

A preliminary version of this paper appears in the proceedings of Public-Key Cryptography 2015. This is the full version.

Interactive message-locked encryption and secure deduplication

Mihir Bellare¹ Sriram Keelveedhi²

January 21, 2015

Abstract

This paper considers the problem of secure storage of outsourced 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.

Keywords: deduplication, cloud storage, message-locked encryption.

¹ Department of Computer Science & Engineering, University of California San Diego, 9500 Gilman Drive, La Jolla, California 92093, USA. Email: mihir@eng.ucsd.edu. URL: <http://cseweb.ucsd.edu/~mihir/>. Supported in part by NSF grants CNS-1228890 and CNS-1116800.

² Work done while author was at UCSD, supported in part by NSF grants CNS-1228890 and CNS-1116800. Email: sriramkr@cs.ucsd.edu. URL: <http://www.cs.ucsd.edu/users/skeelvee/>.

Contents

1	Introduction	3
2	Preliminaries	6
3	Interactive message-locked encryption	7
4	The FCHECK scheme	10
5	Incremental updates	13
6	References	17
A	Interactive protocols	20
B	Deterministic MLE schemes cannot support incremental updates	20
C	Incremental updates: Proofs and extensions	21
	C.1 Proof of Theorem 5.1	21
	C.2 Proof of Theorem 5.2	23
D	The IRCE2 scheme	26
E	Parameter-dependent security: Proofs and extensions	30
	E.1 Proof of Theorem 4.2	31
F	MLEWC: Proofs and extensions	31
	F.1 Proof of Theorem 4.3	31
	F.2 PRV $\$$ -CDA secure MLE from UCE	33

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 [11], CE has long been used in research and commercial systems [2, 4, 5, 17, 25, 26, 28, 35, 39, 46, 47, 52, 54], an indication of practitioners' interest in secure deduplication.

MLE. Bellare, Keelveedhi and Ristenpart (BKR) [11] 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 [11] 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, 11] and for non-correlated but parameter-dependent messages [1]. This past work is summarized in Figure 1 and we now discuss it in a little more detail.

PRIOR SCHEMES. The definition of BKR [11], 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

Scheme(s)	Type	Messages		STD/ROM
		Correlated	Parameter-dependent	
CE, HCE1, HCE2, RCE [11]	MLE	Yes	No	ROM
XtDPKE, XtESPKE, ... [11]	MLE	Yes	No	STD
BHK [10]	MLE	Yes	No	STD
ABMRS [1]	MLE	No	Yes	RO
FCHECK	iMLE	Yes	Yes	STD

Figure 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.

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, 17, 25, 26, 28, 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 [13, 21, 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 [11] 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 Figure 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 [16] and FHE [18–20, 27, 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 [11] 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 [16] 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 [18–20, 27, 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

<p>Run($1^\lambda, P, \text{inp}$)</p> <p>$n \leftarrow 1; i \leftarrow 1; M \leftarrow \epsilon$</p> <p>$\mathbf{a}[1, 1] \leftarrow \text{inp}[1]; \mathbf{a}[2, 1] \leftarrow \text{inp}[2]$</p> <p>While $T[n] = \text{False}$</p> <p style="padding-left: 20px;">$(\mathbf{a}[n, i + 1], M, T[n]) \leftarrow_s P[n, i](1^\lambda, \mathbf{a}[n, i], M)$</p> <p style="padding-left: 20px;">If $n = 2$ then $n \leftarrow 1; i \leftarrow i + 1$ Else $n \leftarrow 2$</p> <p>Ret $\text{last}(\mathbf{a}[1]), \text{last}(\mathbf{a}[2])$</p>	<p>Msgs($1^\lambda, P, \text{inp}, r$)</p> <p>$n \leftarrow 1; i \leftarrow 1; j \leftarrow 1; \mathbf{a}[1, 1] \leftarrow \text{inp}[1]$</p> <p>$\mathbf{a}[2, 1] \leftarrow \text{inp}[2]; M \leftarrow \epsilon$</p> <p>While $T[n] = \text{False}$</p> <p style="padding-left: 20px;">$(\mathbf{a}[n, i + 1], M, T[n]) \leftarrow_s P[n, i](1^\lambda, \mathbf{a}[n, i], M; r[n, i])$</p> <p style="padding-left: 20px;">If $n = 2$ then $n \leftarrow 1; i \leftarrow i + 1$ Else $n \leftarrow 2$</p> <p style="padding-left: 20px;">$M[j] \leftarrow M; j \leftarrow j + 1$</p> <p>Ret M</p>
--	--

Figure 2: **Left:** Running a two player protocol P . **Right:** The Msgs procedure returns the messages exchanged during the protocol when invoked with specified inputs and coins.

MLE schemes turns out to be challenging: deterministic MLE schemes cannot support incremental updates, as we show in Appendix B, 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 [11], 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.

2 Preliminaries

We let $\lambda \in \mathbb{N}$ and 1^λ denote the security parameter and its unary representation. The empty string is denoted by ϵ . 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 $\text{diff}_{\text{HAMM}}(s_1, s_2) = \{i : s_1[i] \neq s_2[i]\}$ and $\text{patch}_{\text{HAMM}}(s_1, \delta)$ be the string s such that $s[i] = s_1[i]$ if $i \notin \delta$ and $s[i] = \neg s_1[i]$ if $i \in \delta$.

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 [15] along with extensions of [53] and [11] 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^λ , 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 (Figure 2) captures the execution of P , and `Msgs` (Figure 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 `Init` sets up server-side state: $\sigma_S \leftarrow \text{Init}(1^\lambda)$. Each protocol P consists of two players - a client $P[1]$ (meaning that the client always initiates), and a server $P[2]$. All server-side algorithms $P[2, \cdot]$ take server-side state σ_S as input, and produce an updated state σ'_S as output. The `Reg` protocol registers new users; here, `Reg[1]` takes no input and returns client parameters $\sigma_C \in \{0, 1\}^*$. The `Put` protocol stores files on the server; here, `Put[1]` takes plaintext $m \in \{0, 1\}^*$ and σ_C as inputs, and outputs an identifier $f \in \{0, 1\}^*$. The `Get` protocol retrieves files from the server; here, `Get[1]` takes identifier f and σ_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.

