, which is the condition for  $A$  to win the game. We associate advantage  $\text{Adv}_{\text{iMLE}, A}^{\text{rec}}(\lambda) = \Pr[\text{REC}_{\text{iMLE}}^A(1^\lambda)]$  with iMLE and  $A$ . For recovery correctness, we require that the advantage should be negligible for all PT  $A$ .

SECURITY. The primary security requirement for iMLE schemes is privacy of unpredictable data. Unpredictability (plaintexts drawn from a distribution with negligible guessing probability) is a prerequisite for privacy in MLE schemes [12], as without unpredictability, a simple brute-force attack can recover the contents of a ciphertext by generating keys from all candidate plaintexts and checking if decrypting the ciphertext with the key leads back to the candidate plaintext. A similar argument extends unpredictability as a requirement to secure deduplication schemes as well. We formalize unpredictability as follows.

A source  $S$  is an algorithm that on input  $1^\lambda$  and a string  $d \in \{0, 1\}^*$  returns a pair of tuples  $(\mathbf{m}_0, \mathbf{m}_1)$ . There exist  $m : \mathbb{N} \rightarrow \mathbb{N}$  and  $\ell : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$  such that  $|\mathbf{m}_0| = |\mathbf{m}_1| = m(\lambda)$ , and  $|\mathbf{m}_0[i]| = |\mathbf{m}_1[i]| = \ell(\lambda, i)$  for all  $i \in [m(\lambda)]$ . All components of  $\mathbf{m}_0$  and  $\mathbf{m}_1$  are unique. The guessing probability  $\mathbf{GP}_S(\lambda)$  of  $S$  is defined as  $\max_{i,b,d}(\mathbf{GP}(\mathbf{m}_b[i]))$  when  $(\mathbf{m}_0, \mathbf{m}_1) \leftarrow S(1^\lambda, d)$ . We say that  $S$  is unpredictable if  $\mathbf{GP}_S(\cdot)$  is negligible. We say that  $S$  is a single source if it only outputs one tuple, but satisfies the other conditions. We say that  $S$  is an auxiliary source if it outputs a string  $z \in \{0, 1\}^*$  along with  $\mathbf{m}_0, \mathbf{m}_1$  and if it holds that guessing probability conditioned on  $z$  is negligible.

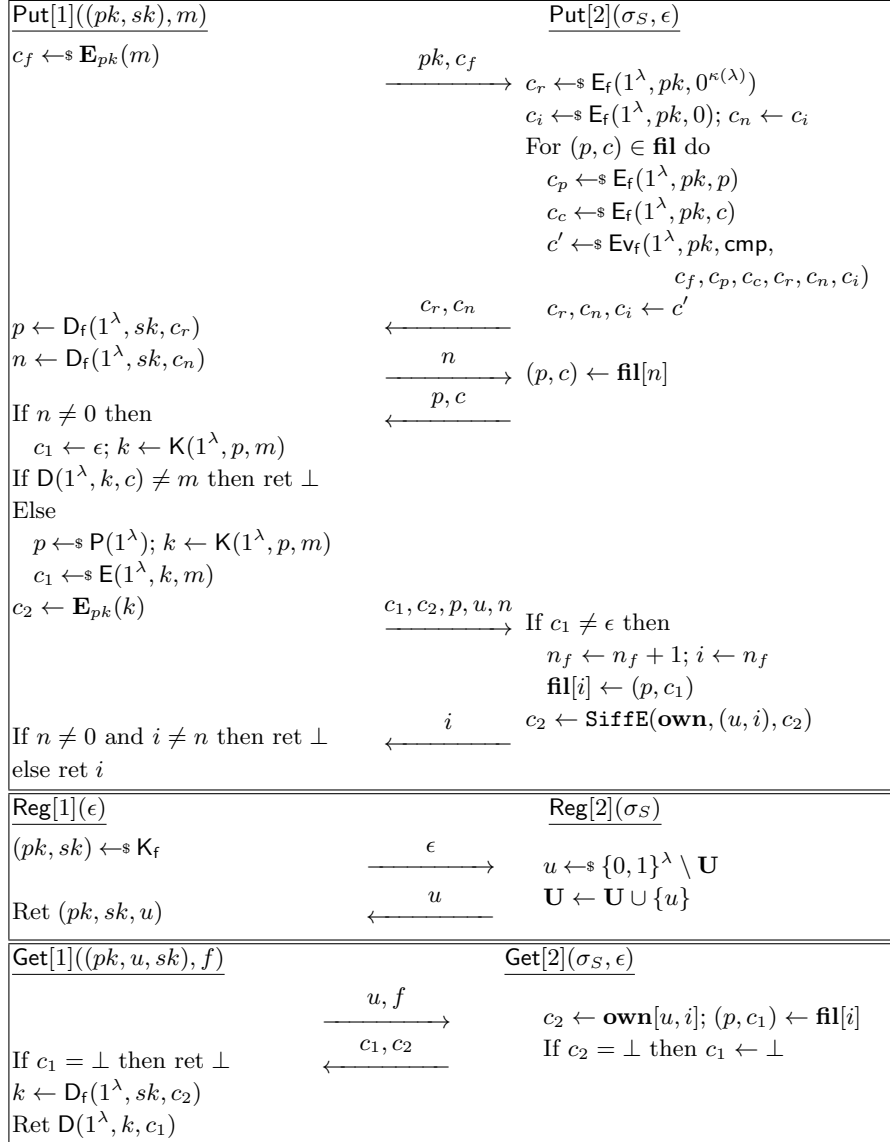
The PRIV game of Fig. 4, associated with iMLE, a source  $S$  and an adversary  $A$ , captures privacy for unpredictable messages independent of the public parameters of the system. As with REC, the game starts by running  $\sigma_S \leftarrow \text{Init}(1^\lambda)$  to set up the server-side state. The game then runs  $S$  to get  $(\mathbf{m}_0, \mathbf{m}_1)$ , picks a random bit  $b$ , and uses  $\mathbf{m}_b$  as messages to be put on the server. Then,  $A$  is invoked with access to REG, PUT, STEP, MSG and STATE. The REG, STATE, and MSG oracles behave in the same way as in REC. The STEP oracle here is similar to that of REC, except that it does not invoke WINCHECK. Adversary  $A$  can initialize an instance of Put with a plaintext  $\mathbf{m}_b[i]$  by calling PUT( $i$ ).