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

Figure 3: The REC game. The STATE procedure returns σ_S, σ_C .

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) σ_C, σ'_C , for all $m \in \{0, 1\}^*$, the expected increase in size of σ''_S over σ'_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 Figure 3, played with an adversary A , which gets access to procedures $\text{REG}, \text{INIT}, \text{STEP}, \text{MSG}, \text{STATE}$. 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 inp and P , where P has to be one of Put, Get , and 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 σ_S, σ_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 $\text{Adv}_{\text{iMLE}, A}^{\text{rec}}(\lambda) = \Pr[\text{REC}_{\text{iMLE}}^{\text{A}}(1^\lambda)]$ with iMLE and A . For recovery correctness, we require that the advantage should be negligible for all $\text{PT } A$.

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 [11], 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

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

Figure 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.

unpredictability as follows.

A source \mathcal{S} is an algorithm that on input 1^λ 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 $\mathbf{GP}_{\mathcal{S}}(\lambda)$ of \mathcal{S} is defined as $\max_{i,b,d} (\mathbf{GP}(\mathbf{m}_b[i]))$ when $(\mathbf{m}_0, \mathbf{m}_1) \leftarrow_{\mathcal{S}} \mathcal{S}(1^\lambda, d)$. We say that \mathcal{S} is unpredictable if $\mathbf{GP}_{\mathcal{S}}(\cdot)$ is negligible. We say that \mathcal{S} is a single source if it only outputs one tuple, but satisfies the other conditions. We say that \mathcal{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 Figure 4, associated with iMLE, a source \mathcal{S} and an adversary \mathcal{A} , captures privacy for unpredictable messages independent of the public parameters of the system. As with REC, the game starts by running $\sigma_S \leftarrow_{\mathcal{S}} \text{Init}(1^\lambda)$ to set up the server-side state. The game then runs \mathcal{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, \mathcal{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 \mathcal{A} can initialize an instance of Put with a plaintext $\mathbf{m}_b[i]$ by calling PUT(i).

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

The PDPRIV game of Figure 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 \mathcal{A} gets access to a PTXT procedure, which runs $\mathcal{S}(1^\lambda, \sigma_S)$ to get $\mathbf{m}_0, \mathbf{m}_1$. The other procedures follow PRIV. A simpler approach is to run \mathcal{S} with σ_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 \mathcal{A} decide, through PTXT, when \mathcal{S} is to be run. We associate advantage $\text{Adv}_{\text{iMLE}, \mathcal{S}, \mathcal{A}}^{\text{ldpriv}}(\lambda) = 2 \Pr[\text{PDPRIV}_{\text{iMLE}}^{\mathcal{S}, \mathcal{A}}(1^\lambda)] - 1$ with a scheme iMLE a source \mathcal{S} and an adversary \mathcal{A} . We