<u>MAIN</u> ( $1^\lambda$ ) // <u>PRIV</u> <sup>S,A</sup> ( $1^\lambda$ )	<u>MAIN</u> ( $1^\lambda$ ) // <u>PDPRIV</u> <sup>S,A</sup> ( $1^\lambda$ )
$b \leftarrow_s \{0, 1\}; \mathbf{p} \leftarrow 0; \sigma_S \leftarrow_s \text{Init}(1^\lambda)$	$b \leftarrow_s \{0, 1\}; \mathbf{p} \leftarrow 0; \sigma_S \leftarrow_s \text{Init}(1^\lambda)$
$\mathbf{m}_0, \mathbf{m}_1 \leftarrow_s \mathbf{S}(1^\lambda, \epsilon)$	$b' \leftarrow_s \mathbf{A}^{\text{PTXT,PUT,UPD,STEP,MSG,REG,STATE}}(1^\lambda)$
$b' \leftarrow_s \mathbf{A}^{\text{PUT,UPD,STEP,MSG,REG,STATE}}(1^\lambda)$	Ret ( $b = b'$ )
Ret ( $b = b'$ )	<u>PTXT</u> ( $d$ )
<u>REG</u>	$\mathbf{m}_0, \mathbf{m}_1 \leftarrow_s \mathbf{S}(1^\lambda, d)$
$(\sigma_C, \sigma_S) \leftarrow_s \text{Run}(\text{Reg}, \epsilon, \sigma_S)$	<u>MSG</u> ( $\mathbf{P}', \mathbf{M}$ ) // Send a message to the server
<u>PUT</u> ( $i$ ) // Start a Put instance	If $\mathbf{P}' \notin \{\text{Reg}, \text{Put}, \text{Get}, \text{Upd}\}$ then ret $\perp$
$\mathbf{p} \leftarrow \mathbf{p} + 1; \mathbf{PS}[\mathbf{p}] = \text{Put}$	$(\sigma_S, \mathbf{M}, \mathbf{T}) \leftarrow_s \mathbf{P}'[2](1^\lambda, \sigma_S, \mathbf{M})$
$\mathbf{a}[\mathbf{p}, 1] \leftarrow \mathbf{m}_b[i]$	Ret ( $\sigma_S, \mathbf{M}, \mathbf{N}, \mathbf{T}$ )
$\mathbf{N}[\mathbf{p}] \leftarrow 1; \mathbf{M}[\mathbf{p}] \leftarrow \epsilon; \text{Ret } \mathbf{p}$	<u>STEP</u> ( $j$ ) // Advance instance by a step.
<u>STATE</u>	$\mathbf{P} \leftarrow \mathbf{PS}[j]; n \leftarrow \mathbf{N}[j]; i \leftarrow \text{rd}[j]$
If <b>cheat</b> = <b>False</b> then	If $\mathbf{T}[j, n]$ or <b>done</b> then return $\perp$
<b>done</b> $\leftarrow$ <b>True</b> ; ret $\sigma_S$	If $n = 2$ then $\text{inp} \leftarrow \sigma_S$ else $\text{inp} \leftarrow \mathbf{a}[j, i]$
else ret $\perp$	$(\text{outp}, \mathbf{M}[j], \mathbf{T}[j, n]) \leftarrow_s \mathbf{P}[n, i](1^\lambda, \text{inp}, \mathbf{M}[j])$
	If $n = 2$ then
	$\sigma_S \leftarrow \text{outp}; \mathbf{N}[j] \leftarrow 1; \text{rd}[j] \leftarrow \text{rd}[j] + 1$
	else $\mathbf{a}[j, i + 1] \leftarrow \text{outp}; \mathbf{N}[j] \leftarrow 2$
	If $n = 1$ and $\mathbf{T}[j, n]$ then
	$T_f[\mathbf{a}[j, 1]] \leftarrow \text{last}(\mathbf{a}[j])$

**Fig. 4.** The PRIV and PDPRIV security games. Apart from MAIN, the games share the same code for all procedures. The PDPRIV game has an additional PTXT procedure.

We associate advantage  $\text{Adv}_{\text{iMLE}, \mathbf{S}, \mathbf{A}}^{\text{priv}}(\lambda) = 2 \Pr[\text{PRIV}_{\text{iMLE}}^{\mathbf{S}, \mathbf{A}}(1^\lambda)] - 1$  with a iMLE a source  $\mathbf{S}$  and an adversary  $\mathbf{A}$ . We require that the advantage should be negligible for all PT  $\mathbf{A}$  for all unpredictable PT  $\mathbf{S}$ .

The PDPRIV game of Fig. 4 extends PRIV-security to messages depending on the public parameters of the system, a notion termed lock-dependent security in [1]. Here, we term this parameter-dependent security. In this game, the adversary  $\mathbf{A}$  gets access to a PTXT procedure, which runs  $\mathbf{S}(1^\lambda, \sigma_S)$  to get  $\mathbf{m}_0, \mathbf{m}_1$ . The other procedures follow PRIV. A simpler approach is to run  $\mathbf{S}$  with  $\sigma_S$  when the game starts (i.e. in main) as in PRIV. However, this leads to trivial constructions where **Init** is a dummy procedure, and the system parameters are generated when the first client registers. This is avoided in PDPRIV by letting  $\mathbf{A}$  decide, through PTXT, when  $\mathbf{S}$  is to be run. We associate advantage  $\text{Adv}_{\text{iMLE}, \mathbf{S}, \mathbf{A}}^{\text{ldpriv}}(\lambda) = 2 \Pr[\text{PDPRIV}_{\text{iMLE}}^{\mathbf{S}, \mathbf{A}}(1^\lambda)] - 1$  with a scheme iMLE a source  $\mathbf{S}$  and an adversary  $\mathbf{A}$ . We require that advantage should be negligible for all PT  $\mathbf{A}$  for all unpredictable PT  $\mathbf{S}$ .



**Fig. 5.** The FCHECK scheme over FHE = (K<sub>f</sub>, E<sub>f</sub>, D<sub>f</sub>, E<sub>v<sub>f</sub></sub>) and MLEWC = (P, E, K, D).

## 4 The FCHECK scheme

In this section, we describe the the FCHECK construction, which achieves soundness as well as security for messages that are both correlated and parameter-dependent, all in the standard model. As we noted in the introduction, prior to our work, achieving parameter-dependent correlated input security was open even in the random oracle model. We are able to exploit interactivity as a new

$\text{MAIN}(1^\lambda) \quad // \quad \text{WPRIV}^{\mathcal{S}, \mathcal{A}}(1^\lambda)$ $(\mathbf{m}_0, \mathbf{m}_1, z) \leftarrow_{\mathcal{S}} \mathcal{S}(1^\lambda, \epsilon); b \leftarrow_{\mathcal{S}} \{0, 1\}$ For $i \in [ \mathbf{m}_b ]$ do $\mathbf{p}[i] \leftarrow_{\mathcal{P}} \mathcal{P}(1^\lambda); \mathbf{k}[i] \leftarrow_{\mathcal{K}} \mathcal{K}(1^\lambda, \mathbf{p}[i], \mathbf{m}_b[i])$ $\mathbf{c}[i] \leftarrow_{\mathcal{E}} \mathcal{E}(1^\lambda, \mathbf{k}[i], \mathbf{m}_b[i])$ $b' \leftarrow_{\mathcal{A}_2} \mathcal{A}_2(1^\lambda, \mathbf{p}, \mathbf{c}, z); \text{Ret } (b = b')$
$\text{MAIN}(1^\lambda) \quad // \quad \text{CDIPFO}_{\text{OS}}^{\mathcal{S}, \mathcal{A}}(1^\lambda)$ $(\mathbf{p}, z) \leftarrow_{\mathcal{S}} \mathcal{S}(1^\lambda); b \leftarrow_{\mathcal{S}} \{0, 1\}$ For $i \in [ \mathbf{p} ]$ do $\text{If } b = 1 \text{ then}$ $(\alpha, \beta) \leftarrow \mathbf{p}[i]; \mathbf{F}[i] \leftarrow_{\mathcal{O}\text{bf}} \mathcal{O}\text{bf}(1^\lambda, (\alpha, \beta))$ $\text{Else}$ $(\alpha', \beta') \leftarrow \mathbf{p}[i]; \alpha \leftarrow_{\mathcal{S}} \{0, 1\}^{ \alpha' }; \beta \leftarrow_{\mathcal{S}} \{0, 1\}^{ \beta' }$ $\mathbf{F}[i] \leftarrow_{\mathcal{O}\text{bf}} \mathcal{O}\text{bf}(1^\lambda, (\alpha, \beta))$ $b' \leftarrow_{\mathcal{A}} \mathcal{A}(1^\lambda, \mathbf{F}[i], z); \text{Ret } (b = b')$

**Fig. 6.** The WPRIV game on the left, and the and CDIPFO game on the right.

ingredient to design a scheme that achieves security for parameter-dependent correlated messages.