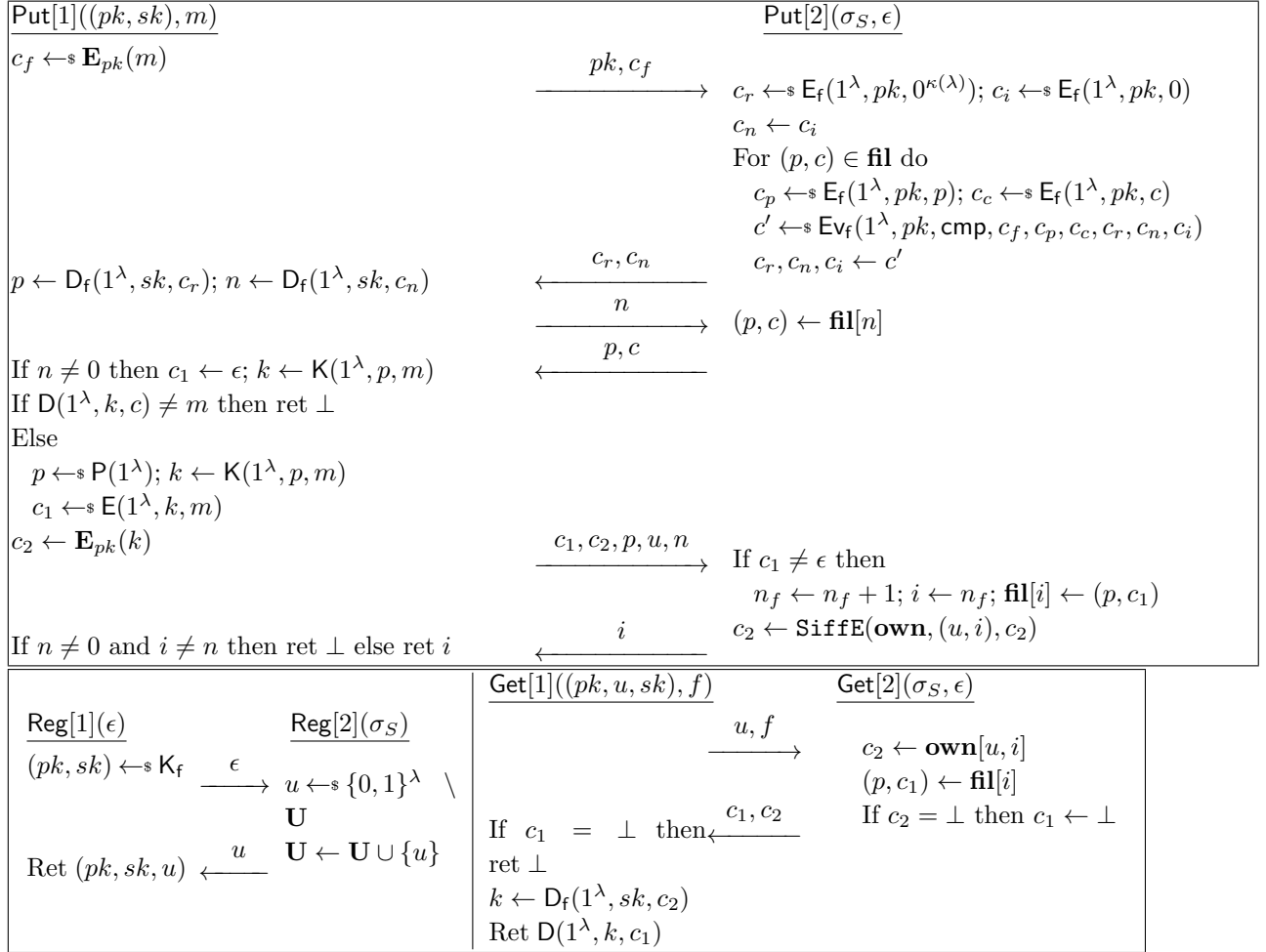


Figure 5: The FCHECK iMLE scheme over FHE = (K_f, E_f, D_f, Ev_f) and MLEWC = (P, E, K, D).

require that advantage should be negligible for all PT A for all unpredictable PT S.

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 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 [16] 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 [18–20, 27, 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 MLEWC = (P, E, K, D) consists of four algorithms. Parameters are generated via $p \leftarrow_s P(1^\lambda)$. Keys are generated via $k \leftarrow_s K(1^\lambda, p, m)$, where

$m \in \{0, 1\}^{\mu(\lambda)}$ is the plaintext. Encryption E takes p, k, m and returns a ciphertext $c \leftarrow_{\$} E(1^\lambda, k, m)$. Decryption D takes input k, c and returns $m \leftarrow D(1^\lambda, k, c)$, or \perp . Correctness requires that $D(1^\lambda, k, c) = m$ for all $k \in [K(1^\lambda, p, m)]$, for all $c \in [E(1^\lambda, k, m)]$, for all $p \in [P(1^\lambda)]$, for all $m \in \{0, 1\}^{\kappa(\lambda)}$ for all $\lambda \in \mathbb{N}$.

The WPRIV game with MLEWC, an auxiliary source S and an adversary A is described in Figure 17. The game runs 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. A should guess which the case is. We associate advantage $\text{Adv}_{\text{MLEWC}, S, A}^{\text{wpriv}}(\lambda) = 2 \Pr[\text{WPRIV}_{\text{MLEWC}}^{A, S}(1^\lambda)] - 1$. For MLEWC to be WPRIV-secure, advantage should be negligible for all PT adversaries 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 $\text{FHE} = (K_f, E_f, D_f, \text{Ev}_f)$ is a 4-tuple of algorithms. Key generation returns $(pk, sk) \leftarrow_{\$} K_f(1^\lambda)$, encryption takes pk , plaintext $m \in \mathbf{M}(\lambda)$, and returns ciphertext $c \leftarrow_{\$} E_f(1^\lambda, pk, m)$, 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 $c' \leftarrow_{\$} \text{Ev}_f(1^\lambda, pk, \langle f \rangle, \mathbf{c})$. 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. More precisely, for all boolean circuits f , when $(pk, sk) \leftarrow_{\$} K_f(1^\lambda)$ and $c' \leftarrow_{\$} [\text{Ev}_f(1^\lambda, pk, \langle f \rangle, \mathbf{c})]$, if $|\mathbf{c}|$ is the input length of f , then it holds that the probability that $\text{Eval}(f, \mathbf{x}) \neq y$ where $y \leftarrow D_f(1^\lambda, sk, c')$ and $\mathbf{x}[i] \leftarrow D_f(1^\lambda, sk, \mathbf{c}[i])$ is negligible for all $\mathbf{c}[1], \dots, \mathbf{c}[|\mathbf{c}|] \in \mathbf{C}(\lambda)^{|\mathbf{c}|}$.

THE FCHECK SCHEME. Let $\text{FHE} = (K_f, E_f, D_f, \text{Ev}_f)$ be an FHE scheme, and let $\text{MLEWC} = (P, E, K, D)$ be a MLEWC scheme where K is deterministic. The FCHECK[FHE, MLEWC] iMLE scheme is described in Figure 5. The Init algorithm is omitted: it lets $\mathbf{U} \leftarrow \emptyset$, and lets \mathbf{fil} and \mathbf{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 \mathbf{cmp} circuit described below with these values.

$\underline{\mathbf{cmp}(m, p, c, r, n, i)}$

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 Figure 5. It can be checked that FCHECK performs deduplication as expected, and we show this formally in Proposition E.1 of Appendix E.

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

Figure 6: 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 4.1. *If MLEWC is a correct MLEWC scheme then FCHECK[MLEWC, FHE] is REC-secure.*

Proof sketch. Observe that 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 4.2. *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 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 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 relegate to Appendix F, is straightforward. We start with a hash function family, $\text{HF} = (K_h, H)$. Parameter generation picks a hash key k_H . Given m , the key is generated as $k \leftarrow H(1^\lambda, k_H, m, 1^\lambda)$, and ciphertext as $c \leftarrow H(1^\lambda, k_H, k, 1^{|m|}) \oplus m$. 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) [16]. 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 [24]. 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 β if $\gamma = \alpha$, and \perp otherwise. We call α the special input, and β the special output. A point

function obfuscator $\text{OS} = (\text{Obf}, \text{Eval})$ is a pair of algorithms. Obfuscation takes (α, β) and outputs $F \leftarrow_s \text{Obf}(1^\lambda, (\alpha, \beta))$, while Eval takes F , and a point γ and returns $y \leftarrow_s \text{Eval}(1^\lambda, F, \gamma)$. 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 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 $\mathbf{GP}_S(\lambda)$ is defined as $\max_i(\mathbf{GP}(\mathbf{p}[i, 0]|z))$ when $(\mathbf{p}, z) \leftarrow_s S(1^\lambda)$. We say that S is unpredictable if $\mathbf{GP}_S(\cdot)$ is negligible.

Distributional indistinguishability for point function obfuscators is captured by the CDIPFO game (Figure 17) associated with OS , an PF source S , and an adversary A . At a high level, the game either provides OS -obfuscations of point functions from S , or from a uniform distribution, and to win, the adversary A should guess which the case is. We associate advantage $\text{Adv}_{\text{OS}, S, A}^{\text{cdipfo}}(\lambda) = 2 \Pr[\text{CDIPFO}_{\text{OS}}^{A, S}(1^\lambda)] - 1$ with OS, S and A and say that OS is CDIPFO-secure if advantage is negligible for all PT A for all unpredictable PT S . Bitansky and Canetti show that CDIPFOs can be built in the standard model, from the t -Strong Vector Decision Diffie Hellman assumption [16].

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 HF and OS as in Figure 6, restricting the message space to ℓ -bit strings. At a high level, a key is generated by hashing the plaintext m with 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 OS , and the following theorem shows WPRIV-security.

Theorem 4.3. *If HF is CR-secure, and OS is CDIPFO-secure, then $\text{HtO}[\text{HF}, \text{OS}]$ is WPRIV-secure.*

The proof of the theorem and some remarks on HtO are provided in Appendix F.

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 Appendix B. . 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 Appendix C, 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, $\text{Upd}[1]$ (i.e. the client-side algorithm) takes inputs f , σ_C , and two plaintexts m_1, m_2 , and outputs a new identifier $f_2 \in \{0, 1\}^*$.

Now, the REC game (Figure 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 $\text{HAMM}(m_1, m_2)$

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

Figure 7: The Init algorithm, and Reg and Get protocols of the IRCE iMLE scheme.

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 σ_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 SE supports incremental updates w.r.t Hamming distance if there exists an algorithm \mathbf{U} such that $\mathbf{U}(1^\lambda, \text{E}(1^\lambda, k, m_1), \text{diff}(m_1, m_2)) = \text{E}(1^\lambda, k, m_2)$ 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 A query an oracle ENC with a plaintext m then picks $k \leftarrow_s \{0, 1\}^{\kappa(\lambda)}$ and returns $\text{E}(1^\lambda, k, m)$; 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 A with access to ENC , and responds to queries m by returning $c \leftarrow \text{E}(k, m)$ 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 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 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 H with $\kappa(\lambda)$ -bit keys is a PT algorithm that takes $p \in \{0, 1\}^{\kappa(\lambda)}$ and a plaintext m returns hash $h \leftarrow \text{H}(p, m)$. Collision resistance is defined through game $\text{CR}_{\text{H}}^{\text{A}}(1^\lambda)$, 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 $\text{H}(p, m_0) = \text{H}(p, m_1)$. We say that H is collision resistant if $\text{Adv}_{\text{H}, \text{A}}^{\text{cr}}(\lambda) = \Pr[\text{CR}_{\text{H}}^{\text{A}}(1^\lambda)]$ is negligible for all PT 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 H denote a hash function with $\kappa(\lambda)$ -bit keys and $\kappa(\lambda)$ -bit outputs, and let $\text{SE} = (\text{E}, \text{D})$ 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 $\text{IRCE}[\text{SE}, \text{H}]$ is described in figures 7 and 8. We call the construction IRCE, expanding to interactive randomized convergent encryption. since it resembles the randomized convergent encryption (RCE) scheme of [11].

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 ℓ to get c_1 , and then encrypts ℓ with $k_m = \text{H}(p, m)$ to get c_2 , where p is a system-wide public parameter. Then, k_m is hashed once more to get the tag $t = \text{H}(p, k_m)$. The client sends t, c_1, c_2 and the server stores c_1, c_2 in a table \mathbf{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 ℓ is picked at random. However, when this client sends t, c'_1, c'_2 , the server knows that $\mathbf{fil}[t]$ 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 ℓ 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 $\mathbf{fil}[t] \neq \perp$ then it responds with h, c'_2 , where $(c'_1, c'_2) \leftarrow \mathbf{fil}[t]$ and $h \leftarrow \mathbf{H}(p, c'_1)$. Now, the client can check that $\mathbf{H}(p, \mathbf{E}(1^\lambda, c'_2 \oplus k_m, m)) = h$ 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 Figure 8, 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 $\mathbf{H}(p, c_1), c_2, c_3$ so that external adversaries cannot learn if deduplication occurred. We note that in Figure 8, the **fil** and **own** tables are immutable, and this will help in arguing soundness of the scheme.

The Init algorithm (Figure 7) sets up the **fil** and **own** tables, and additional server-side state, and picks a key p for \mathbf{H} , which becomes the public-parameter of the system. The Reg protocol (Figure 8) 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 (Figure 8) recovers a plaintext from the identifier, which in the case of IRCE is the tag.

IRCE supports incremental updates, as described in Figure 8. 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 ℓ 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 $\mathbf{fil}[t_2]$ is not empty, the same check as in Put is performed.

Propositions C.1 and C.2 of Appendix C show that IRCE performs deduplication, and supports incremental updates, and their proofs proceed in a straightforward manner. The following theorem, with proof in Appendix C shows that IRCE is REC-secure (which, along with deduplication, establishes soundness).

Theorem 5.1. *If \mathbf{H} is collision resistant and SE is a correct D-SE scheme, then $\text{IRCE}[\mathbf{H}, \text{SE}]$ 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 \mathbf{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 \mathbf{B} can be built which has the same advantage as the REC-advantage of \mathbf{A} .

The following theorem (with proof in Appendix C) shows that IRCE is PRIV-secure in the ROM, assuming that SE is secure. Let IRCE_{RO} denote the ROM analogue of IRCE, formed by modelling \mathbf{H} as a random oracle.