Our approach starts by going after a new, seemingly weak primitive, one we call MLE-Without-Comparison (MLEWC). As the name indicates, MLEWC schemes are similar to MLE schemes in syntax and functionality, except that they do not support comparison between ciphertexts. We show that MLEWC can be realized in the standard model, starting from point function obfuscation [17] or, alternatively, UCE-secure hash function families [10]. However, comparison is essential to enable deduplication. To enable comparison, FCHECK employs an interactive protocol using a fully homomorphic encryption (FHE) scheme [19–21, 28, 36, 38, 55], transforming the MLEWC scheme into an iMLE scheme. We view FCHECK as a theoretical construction, and not an immediately practical iMLE scheme.

**MLE WITHOUT COMPARISON (MLEWC).** A scheme  $\text{MLEWC} = (\mathcal{P}, \mathcal{E}, \mathcal{K}, \mathcal{D})$  consists of four algorithms. Parameters are generated via  $p \leftarrow_{\mathcal{P}} \mathcal{P}(1^\lambda)$ . Keys are generated via  $k \leftarrow_{\mathcal{K}} \mathcal{K}(1^\lambda, p, m)$ , where  $m \in \{0, 1\}^{\mu(\lambda)}$  is the plaintext. Encryption  $\mathcal{E}$  takes  $p, k, m$  and returns a ciphertext  $c \leftarrow_{\mathcal{E}} \mathcal{E}(1^\lambda, k, m)$ . Decryption  $\mathcal{D}$  takes input  $k, c$  and returns  $m \leftarrow \mathcal{D}(1^\lambda, k, c)$ , or  $\perp$ . Correctness requires that  $\mathcal{D}(1^\lambda, k, c) = m$  for all  $k \in [\mathcal{K}(1^\lambda, p, m)]$ , for all  $c \in [\mathcal{E}(1^\lambda, k, m)]$ , for all  $p \in [\mathcal{P}(1^\lambda)]$ , for all  $m \in \{0, 1\}^{\kappa(\lambda)}$  for all  $\lambda \in \mathbb{N}$ .

The WPRIV game with MLEWC, an auxiliary source  $\mathcal{S}$  and an adversary  $\mathcal{A}$  is described in Fig. 6. The game runs  $\mathcal{S}$  to get two vectors  $\mathbf{m}_0, \mathbf{m}_1$ , and forms  $\mathbf{c}$  by encrypting one of the two vectors, using a fresh parameter for each component, or by picking random strings.  $\mathcal{A}$  should guess which the case is. We associate advantage  $\text{Adv}_{\text{MLEWC}, \mathcal{S}, \mathcal{A}}^{\text{wpriv}}(\lambda) = 2 \Pr[\text{WPRIV}_{\text{MLEWC}}^{\mathcal{A}, \mathcal{S}}(1^\lambda)] - 1$ . For MLEWC to be WPRIV-secure, advantage should be negligible for all PT adversaries  $\mathcal{A}$  for all

unpredictable PT auxiliary sources  $S$ . Note that unlike PRIV, here, a fresh parameter is picked for each encryption, and although we will end up using WPRIV-secure schemes to build parameter-dependent iMLE, in the WPRIV game, the source  $S$  is not provided the parameters.

**FULLY HOMOMORPHIC ENCRYPTION (FHE)** [36]. An FHE scheme  $FHE = (K_f, E_f, D_f, Ev_f)$  is a 4-tuple of algorithms. Key generation  $K_f(1^\lambda)$  returns  $(pk, sk)$ , encryption  $E_f(1^\lambda, \cdot)$  takes  $pk$ , plaintext  $m \in \mathbf{M}(\lambda)$ , and returns ciphertext  $c$ , and decryption returns  $m' \leftarrow D_f(1^\lambda, sk, c)$  on input  $sk$  and ciphertext  $c$ , where  $m = \perp$  indicates an error. The set of valid ciphertexts is denoted by  $\mathbf{C}(\lambda) = \{c : D_f(1^\lambda, sk, c) \neq \perp, (pk, sk) \in [K_f(1^\lambda)]\}$ . Decryption correctness requires that  $D_f(1^\lambda, sk, E_f(1^\lambda, pk, m)) = m$  for all  $(pk, sk) \in [K_f(1^\lambda)]$ , for all  $m \in \mathbf{M}(\lambda)$  for all  $\lambda \in \mathbb{N}$ .

Let  $\langle \cdot \rangle$  denote an encoding which maps boolean circuits  $f$  to strings denoted by  $\langle f \rangle$  such that there exists PT Eval which satisfies  $\text{Eval}(\langle f \rangle, x) = f(x)$  for every valid input  $x \in \{0, 1\}^n$ , where  $n$  is the input length of  $f$ . Evaluation  $Ev_f$  takes input a public key  $pk$ , a circuit encoding  $\langle f \rangle$  and a tuple of ciphertexts  $\mathbf{c}$  such that  $|\mathbf{c}|$  is the input length of  $f$  and returns  $c' \leftarrow Ev_f(1^\lambda, pk, \langle f \rangle, \mathbf{c})$ . Evaluation correctness requires that for random keys, on all functions and all inputs,  $Ev_f$  must compute the correct output when run on random coins, except with negligible error.

**THE FCHECK SCHEME.** Let  $FHE = (K_f, E_f, D_f, Ev_f)$  be an FHE scheme, and let MLEWC = (P, E, K, D) be a MLEWC scheme where K is deterministic. The FCHECK[FHE, MLEWC] iMLE scheme is described in Fig. 5. The Init algorithm is omitted: it lets  $\mathbf{U} \leftarrow \emptyset$ , and lets **fil** and **own** be empty tables.

In FCHECK, clients encrypt their plaintexts with MLEWC to be stored on the server, but pick a fresh parameter each time. The server's storage consists of a list of ciphertext-parameter pairs  $\mathbf{c}[i], \mathbf{p}[i]$ . When a client wants to put  $m$ , for each such  $\mathbf{c}[i], \mathbf{p}[i]$ , the server should generate a key  $\mathbf{k}[i] \leftarrow K(1^\lambda, \mathbf{p}[i], m)$  and check if  $D(1^\lambda, \mathbf{k}[i], \mathbf{c}[i]) = m$ .

A match means that a duplicate ciphertext already exists on the server, while no match means that  $m$  is a fresh plaintext. The search for a match should be carried without the server learning  $m$  and is hence done over FHE ciphertexts of the components. The client sends  $pk$  and  $c_f \leftarrow E_f(1^\lambda, pk, m)$  and the server encrypts each  $\mathbf{c}[i], \mathbf{p}[i]$  to get  $c_c$  and  $c_p$  and runs  $Ev_f$  on the **cmp** circuit described below with these values.

**cmp** $(m, p, c, r, n, i)$

If  $D(1^\lambda, K(1^\lambda, p, m), c) = m$  and  $r = 0^{\kappa(\lambda)}$  then return  $p, i + 1, i + 1$

Else return  $r, n, i + 1$

The client is provided the encryptions of  $r$  and  $n$  in the end. If  $n = 0$ , no match was found, and the client picks  $p \leftarrow P(1^\lambda)$ , computes  $c \leftarrow E(1^\lambda, K(1^\lambda, p, m), m)$ , and sends  $p, c$  to be stored on the server. Otherwise,  $n$  refers to the index of the match, and serves as the tag, and  $r$  refers to the parameter in the match. Now the client computes  $k \leftarrow K(1^\lambda, r, m)$ , encrypts it under its private key, and stores the result on the server. The Reg and Get protocols proceed in a simple

manner, and are described in Fig. 5. It can be checked that FCHECK performs deduplication as expected, and we show this formally in the full version.