Theorem 5.2. *If SE is CPA-secure and KR-secure, then $\text{IRCE}_{\text{RO}}[\text{SE}]$ is PRIV-secure.*

Proof sketch. In PRIV, the source \mathbf{S} outputs $\mathbf{m}_0, \mathbf{m}_1$, the game picks $b \leftarrow_{\$} \{0, 1\}$ and adversary \mathbf{A} can

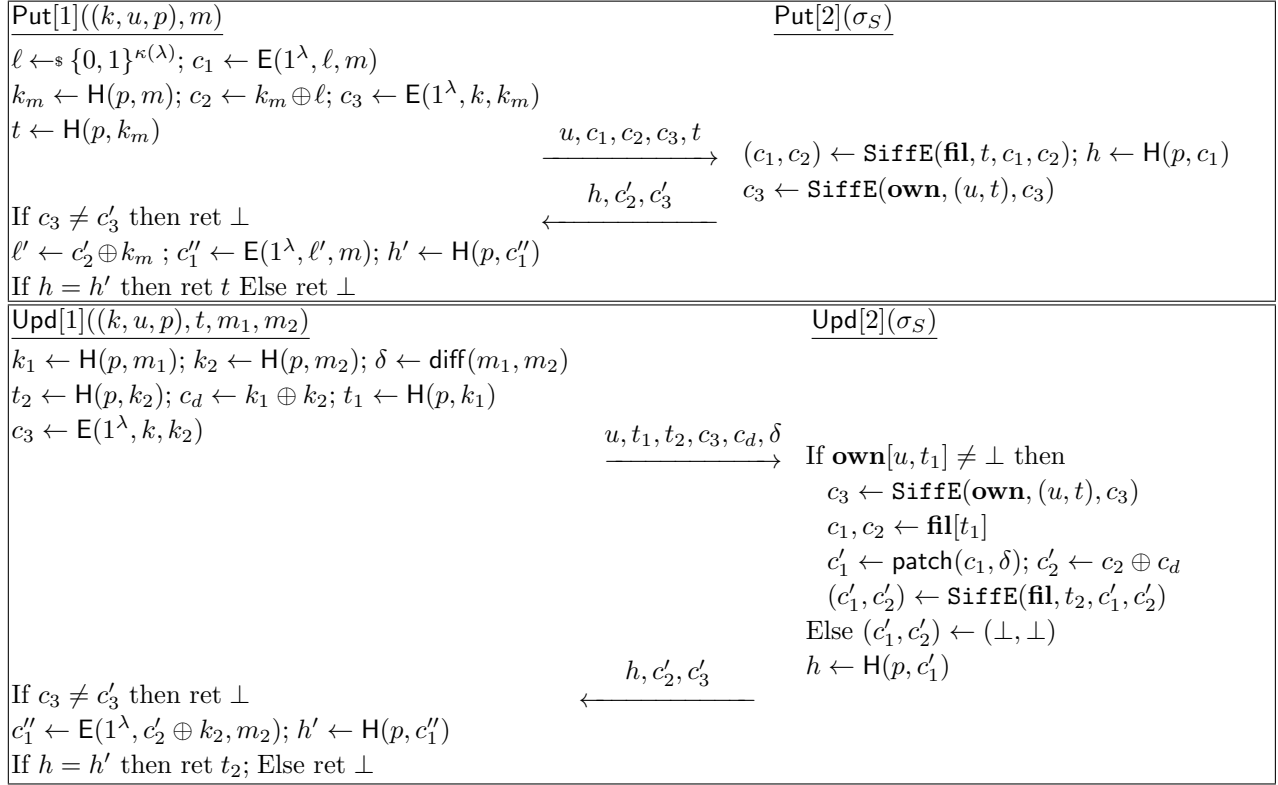


Figure 8: The Put and Upd protocols of the IRCE iMLE scheme. The **fil** and **own** tables are immutable, and support the set-iff-empty operation (**SiffE**) explained in text.

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 ℓ 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**.

6 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. (Cited on page 3, 4, 5, 6, 9, 13.)
- [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. (Cited on page 3, 4.)
- [3] Amazon. S3. <http://aws.amazon.com/s3/pricing/>. (Cited on page 4.)
- [4] P. Anderson and L. Zhang. Fast and secure laptop backups with encrypted de-duplication. In *Proc. of USENIX LISA*, 2010. (Cited on page 3, 4.)
- [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. (Cited on page 3, 4.)
- [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. (Cited on page 3.)
- [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. (Cited on page 7.)
- [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. (Cited on page 4, 13.)
- [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. (Cited on page 4, 13.)
- [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*. (Cited on page 3, 4, 5, 10, 12, 34.)
- [11] 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. (Cited on page 3, 4, 5, 6, 8, 14, 20.)
- [12] 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. (Cited on page 10.)
- [13] 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. (Cited on page 4, 13.)
- [14] 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. (Cited on page 7.)
- [15] 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. (Cited on page 6, 31.)
- [16] 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. (Cited on page 5, 10, 12, 13, 33.)
- [17] Bitcasa. Bitcasa infinite storage. <http://blog.bitcasa.com/tag/patented-de-duplication/>. (Cited on page 3, 4.)
- [18] 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. (Cited on page 5, 10.)