$E(1^\lambda, k_H, k, m)$	$D(1^\lambda, k_H, k, c_0, \dots, c_n)$
$c_0 \leftarrow \text{Obf}(1^\lambda, k, 0)$	If $\text{Eval}(1^\lambda, c_0, k) = \perp$ then ret $\perp$
For $i \in [m]$ do	For $i \in [n]$ do
$c_i \leftarrow \text{Obf}(1^\lambda, k \  \langle i, \ell \rangle \  m[i], 0)$	If $\text{Eval}(1^\lambda, c_i, k \  \langle i, \ell \rangle \  0) = 1$ then
Ret $c_0, \dots, c_{ m }$	$m_i \leftarrow 0$
	else $m_i \leftarrow 1$
	Ret $m_1 \  \dots \  m_n$

**Fig. 7.** The HtO MLEWC scheme, with a CR hash HF and a point obfuscation scheme OS. Here, parameters are generated via  $P(1^\lambda)$  which runs  $K_h(1^\lambda)$  and returns the output, while message-derived keys are generated by letting  $K(1^\lambda, k_H, m)$  return  $k \leftarrow H(1^\lambda, k_H, m)$ .

**Theorem 1.** *If MLEWC is a correct MLEWC scheme then FCHECK[MLEWC, FHE] is REC-secure.*

**Proof sketch.** Observe that whenever a client puts  $m$ , and a match is found in  $\text{Put}[2, 1]$ , the client asks for the  $p, c$  pair corresponding to the index with the match, and checks by itself that this pair is a valid ciphertext for  $m$ . This, combined with the immutability of **fil** and **own** leads to perfect recovery correctness.

**Theorem 2.** *If MLEWC is WPRIV-secure and FHE is CPA-secure, then FCHECK[MLEWC, FHE] is PDPRIV-secure.*

**Proof sketch.** We replace the  $c_2$  components with encryptions of random strings, and use the CPA security of FHE to justify this. Now, only the  $\mathbf{p}, \mathbf{c}$  pairs of the plaintexts reside on the server, and hence we can hope to show that if there exists an adversary  $A$  that can guess the challenge bit from only the  $\mathbf{p}, \mathbf{c}$  values, then such an  $A$  can be used to build another adversary  $B$  which breaks WPRIV security of MLEWC.

But this cannot be accomplished right away. When  $A$  asks the game to run  $\text{Put}$  with some  $\mathbf{m}_b[i]$ , then  $B$  cannot simulate  $\text{Put}[2, 1]$  which looks through  $\mathbf{p}, \mathbf{c}$  for a match for  $\mathbf{m}_b[i]$  without knowing  $\mathbf{m}_b[i]$ . The proof first gets rid of the search step in  $\text{Put}[2, 1]$  and then builds  $B$ . We argue that the search step can be avoided. The adversary  $A$ , with no knowledge of the messages that the unpredictable source  $S$  produced, would have been able to use  $\text{MSG}$  to put a ciphertext for a  $\mathbf{m}_b[i]$  only with negligible probability.

CONSTRUCTING MLEWC SCHEMES. To get an iMLE scheme via FCHECK, we still need to construct a MLEWC scheme. The lack of comparison means that MLEWC schemes should be easier to construct compared to MLE schemes, but

constructions must still overcome two technical challenges: encrypting messages with keys derived from the messages themselves, and dealing with correlated messages. We explore two approaches to overcoming these two challenges. The first utilizes a special kind of point-function obfuscation scheme, and the second uses a UCE-secure [10] hash function. This construction, which we describe in the full version, is straightforward. We start with a hash function family,  $\text{HF} = (\text{K}_h, \text{H})$ . Parameter generation picks a hash key  $k_H$ . Given  $m$ , the key is generated as  $k \leftarrow \text{H}(1^\lambda, k_H, m, 1^\lambda)$ , and ciphertext as  $c \leftarrow \text{H}(1^\lambda, k_H, k, 1^{|m|}) \oplus m$ . Decryption, on input  $k, c$  removes the mask to recover  $m$ .

We now elaborate on the first approach, which builds a MLEWC scheme from a composable distributional indistinguishable point-function obfuscation scheme (CDIPFO) [17]. To give a high level idea for why CDIPFOs are useful, we note that point-function obfuscation is connected to encryption secure when keys and messages are related [25]. Moreover, CDIPFOs, due to their composability, remain secure even when obfuscations of several correlated points are provided and thus enable overcoming the two challenges described above.

Let  $\alpha, \beta \in \{0, 1\}^*$ . We let  $\phi_{\alpha, \beta} : \{0, 1\}^* \rightarrow \{\beta, \perp\}$  denote the function that on input  $\gamma \in \{0, 1\}^*$  returns  $\beta$  if  $\gamma = \alpha$ , and  $\perp$  otherwise. We call  $\alpha$  the special input, and  $\beta$  the special output. A point function obfuscator  $\text{OS} = (\text{Obf}, \text{Eval})$  is a pair of algorithms. Obfuscation takes  $(\alpha, \beta)$  and outputs  $\text{F} \leftarrow^s \text{Obf}(1^\lambda, (\alpha, \beta))$ , while  $\text{Eval}$  takes  $\text{F}$ , and a point  $\gamma$  and returns  $y \leftarrow^s \text{Eval}(1^\lambda, \text{F}, \gamma)$ . Correctness requires that  $\text{Eval}(1^\lambda, \text{Obf}(1^\lambda, \alpha, \beta), \alpha) = \beta$  for all  $\alpha, \beta \in \{0, 1\}^*$ , for all  $\lambda \in \mathbb{N}$ .

A PF source  $\text{S}$  outputs a tuple of point pairs  $\mathbf{p}$ , along with auxiliary information  $z$ . There exist  $m : \mathbb{N} \rightarrow \mathbb{N}$  and  $\ell : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$  such that  $|\mathbf{p}| = m(\lambda)$ , and  $|\mathbf{p}[i, 0]| = \ell(\lambda, 0)$  and  $|\mathbf{p}[i, 1]| = \ell(\lambda, 1)$  for all  $i \in [m(\lambda)]$ . Guessing probability  $\mathbf{GP}_\text{S}(\lambda)$  is defined as  $\max_i(\mathbf{GP}(\mathbf{p}[i, 0]|z))$  when  $(\mathbf{p}, z) \leftarrow^s \text{S}(1^\lambda)$ . We say that  $\text{S}$  is unpredictable if  $\mathbf{GP}_\text{S}(\cdot)$  is negligible.

Distributional indistinguishability for point function obfuscators is captured by the CDIPFO game (Fig. 6) associated with  $\text{OS}$ , an PF source  $\text{S}$ , and an adversary  $\text{A}$ . At a high level, the game either provides  $\text{OS}$ -obfuscations of point functions from  $\text{S}$ , or from a uniform distribution, and to win, the adversary  $\text{A}$  should guess which the case is. We associate advantage  $\text{Adv}_{\text{OS}, \text{S}, \text{A}}^{\text{cdipfo}}(\lambda) = 2 \Pr[\text{CDIPFO}_{\text{OS}}^{\text{A}, \text{S}}(1^\lambda)] - 1$  with  $\text{OS}, \text{S}$  and  $\text{A}$  and say that  $\text{OS}$  is CDIPFO-secure if advantage is negligible for all PT  $\text{A}$  for all unpredictable PT  $\text{S}$ . Bitansky and Canetti show that CDIPFOs can be built in the standard model, from the  $t$ -Strong Vector Decision Diffie Hellman assumption [17].

Let  $\text{HF} = (\text{K}_h, \text{H})$  denote a family of CR hash functions. The Hash-then-Obfuscate transform  $\text{HtO}[\text{HF}, \text{OS}] = (\text{P}, \text{E}, \text{K}, \text{D})$  associates an MLEWC scheme with  $\text{HF}$  and  $\text{OS}$  as in Fig. 7, restricting the message space to  $\ell$ -bit strings. At a high level, a key is generated by hashing the plaintext  $m$  with  $\text{HF}$ , and  $m$  is obfuscated bit-by-bit, with the hash as the special input. Decryption, given the hash, can recover  $m$  from the obfuscations. Correctness follows from the correctness of  $\text{OS}$ , and the following theorem shows WPRIV-security.