- [19] 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. (Cited on page 5, 10.)
- [20] 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. (Cited on page 5, 10.)
- [21] 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. (Cited on page 4.)
- [22] R. Canetti. Universally composable security: A new paradigm for cryptographic protocols. In *42nd FOCS*, pages 136–145. IEEE Computer Society Press, Oct. 2001. (Cited on page 7.)
- [23] 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. (Cited on page 10.)
- [24] 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. (Cited on page 10, 12.)
- [25] Ciphertite. Ciphertite data backup. <https://www.cypfertite.com/faq.php>. (Cited on page 3, 4.)
- [26] J. Cooley, C. Taylor, and A. Peacock. ABS: the apportioned backup system. *MIT Laboratory for Computer Science*, 2004. (Cited on page 3, 4.)
- [27] 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. (Cited on page 5, 10.)
- [28] 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. (Cited on page 3, 4.)
- [29] Y. Dodis, T. Ristenpart, and S. P. Vadhan. Randomness condensers for efficiently samplable, seed-dependent sources. In R. Cramer, editor, *TCC 2012*, volume 7194 of *LNCS*, pages 618–635. Springer, Mar. 2012. (Cited on page 33.)
- [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. (Cited on page 3.)
- [31] Dropbox. Deduplication in Dropbox. <https://forums.dropbox.com/topic.php?id=36365>. (Cited on page 3.)
- [32] T. Duong and J. Rizzo. Here come the ninjas. *Unpublished manuscript*, 2011. (Cited on page 4.)
- [33] M. Dutch. Understanding data deduplication ratios. In *SNIA Data Management Forum*, 2008. (Cited on page 7.)
- [34] M. Fischlin. Incremental cryptography and memory checkers. In W. Fumy, editor, *EUROCRYPT'97*, volume 1233 of *LNCS*, pages 293–408. Springer, May 1997. (Cited on page 4.)
- [35] Flud. The Flud backup system. <http://flud.org/wiki/Architecture>. (Cited on page 3, 4.)
- [36] C. Gentry. Fully homomorphic encryption using ideal lattices. In M. Mitzenmacher, editor, *41st ACM STOC*, pages 169–178. ACM Press, May / June 2009. (Cited on page 5, 10, 11.)
- [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. (Cited on page 10.)
- [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. (Cited on page 5, 10.)
- [39] GUNet. GUNet, a framework for secure peer-to-peer networking. <https://gnunet.org/>. (Cited on page 3, 4.)

- [40] Google. Blob store. <https://developers.google.com/appengine/docs/pricing>. (Cited on page 4.)
- [41] Google. Google Drive. <http://drive.google.com>. (Cited on page 3.)
- [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. (Cited on page 10.)
- [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. (Cited on page 7.)
- [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. (Cited on page 7.)
- [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. (Cited on page 7.)
- [46] M. Killijian, L. Courtès, D. Powell, et al. A survey of cooperative backup mechanisms, 2006. (Cited on page 3, 4.)
- [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. (Cited on page 3, 4.)
- [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. (Cited on page 7.)
- [49] Microsoft. Windows Azure. <http://www.windowsazure.com/en-us/pricing/details/storage/>. (Cited on page 4.)
- [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. (Cited on page 4, 13, 21.)
- [51] NetApp. NetApp. <http://www.netapp.com/us/products/platform-os/dedupe.aspx>. (Cited on page 3.)
- [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. (Cited on page 3, 4.)
- [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. (Cited on page 6.)
- [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. (Cited on page 3, 4.)
- [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. (Cited on page 5, 10.)

<p>Run($1^\lambda, P, \text{inp}$)</p> <p>$T \leftarrow \emptyset; n \leftarrow 1; \mathbf{M} \leftarrow \epsilon$</p> <p>For $i = 1$ to n do $\mathbf{a}[i, 1] \leftarrow \text{inp}[i]; \mathbf{rd}[i] \leftarrow 1$</p> <p>While $T \neq [n]$ do</p> <p style="padding-left: 2em;">If $n \in T$ then return \perp</p> <p style="padding-left: 2em;">$i \leftarrow \mathbf{rd}[n]$</p> <p style="padding-left: 2em;">$(\mathbf{a}[n, i + 1], \mathbf{M}, \mathbf{N}, \mathbf{T}) \leftarrow_s P[n, i](1^\lambda, \mathbf{a}[n, i], \mathbf{M})$</p> <p style="padding-left: 2em;">If $\mathbf{T} = \text{True}$ then $T \leftarrow T \cup \{n\}$</p> <p style="padding-left: 2em;">$\mathbf{rd}[n] \leftarrow \mathbf{rd}[n] + 1; n \leftarrow \mathbf{N}$</p> <p>For $i = 1$ to n do $\text{outp}[i] \leftarrow \text{last}(\mathbf{a}[i])$</p> <p>Ret outp</p>	<p>Msgs($1^\lambda, P, \mathbf{a}, r$)</p> <p>$T \leftarrow \emptyset; n \leftarrow 1; \mathbf{M} \leftarrow \epsilon; j \leftarrow 1$</p> <p>For $i = 1$ to n do $\mathbf{a}[i, 1] \leftarrow \text{inp}[i]; \mathbf{rd}[i] \leftarrow 1$</p> <p>While $T \neq [n]$ do</p> <p style="padding-left: 2em;">If $n \in T$ then return \perp</p> <p style="padding-left: 2em;">$i \leftarrow \mathbf{rd}[n]$</p> <p style="padding-left: 2em;">$(\mathbf{a}[n, i + 1], \mathbf{M}, \mathbf{N}, \mathbf{T}) \leftarrow_s P[n, i](1^\lambda, \mathbf{a}[n, i], \mathbf{M}; r[n, i])$</p> <p style="padding-left: 2em;">If $\mathbf{T} = \text{True}$ then $T \leftarrow T \cup \{n\}$</p> <p style="padding-left: 2em;">$\mathbf{rd}[n] \leftarrow \mathbf{rd}[n] + 1; n \leftarrow \mathbf{N}; \mathbf{M}[j] \leftarrow \mathbf{M}; j \leftarrow j + 1$</p> <p>Ret M</p>
--	--

Figure 9: **Left:** Running a protocol P . **Right:** The **Msgs** procedure returns the messages exchanged during the protocol when invoked with specified inputs and coins.