**Theorem 3.** *If  $\text{HF}$  is CR-secure, and  $\text{OS}$  is CDIPFO-secure, then  $\text{HtO}[\text{HF}, \text{OS}]$  is WPRIV-secure.*

The proof of the theorem and some remarks on HtO appear in the full version.

## 5 Incremental updates

In this section, we define iMLE with incremental updates, and provide a construction which achieves this goal. Building MLE schemes which can support incremental updates turns out to be challenging. On the one hand, it is easy to show that deterministic MLE schemes cannot support incremental updates. We elaborate on this in full version. On the other hand, randomized MLE schemes seem to need complex machinery such as NIZK proofs of consistency [1] to support incremental updates while retaining the same level of security as deterministic schemes, which makes them unfit for practical usage. We show how interaction can be exploited to achieve incremental updates in a practical manner, by building an efficient ROM iMLE scheme IRCE that supports incremental updates. We fix Hamming distance as the metric. In the full version, we define incremental updates w.r.t edit distance, and extend IRCE to work in this setting.

An interactive message-locked encryption scheme iMLE with updates supports an additional protocol **Upd** along with the usual three protocols **Reg**, **Put**, and **Get**. The **Upd** protocol updates a ciphertext of a file  $m_1$  stored on the server to a ciphertext of an updated file  $m_2$ . Here, **Upd**[1] (i.e. the client-side algorithm) takes inputs  $f$ ,  $\sigma_C$ , and two plaintexts  $m_1, m_2$ , and outputs a new identifier  $f_2 \in \{0, 1\}^*$ .

Now, the REC game (Fig. 3) which asks for correct recovery of files also imposes conditions on update, namely that if a legitimate client puts a file on the server, it should be able to get the file later along with updates made to the file. This is captured by letting the adversary pick **Upd** as the protocol in the INIT procedure. The WINCHECK procedure, which checks if the adversary has won, is now invoked at successful runs of **Upd** additionally. It infers the value of  $f$  used in the update protocol as well as the updated plaintext  $m_2$  and sets  $T[f] \leftarrow m_2$ , thus letting the adversary to win if a get at  $f$  does not return  $m_2$ .

We say that a scheme iMLE has incremental updates if the communication cost of updating a ciphertext for  $m_1$  stored on the server to a ciphertext for  $m_2$  is a linear function of  $\text{HAMM}(m_1, m_2)$  and  $\log |m_2|$ . More formally, there exists a linear function  $u : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$  such that for all client parameters  $\sigma_C$ , for all server-side states  $\sigma_S \in \{0, 1\}^*$ , for all plaintexts  $m_1, m_2 \in \{0, 1\}^*$  such that  $|m_1| = |m_2|$ , for all coins  $r_1, r_2$ , for all  $f \in \{0, 1\}^*$ , if  $(m_1, \sigma'_S) \leftarrow \text{Run}(\text{Get}, (\sigma_C, f), \sigma_S; r_1)$ , and  $(f', \sigma''_S) \leftarrow \text{Run}(\text{Upd}, (\sigma_C, m_1, m_2), \sigma_S; r_2)$ , and  $f' \neq \perp$ , then  $|\text{Msgs}(\text{Upd}, (\sigma_C, m_1, m_2), \sigma_S; r_2)| \leq \text{HAMM}(m_1, m_2)u(\log |m_1|, \lambda)$ .

PRELIMINARIES. A deterministic symmetric encryption (D-SE) scheme  $\text{SE} = (\text{E}, \text{D})$  is a pair of algorithms, where encryption returns  $c \leftarrow \text{E}(1^\lambda, k, m)$  on input plaintext  $m \in \{0, 1\}^*$  and key  $k \in \{0, 1\}^{\kappa(\lambda)}$ , and decryption returns  $m \leftarrow \text{D}(1^\lambda, k, c)$ . Correctness requires  $\text{D}(1^\lambda, k, \text{E}(1^\lambda, k, m)) = m$  for all plaintexts  $m \in \{0, 1\}^*$  for all keys  $k \in \{0, 1\}^{\kappa(\lambda)}$  for all  $\lambda \in \mathbb{N}$ . We say that  $\text{SE}$  supports incremental updates w.r.t Hamming distance if there exists an algorithm  $\text{U}$  such



<b>Init</b> ( $1^\lambda$ )	
$p \leftarrow_s \{0, 1\}^{\kappa(\lambda)}$ ; $\mathbf{U} \leftarrow \emptyset$ ; $\mathbf{fil} \leftarrow \emptyset$ ; $\mathbf{own} \leftarrow \emptyset$ ; Ret $\sigma_S = (p, \mathbf{U}, \mathbf{fil}, \mathbf{own})$	
<b>Reg</b> [1]( $\epsilon$ )	<b>Reg</b> [2]( $\sigma_S$ )
$k \leftarrow_s \{0, 1\}^{\kappa(\lambda)}$	$u \leftarrow_s \{0, 1\}^\lambda \setminus \mathbf{U}$ ; $\mathbf{U} \leftarrow \mathbf{U} \cup \{u\}$
Ret $(k, u, p)$	
<b>Get</b> [1]( $(k, u, p), t$ )	<b>Get</b> [2]( $\sigma_S$ )
If $c_1 = \perp$ then ret $\perp$	$(c_1, c_2) \leftarrow \mathbf{fil}[t]$ ; $c_3 \leftarrow \mathbf{own}[u, t]$
$k_2 \leftarrow \mathbf{D}(1^\lambda, k, c_3)$	If $c_3 = \perp$ then $(c'_1, c'_2) \leftarrow (\perp, \perp)$
Ret $\mathbf{D}(1^\lambda, k_2 \oplus c_2, c_1)$	

**Fig. 8.** The Init algorithm, and Reg and Get protocols of the IRCE iMLE scheme.

that  $\mathbf{U}(1^\lambda, \mathbf{E}(1^\lambda, k, m_1), \text{diff}(m_1, m_2)) = \mathbf{E}(1^\lambda, k, m_2)$  for all plaintexts  $m_1, m_2 \in \{0, 1\}^*$  for all keys  $k \in \{0, 1\}^{\kappa(\lambda)}$  for all  $\lambda \in \mathbb{N}$ .

Key-recovery security is defined through game  $\text{KR}_{\text{SE}}^{\text{A}}(1^\lambda)$  which lets adversary  $\text{A}$  query an oracle  $\text{ENC}$  with a plaintext  $m$  then picks  $k \leftarrow_s \{0, 1\}^{\kappa(\lambda)}$  and returns  $\mathbf{E}(1^\lambda, k, m)$ ;  $\text{A}$  wins if it can guess  $k$ .

The CPA security game  $\text{CPA}_{\text{SE}}^{\text{A}}(1^\lambda)$ , picks  $b \leftarrow_s \{0, 1\}$  and  $k \leftarrow_s \kappa(\lambda)$ , runs  $\text{A}$  with access to  $\text{ENC}$ , and responds to queries  $m$  by returning  $c \leftarrow \mathbf{E}(k, m)$  if  $b = 1$  and returning a random  $|c|$ -bit string if  $b = 0$ . To win, the adversary should guess  $b$ . We define advantages  $\text{Adv}_{\text{SE}, \text{A}}^{\text{kr}}(\lambda) = \Pr[\text{KR}_{\text{SE}}^{\text{A}}(1^\lambda)]$  and  $\text{Adv}_{\text{SE}, \text{A}}^{\text{cpa}}(\lambda) = 2 \cdot \Pr[\text{CPA}_{\text{SE}}^{\text{A}}(1^\lambda)] - 1$  and say that  $\text{SE}$  is KR-secure (resp. CPA-secure) if  $\text{Adv}_{\text{SE}, \text{A}}^{\text{kr}}(\cdot)$  (resp.  $\text{Adv}_{\text{SE}, \text{A}}^{\text{cpa}}(\cdot)$ ) is negligible for all PT  $\text{A}$ . The CTR mode of operation over a blockcipher, with a fixed IV is an example of a D-SE scheme with incremental updates, and KR and CPA security.

A hash function  $\mathbf{H}$  with  $\kappa(\lambda)$ -bit keys is a PT algorithm that takes  $p \in \{0, 1\}^{\kappa(\lambda)}$  and a plaintext  $m$  returns hash  $h \leftarrow \mathbf{H}(p, m)$ . Collision resistance is defined through game  $\text{CR}_{\mathbf{H}}^{\text{A}}(1^\lambda)$ , which picks  $p \leftarrow_s \{0, 1\}^{\kappa(\lambda)}$ , runs adversary  $\text{A}(1^\lambda, p)$  to get  $m_0, m_1$ , and returns True if  $m_0 \neq m_1$  and  $\mathbf{H}(p, m_0) = \mathbf{H}(p, m_1)$ . We say that  $\mathbf{H}$  is collision resistant if  $\text{Adv}_{\mathbf{H}, \text{A}}^{\text{cr}}(\lambda) = \Pr[\text{CR}_{\mathbf{H}}^{\text{A}}(1^\lambda)]$  is negligible for all PT  $\text{A}$ .

A table  $T$  is immutable if each entry  $T[t]$  can be assigned only one value after initialization. Immutable tables supports the Set-iff-empty, or **SiffE** operation, which takes inputs a table  $T$ , an index  $f$ , and a value  $m$ . If  $T[f] = \perp$  then  $T[f] \leftarrow m$  and  $m$  is returned; otherwise  $T[f]$  is returned.

**THE IRCE SCHEME.** Let  $\mathbf{H}$  denote a hash function with  $\kappa(\lambda)$ -bit keys and  $\kappa(\lambda)$ -bit outputs, and let  $\text{SE} = (\mathbf{E}, \mathbf{D})$  denote a D-SE scheme with  $\kappa(\lambda)$ -bit keys, where ciphertexts have same lengths as plaintexts and incremental updates are supported through an algorithm  $\mathbf{U}$ . The IMLE scheme  $\text{IRCE}[\text{SE}, \mathbf{H}]$  is described in figures 8 and 9. We call the construction IRCE, expanding to interactive ran-

domized convergent encryption. since it resembles the randomized convergent encryption (RCE) scheme of [12].

To describe how IRCE works, let us consider a IMLE scheme built around RCE. In RCE, to put  $m$  on the server, the client encrypts  $m$  with a random key  $\ell$  to get  $c_1$ , and then encrypts  $\ell$  with  $k_m = H(p, m)$  to get  $c_2$ , where  $p$  is a system-wide public parameter. Then,  $k_m$  is hashed once more to get the tag  $t = H(p, k_m)$ . The client sends  $t, c_1, c_2$  and the server stores  $c_1, c_2$  in a table **fil** at index  $t$ . If another client starts with  $m$ , it will end up with the same  $t$ , although it will derive a different  $c'_1, c'_2$ , as  $\ell$  is picked at random. However, when this client sends  $t, c'_1, c'_2$ , the server knows that **fil**[ $t$ ] is not empty, meaning a duplicate exists, and hence will drop  $c'_1, c'_2$ , thereby achieving deduplication. The second client should be able to recover  $m$  by sending  $t$  to the server, receiving  $c_1, c_2$ , recovering  $\ell$  from  $c_2$  and decrypting  $c_1$ . However, the problem with RCE is that, when the first client sends  $t, c_1, c_2$ , the server has no way of checking whether  $c_1, c_2$  is a proper ciphertext of  $m$ , or a corrupted one. Thus, the second client, in spite of storing a ciphertext of  $m$  on the server might not be able to recover  $m$  — this violates our soundness requirement. We will now fix this issue with interaction.

The Put protocol in IRCE differs in that, if the server finds that **fil**[ $t$ ]  $\neq \perp$  then it responds with  $h, c'_2$ , where  $(c'_1, c'_2) \leftarrow \mathbf{fil}[t]$  and  $h \leftarrow H(p, c'_1)$ . Now, the client can check that  $H(p, E(1^\lambda, c'_2 \oplus k_m, m)) = h$  which means that whenever deduplication happens, the client can check the validity of the duplicate ciphertext, which in turn guarantees soundness. The Put protocol is specified in Fig. 9, and is a bit more involved than our sketch here. Specifically, the clients are assigned unique identifiers which are provided during Put. The message-derived key  $k_m$  is also encrypted to get  $c_3$  (under per-client keys) and stored on the server, in a separate table **own**, which enables checking that a client starting a get protocol with an identifier did put the file earlier. If the client is the first to put a ciphertext with tag  $t$ , then the server still returns  $H(p, c_1), c_2, c_3$  so that external adversaries cannot learn if deduplication occurred. We note that in Fig. 9, the **fil** and **own** tables are immutable, and this will help in arguing soundness of the scheme.

The Init algorithm (Fig. 8) sets up the **fil** and **own** tables, and additional server-side state, and picks a key  $p$  for  $H$ , which becomes the public-parameter of the system. The Reg protocol (Fig. 9) sets up a new client by creating a unique client identifier  $u$ , and providing the client  $p$ . The client also picks a secret key  $k$  without the involvement of the server. The Get protocol (Fig. 9) recovers a plaintext from the identifier, which in the case of IRCE is the tag.

IRCE supports incremental updates, as described in Fig. 9. If the client wants to update  $m$  to  $m_2$ , it does not have to resend all of  $c_1, c_2, c_3$ . Instead, it can use the same key  $\ell$  and incrementally update  $c_1$ , and compute new values for  $c_2$  and  $c_3$ , along with the new tag  $t_2$ . If the server finds that **fil**[ $t_2$ ] is not empty, the same check as in Put is performed.

It is easy to see that IRCE performs deduplication, and supports incremental updates. In the full version, we formally state and prove the deduplication and

<u>Put[1]((k, u, p), m)</u>	<u>Put[2](σ<sub>S</sub>)</u>
$\ell \leftarrow_{\$} \{0, 1\}^{\kappa(\lambda)}; c_1 \leftarrow \mathbf{E}(1^\lambda, \ell, m)$ $k_m \leftarrow \mathbf{H}(p, m); c_2 \leftarrow k_m \oplus \ell; c_3 \leftarrow \mathbf{E}(1^\lambda, k, k_m)$ $t \leftarrow \mathbf{H}(p, k_m)$	$(c_1, c_2) \leftarrow \mathbf{SiffE}(\mathbf{fil}, t, c_1, c_2)$ $h \leftarrow \mathbf{H}(p, c_1)$ $c_3 \leftarrow \mathbf{SiffE}(\mathbf{own}, (u, t), c_3)$
If $c_3 \neq c'_3$ then ret $\perp$ $\ell' \leftarrow c'_2 \oplus k_m$ $c'_1 \leftarrow \mathbf{E}(1^\lambda, \ell', m)$ $h' \leftarrow \mathbf{H}(p, c'_1)$ If $h = h'$ then ret $t$ Else ret $\perp$	$\xrightarrow{u, c_1, c_2, c_3, t}$ $\xleftarrow{h, c'_2, c'_3}$
<u>Upd[1]((k, u, p), t, m<sub>1</sub>, m<sub>2</sub>)</u>	<u>Upd[2](σ<sub>S</sub>)</u>
$k_1 \leftarrow \mathbf{H}(p, m_1); k_2 \leftarrow \mathbf{H}(p, m_2)$ $\delta \leftarrow \mathbf{diff}(m_1, m_2); t_2 \leftarrow \mathbf{H}(p, k_2)$ $c_d \leftarrow k_1 \oplus k_2; t_1 \leftarrow \mathbf{H}(p, k_1)$ $c_3 \leftarrow \mathbf{E}(1^\lambda, k, k_2)$	$\xrightarrow{u, t_1, t_2, c_3, c_d, \delta}$ If $\mathbf{own}[u, t_1] \neq \perp$ then $c_3 \leftarrow \mathbf{SiffE}(\mathbf{own}, (u, t), c_3)$ $c_1, c_2 \leftarrow \mathbf{fil}[t_1]$ $c'_1 \leftarrow \mathbf{patch}(c_1, \delta)$ $c'_2 \leftarrow c_2 \oplus c_d$ $(c'_1, c'_2) \leftarrow \mathbf{SiffE}(\mathbf{fil}, t_2, c'_1, c'_2)$ Else $(c'_1, c'_2) \leftarrow (\perp, \perp)$ $h \leftarrow \mathbf{H}(p, c'_1)$
If $c_3 \neq c'_3$ then ret $\perp$ $c'_1 \leftarrow \mathbf{E}(1^\lambda, c'_2 \oplus k_2, m_2); h' \leftarrow \mathbf{H}(p, c'_1)$ If $h = h'$ then ret $t_2$ ; Else ret $\perp$	$\xleftarrow{h, c'_2, c'_3}$