A Interactive protocols

Consider a protocol P with n -players, where each player is invoked for a maximum of q -times. We represent such a protocol as a $n \times q$ -tuple $(P[i, j])_{i \in [n], j \in [q]}$ of algorithms. The algorithm $P[i, j]$ represents the action of the i -th player, when invoked for the j -th time. Each algorithm is invoked with the security parameter 1^λ , an input \mathbf{a} , and a message $\mathbf{M} \in \{0, 1\}^*$, and returns a 4-tuple consisting of an output \mathbf{a}' , an outgoing message $\mathbf{M}' \in \{0, 1\}^*$, the index $\mathbf{N} \in \mathbb{N}$ of the next algorithm to run, and a boolean \mathbf{T} to indicate termination. When all algorithms of a protocol have terminated, the protocol is said to have terminated. We denote the n -players of the protocol by $P[1], \dots, P[n]$. The execution of a protocol P is captured by the **Run** algorithm, which takes inputs **inp** and returns **outp**, both tuples of n -elements, is described in Figure 9. We say that P is run on **inp**, or that $P[i]$ is invoked with **inp** $[i]$ for $i \in [n]$ to mean that the input to $P[i, 1]$ is set to **inp** $[i]$ for $i \in [n]$. We say that P returns **outp** or that $P[i]$ gets output **outp** $[i]$ for $i \in [n]$ to mean that **Run**($1^\lambda, P, \text{inp}$) returns **outp**. The **Msgs** procedure of Figure 9 returns the protocol messages when invoked with specified inputs and coins.

B Deterministic MLE schemes cannot support incremental updates

An updatable MLE scheme $\text{MLE} = (Pg, K, E, D, T, U_1, U_2)$ is a seven-tuple of PT algorithms. The first five algorithms work as in regular MLE schemes. The two update algorithms U_1 and U_2 work as follows. On input parameters p and two messages m_1, m_2 where m_2 is the newer version of m_1 , the U_1 algorithm outputs a string $c_u \in \{0, 1\}^*$. On input parameters p , update string c_u , and original ciphertext c , the U_2 algorithm produces updated ciphertext c' . We say that **MLE** is deterministic if K, E, U_1 and U_2 are deterministic.

We say that **MLE** is incremental if there exist functions $a : \mathbb{N} \rightarrow \mathbb{N}$ and $b : \mathbb{N} \rightarrow \mathbb{N}$ such that for for all $p \in [P(1^\lambda)]$, for all $m_1, m_2 \in \{0, 1\}^*$, it holds that $|c_u| \leq a(\lambda) \log(|m_1| + |m_2|) \delta + b(\lambda)$ for all $c_u \in [U_1(1^\lambda, p, m_1, m_2)]$, where $\delta = \text{HAMM}(m_1, m_2)$.

We now show that if a deterministic MLE scheme supports incremental updates, then it cannot even satisfy PRV-CDA security, the weakest among the security notions of [11]. Note that the attacker does not have the ability to update ciphertexts. Informally, the result is achieved using the fact that highly correlated plaintexts will produce similar ciphertexts, while independently chosen plaintexts will produce different-looking ciphertexts, and this can be used to guess the bit in the PRV-CDA game.

Theorem B.1. *Let **MLE** denote a deterministic MLE scheme which supports incremental updates with bound $u : \mathbb{N} \rightarrow \mathbb{N}$. Then **MLE** is not PRV-CDA secure.*

Proof. Consider $A = (A_1, A_2)$ where A_1 picks m_1, m_3, m_4 at random from $\{0, 1\}^{\mu(\lambda)}$ and picks m_2 such that $\text{dist}(m_1, m_2) = 1$. Now, A_1 outputs $(m_1, m_2), (m_3, m_4)$ as its output tuples. Informally, the two components of \mathbf{m}_0 are distance 1-apart and hence their ciphertexts should be close. On the other hand the two components of \mathbf{m}_1 are unlikely to be close, as they are picked independently at random, and A_2 uses this difference to infer the bit of the game. Specifically, $A_2(1^\lambda, p, \mathbf{c})$ returns 0 if $\text{HAMM}(\mathbf{c}[0], \mathbf{c}[1]) \leq u(1)$, and 1 otherwise. Clearly, if A_2 outputs 0 whenever $b = 0$ in the game. The probability that the ciphertexts of m_3 and m_4 are less than $u(1)$ units apart is negligible. Moreover, $\text{GP}_{A_1} = 2^{1-\mu(\lambda)}$, making it an unpredictable source. Thus, A is a valid PRV-CDA adversary with advantage negligibly away from 1, meaning that MLE is not PRV-CDA secure. \square

We remark that a similar result applies for deterministic PKE schemes, with regards to PRIV security. Mironov, Pandey, Reingold, and Segev [50] model incremental deterministic PKE schemes, but they restrict to attention to PRIV1 security, which does not consider correlated messages.

C Incremental updates: Proofs and extensions

Proposition C.1. *Let H denote a deterministic hash function and $\text{SE} = (\text{E}, \text{D})$ denote a deterministic symmetric encryption scheme. Then $\text{IRCE}[H, \text{SE}]$ supports deduplication.*

Proof. When a client with id u and parameters σ_C puts a plaintext m on the server, then, in **fil**, an entry c_1, c_2 , is added at index $t = H(p, H(p, m))$ where $c_1 \leftarrow \text{E}(1^\lambda, \ell, m)$, and $k_m \leftarrow H(p, m)$, and $c_2 \leftarrow k_m \oplus \ell$. Now, if another client with id u' and parameters σ'_C tries to put m , and sends across a c'_1, c'_2, t (we need H to be deterministic, to ensure that the same tag is generated both times) to the server, the server detects a duplicate at **fil**[t] and drops c'_1, c'_2 . A fresh copy of c_3 is still stored at **own**[u', t], but the size of c_3 is bounded by $\kappa(\lambda)$, and the increase in σ_S is bounded by $|u'| + |t| + |\kappa(\lambda)|$, which is independent of $|m|$. \square

Proposition C.2. *Let H denote a cryptographic hash function and $\text{SE} = (\text{E}, \text{D})$ denote a deterministic symmetric encryption scheme, which supports incremental updates w.r.t Hamming distance. Then $\text{IRCE}[H, \text{SE}]$ also supports incremental updates w.r.t Hamming distance.*