**Fig. 9.** The Put and Upd protocols of the IRCE iMLE scheme. The **fil** and **own** tables are immutable, and support the set-iff-empty operation (**SiffE**) explained in text.

incremental updates properties. We also provide a proof of the following theorem which shows that IRCE is REC-secure (which, along with deduplication, establishes soundness).

**Theorem 4.** *If  $\mathbf{H}$  is collision resistant and  $\mathbf{SE}$  is a correct D-SE scheme, then  $\mathbf{IRCE}[\mathbf{H}, \mathbf{SE}]$  is REC-secure.*

**Proof sketch.** To win the REC game, the adversary  $\mathbf{A}$  must put a plaintext  $m$  on the server, possibly update it to some  $m'$ , complete a **Get** instance with the identifier for  $m$  or  $m'$  and show that the result is incorrect.

The proof uses the immutability of **fil** and **own** to argue that the ciphertext stored in the server could not have changed between the failed **Get** instance and the last time the plaintext was put/updated. However, **Put** and **Upd** both ensure that the hash of the ciphertext stored on the server matches with the hash of a

correctly formed ciphertext for the plaintext being put/updated. Consequently, whenever  $A$  breaks REC-security, it is in effect finding a pair of colliding inputs, namely the hash inputs involved in the comparison. A CR adversary  $B$  can be built which has the same advantage as the REC-advantage of  $A$ .

The following theorem (the proof of which appears in the full version) shows that IRCE is PRIV-secure in the ROM, assuming that SE is secure. Let  $\text{IRCE}_{\text{RO}}$  denote the ROM analogue of IRCE, formed by modelling  $H$  as a random oracle.

**Theorem 5.** *If SE is CPA-secure and KR-secure, then  $\text{IRCE}_{\text{RO}}[\text{SE}]$  is PRIV-secure.*

**Proof sketch.** In PRIV, the source  $S$  outputs  $\mathbf{m}_0, \mathbf{m}_1$ , the game picks  $b \leftarrow_{\$} \{0, 1\}$  and adversary  $A$  can put and update components of  $\mathbf{m}_b$ , and finally gets to learn the server-side state. To win,  $A$  should guess  $b$ .

First, the  $c_3$  components are changed to encrypt random strings instead of message-derived keys  $\mathbf{k}_m[i]$ ; CPA security of SE makes this change indistinguishable by  $A$ . The proof then moves to a game where RO responses are no longer consistent with the keys and tags being generated. For instance, if  $S$  or  $A$  queries the RO at  $p \parallel \mathbf{m}_b[i]$ , it gets a response different from  $\mathbf{k}_m[i]$ . The remainder of the proof involves two steps. First, we show that once we stop maintaining RO consistency, the adversary gets no information about the  $\ell$  values used to encrypt the messages, and hence guessing  $b$  means breaking either the CPA security or key recovery security of SE. Second, we argue that neither  $S$  nor  $A$  can detect that RO responses are inconsistent. This is because  $S$  does not know  $p$ , a prefix to the key and tag generation queries. An  $A$  that detects the inconsistency will break the CPA security of SE.

## Acknowledgments

We thank the PKC 2015 reviewers for their valuable comments. Bellare was supported in part by NSF grants CNS-1228890 and CNS-1116800. Work done while Keelveedhi was at UCSD, supported in part by NSF grants CNS-1228890 and CNS-1116800.

## References

1. M. Abadi, D. Boneh, I. Mironov, A. Raghunathan, and G. Segev. Message-locked encryption for lock-dependent messages. In R. Canetti and J. A. Garay, editors, *CRYPTO 2013, Part I*, volume 8042 of *LNCS*, pages 374–391. Springer, Aug. 2013.
2. A. Adya, W. Bolosky, M. Castro, G. Cermak, R. Chaiken, J. Douceur, J. Howell, J. Lorch, M. Theimer, and R. Wattenhofer. Farsite: Federated, available, and reliable storage for an incompletely trusted environment. *ACM SIGOPS Operating Systems Review*, 36(SI):1–14, 2002.
3. Amazon. S3. <http://aws.amazon.com/s3/pricing/>.
4. P. Anderson and L. Zhang. Fast and secure laptop backups with encrypted deduplication. In *Proc. of USENIX LISA*, 2010.

5. C. Batten, K. Barr, A. Saraf, and S. Trepetin. pStore: A secure peer-to-peer backup system. *Unpublished report, MIT Laboratory for Computer Science*, 2001.
6. M. Bellare, A. Boldyreva, and A. O’Neill. Deterministic and efficiently searchable encryption. In A. Menezes, editor, *CRYPTO 2007*, volume 4622 of *LNCS*, pages 535–552. Springer, Aug. 2007.
7. M. Bellare, R. Canetti, and H. Krawczyk. A modular approach to the design and analysis of authentication and key exchange protocols (extended abstract). In *30th ACM STOC*, pages 419–428. ACM Press, May 1998.
8. M. Bellare, O. Goldreich, and S. Goldwasser. Incremental cryptography: The case of hashing and signing. In Y. Desmedt, editor, *CRYPTO’94*, volume 839 of *LNCS*, pages 216–233. Springer, Aug. 1994.
9. M. Bellare, O. Goldreich, and S. Goldwasser. Incremental cryptography and application to virus protection. In *27th ACM STOC*, pages 45–56. ACM Press, May / June 1995.
10. M. Bellare, V. T. Hoang, and S. Keelveedhi. Instantiating random oracles via uces. Cryptology ePrint Archive, Report 2013/424, 2013. Preliminary version in Crypto 2013.
11. M. Bellare and S. Keelveedhi. Interactive message-locked encryption and secure deduplication. Cryptology ePrint Archive, 2015. Preliminary version in PKC 2015.
12. M. Bellare, S. Keelveedhi, and T. Ristenpart. Message-locked encryption and secure deduplication. In T. Johansson and P. Q. Nguyen, editors, *EUROCRYPT 2013*, volume 7881 of *LNCS*, pages 296–312. Springer, May 2013.
13. M. Bellare and T. Kohno. A theoretical treatment of related-key attacks: RKA-PRPs, RKA-PRFs, and applications. In E. Biham, editor, *EUROCRYPT 2003*, volume 2656 of *LNCS*, pages 491–506. Springer, May 2003.
14. M. Bellare and D. Micciancio. A new paradigm for collision-free hashing: Incrementality at reduced cost. In W. Fumy, editor, *EUROCRYPT’97*, volume 1233 of *LNCS*, pages 163–192. Springer, May 1997.
15. M. Bellare and P. Rogaway. Entity authentication and key distribution. In D. R. Stinson, editor, *CRYPTO’93*, volume 773 of *LNCS*, pages 232–249. Springer, Aug. 1993.
16. M. Bellare and P. Rogaway. The security of triple encryption and a framework for code-based game-playing proofs. In S. Vaudenay, editor, *EUROCRYPT 2006*, volume 4004 of *LNCS*, pages 409–426. Springer, May / June 2006.
17. N. Bitansky and R. Canetti. On strong simulation and composable point obfuscation. In T. Rabin, editor, *CRYPTO 2010*, volume 6223 of *LNCS*, pages 520–537. Springer, Aug. 2010.
18. Bitcasa. Bitcasa infinite storage. <http://blog.bitcasa.com/tag/patented-de-duplication/>.
19. Z. Brakerski. Fully homomorphic encryption without modulus switching from classical GapSVP. In R. Safavi-Naini and R. Canetti, editors, *CRYPTO 2012*, volume 7417 of *LNCS*, pages 868–886. Springer, Aug. 2012.
20. Z. Brakerski and V. Vaikuntanathan. Efficient fully homomorphic encryption from (standard) LWE. In R. Ostrovsky, editor, *52nd FOCS*, pages 97–106. IEEE Computer Society Press, Oct. 2011.
21. Z. Brakerski and V. Vaikuntanathan. Fully homomorphic encryption from ring-LWE and security for key dependent messages. In P. Rogaway, editor, *CRYPTO 2011*, volume 6841 of *LNCS*, pages 505–524. Springer, Aug. 2011.
22. E. Buonanno, J. Katz, and M. Yung. Incremental unforgeable encryption. In M. Matsui, editor, *FSE 2001*, volume 2355 of *LNCS*, pages 109–124. Springer, Apr. 2001.

23. R. Canetti. Universally composable security: A new paradigm for cryptographic protocols. In *42nd FOCS*, pages 136–145. IEEE Computer Society Press, Oct. 2001.
24. R. Canetti and R. R. Dakdouk. Obfuscating point functions with multibit output. In N. P. Smart, editor, *EUROCRYPT 2008*, volume 4965 of *LNCS*, pages 489–508. Springer, Apr. 2008.
25. R. Canetti, Y. T. Kalai, M. Varia, and D. Wichs. On symmetric encryption and point obfuscation. In D. Micciancio, editor, *TCC 2010*, volume 5978 of *LNCS*, pages 52–71. Springer, Feb. 2010.
26. Ciphertite. Ciphertite data backup. <https://www.cyphertite.com/faq.php>.
27. J. Cooley, C. Taylor, and A. Peacock. ABS: the apportioned backup system. *MIT Laboratory for Computer Science*, 2004.
28. J.-S. Coron, A. Mandal, D. Naccache, and M. Tibouchi. Fully homomorphic encryption over the integers with shorter public keys. In P. Rogaway, editor, *CRYPTO 2011*, volume 6841 of *LNCS*, pages 487–504. Springer, Aug. 2011.
29. L. P. Cox, C. D. Murray, and B. D. Noble. Pastiche: making backup cheap and easy. *SIGOPS Oper. Syst. Rev.*, 36:285–298, Dec. 2002.
30. J. Douceur, A. Adya, W. Bolosky, D. 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. IEEE, 2002.
31. Dropbox. Deduplication in Dropbox. <https://forums.dropbox.com/topic.php?id=36365>.
32. T. Duong and J. Rizzo. Here come the ninjas. *Unpublished manuscript*, 2011.
33. M. Dutch. Understanding data deduplication ratios. In *SNIA Data Management Forum*, 2008.
34. M. Fischlin. Incremental cryptography and memory checkers. In W. Fumy, editor, *EUROCRYPT'97*, volume 1233 of *LNCS*, pages 293–408. Springer, May 1997.
35. Flud. The Flud backup system. <http://flud.org/wiki/Architecture>.
36. C. Gentry. Fully homomorphic encryption using ideal lattices. In M. Mitzenmacher, editor, *41st ACM STOC*, pages 169–178. ACM Press, May / June 2009.
37. C. Gentry and S. Halevi. Implementing Gentry’s fully-homomorphic encryption scheme. In K. G. Paterson, editor, *EUROCRYPT 2011*, volume 6632 of *LNCS*, pages 129–148. Springer, May 2011.
38. C. Gentry, S. Halevi, and N. P. Smart. Fully homomorphic encryption with polylog overhead. In D. Pointcheval and T. Johansson, editors, *EUROCRYPT 2012*, volume 7237 of *LNCS*, pages 465–482. Springer, Apr. 2012.
39. GUNet. GUNet, a framework for secure peer-to-peer networking. <https://gnunet.org/>.
40. Google. Blob store. <https://developers.google.com/appengine/docs/pricing>.
41. Google. Google Drive. <http://drive.google.com>.
42. V. Goyal, A. O’Neill, and V. Rao. Correlated-input secure hash functions. In Y. Ishai, editor, *TCC 2011*, volume 6597 of *LNCS*, pages 182–200. Springer, Mar. 2011.
43. S. Halevi, D. Harnik, B. Pinkas, and A. Shulman-Peleg. Proofs of ownership in remote storage systems. In *Proceedings of the 18th ACM conference on Computer and communications security*, pages 491–500. ACM, 2011.
44. D. Harnik, B. Pinkas, and A. Shulman-Peleg. Side channels in cloud services: Deduplication in cloud storage. *Security & Privacy, IEEE*, 8(6):40–47, 2010.

45. J. Katz and V. Vaikuntanathan. Round-optimal password-based authenticated key exchange. In Y. Ishai, editor, *TCC 2011*, volume 6597 of *LNCS*, pages 293–310. Springer, Mar. 2011.
46. M. Killijian, L. Courtès, D. Powell, et al. A survey of cooperative backup mechanisms, 2006.
47. L. Marques and C. Costa. Secure deduplication on mobile devices. In *Proceedings of the 2011 Workshop on Open Source and Design of Communication*, pages 19–26. ACM, 2011.
48. D. Meister and A. Brinkmann. Multi-level comparison of data deduplication in a backup scenario. In *Proceedings of SYSTOR 2009: The Israeli Experimental Systems Conference*, page 8. ACM, 2009.
49. Microsoft. Windows Azure. <http://www.windowsazure.com/en-us/pricing/details/storage/>.
50. I. Mironov, O. Pandey, O. Reingold, and G. Segev. Incremental deterministic public-key encryption. In D. Pointcheval and T. Johansson, editors, *EUROCRYPT 2012*, volume 7237 of *LNCS*, pages 628–644. Springer, Apr. 2012.
51. NetApp. NetApp. <http://www.netapp.com/us/products/platform-os/dedupe.aspx>.
52. A. Rahumed, H. Chen, Y. Tang, P. Lee, and J. Lui. A secure cloud backup system with assured deletion and version control. In *Parallel Processing Workshops (ICPPW), 2011 40th International Conference on*, pages 160–167. IEEE, 2011.
53. T. Ristenpart, H. Shacham, and T. Shrimpton. Careful with composition: Limitations of the indistinguishability framework. In K. G. Paterson, editor, *EUROCRYPT 2011*, volume 6632 of *LNCS*, pages 487–506. Springer, May 2011.
54. M. Storer, K. Greenan, D. Long, and E. Miller. Secure data deduplication. In *Proceedings of the 4th ACM international workshop on Storage security and survivability*, pages 1–10. ACM, 2008.
55. M. van Dijk, C. Gentry, S. Halevi, and V. Vaikuntanathan. Fully homomorphic encryption over the integers. In H. Gilbert, editor, *EUROCRYPT 2010*, volume 6110 of *LNCS*, pages 24–43. Springer, May 2010.