Proof. By inspecting the **Upd** protocol when invoked on plaintexts m_1, m_2 , we can check that the total length of the transmitted messages is $\lambda + 6\kappa(\lambda) + |\delta|$, where $\delta = \text{diff}(m_1, m_2)$. Letting $\kappa(\lambda) = \lambda$, and noting that δ is the list of positions where m_1 and m_2 differ. Since $\log |m_1|$ bits are needed to represent one position, the total size of δ can be bounded by $\text{HAMM}(m_1, m_2) \log |m_1|$. The total length of messages is bounded by $\text{HAMM}(m_1, m_2)(\log |m_1|) + 7\lambda$, proving the proposition. \square

C.1 Proof of Theorem 5.1

Proof. Consider games G_1 and G_2 of Figure 10. Here G_1 is essentially the **REC** game with **IRCE**, except that G_1 maintains tables \mathbf{C}_v , \mathbf{C}_s , and \mathbf{C}_r . When the adversary completes a instantiation of **Put**(m) through calls to **INIT** and **STEP**, the ciphertexts c'_1, c'_2, c_3 are stored in $\mathbf{C}_v[m]$. Note that c'_1, c'_2, c_3 is always valid encryption of m , and hence we store the tuple in \mathbf{C}_v , the set of valid ciphertexts. The tuple (c'_1, c'_2, c'_3) , all values returned by the server through **Put**[2, 1] are stored in $\mathbf{C}_s[m]$, the table of server ciphertexts. The same steps are also performed during **Upd**, except that here the index into \mathbf{C}_v and \mathbf{C}_s is the updated ciphertext. During **Get**[1, 2], the c_1, c_2, c_3 tuple returned by the server are added to $\mathbf{C}_r[m]$, the table of recovered ciphertexts, where m is the recovered plaintext value.

Note that, in G_1 , the **Put**[2, 1] and **Upd**[2, 1] algorithms return the ciphertexts c'_1 along with the hashes. This change does not affect the outcome of the game as these values are ignored by the **Put**[1, 2]

<p><u>MAIN</u>(1^λ) // $G_1^A(1^\lambda), G_2^A(1^\lambda)$</p> <p>$\sigma_S \leftarrow \text{Init}(1^\lambda); (\sigma_C, \sigma_S) \leftarrow \text{Run}(\text{Reg}, \epsilon, \sigma_S)$</p> <p>$\mathbf{A}_{\text{INIT,STEP,MSG,STATE}}(1^\lambda, \sigma_C); \text{Ret win}$</p> <p><u>Get</u>[1, 2](\mathbf{a})</p> <p>$(c_1, c_2, c_3) \leftarrow \mathbf{M}; \text{If } c_1 = \perp \text{ then ret } \perp; k_2 \leftarrow \text{D}(1^\lambda, k, c_3)$</p> <p>$m' \leftarrow \text{D}(1^\lambda, k_2 \oplus c_2, c_1); \mathbf{C}_r[m'] \leftarrow (c_1, c_2, c_3)$</p> <p>$\text{Ret } m', \epsilon, \text{True}$</p> <p><u>Upd</u>[1, 1]($\mathbf{a} = (t, m_1, m_2), \mathbf{M}$)</p> <p>$k_1 \leftarrow \text{H}(p, m_1); k_2 \leftarrow \text{H}(p, m_2); t_2 \leftarrow \text{H}(p, k_2)$</p> <p>$\delta \leftarrow \text{diff}(m_1, m_2); c_3 \leftarrow \text{E}(1^\lambda, k, k_2)$</p> <p>$c' \leftarrow (k_1 \oplus k_2, c_3, \delta); \mathbf{M} \leftarrow (u, t_1, t_2, c')$</p> <p>$\mathbf{a} \leftarrow (m_1, m_2, k_1, k_2, t_1, t_2, c_1, c_2, c_3); \text{Ret } (\mathbf{a}, \mathbf{M}, \text{False})$</p> <p><u>Upd</u>[2, 1](σ_S, \mathbf{M})</p> <p>$(u, t_1, t_2, c_d, c_3, \delta) \leftarrow \mathbf{M}; (p, \mathbf{U}, \mathbf{fil}, \mathbf{own}) \leftarrow \sigma_S$</p> <p>$\text{If } \mathbf{own}[u, t_1] = \perp \text{ then } (c'_1, c'_2) \leftarrow (\perp, \perp)$</p> <p>Else</p> <p style="padding-left: 2em;">$\text{If } \mathbf{fil}[t_2] = \perp \text{ then}$</p> <p style="padding-left: 4em;">$(c_1, c_2) \leftarrow \mathbf{fil}[t_1]; c'_1 \leftarrow \text{patch}(c_1, \delta)$</p> <p style="padding-left: 4em;">$c'_2 \leftarrow c_2 \oplus c_d; \mathbf{fil}[t_2] \leftarrow (c'_1, c'_2)$</p> <p style="padding-left: 2em;">Else $(c'_1, c'_2) \leftarrow \mathbf{fil}[t_2]$</p> <p style="padding-left: 2em;">$\text{If } \mathbf{own}[t_2, u] = \perp \text{ then } \mathbf{own}[t_2, u] = c_3$</p> <p style="padding-left: 2em;">Else $c_3 \leftarrow \mathbf{own}[t_2, u]$</p> <p>$h \leftarrow \text{H}(p, c'_1); \mathbf{M} \leftarrow (c'_1, h, c'_2, c_3)$</p> <p>$\text{Ret } (p, \mathbf{U}, \mathbf{fil}, \mathbf{own}), \mathbf{M}, \text{True}$</p> <p><u>Upd</u>[1, 2](\mathbf{a}, \mathbf{M})</p> <p>$(m_1, m_2, k_1, k_2, t_1, t_2, c_1, c_2, c_3) \leftarrow \mathbf{a}; (h, c'_1, c'_2, c'_3) \leftarrow \mathbf{M}$</p> <p>$\text{If } c_3 \neq c'_3 \text{ then } \mathbf{a} \leftarrow \perp$</p> <p>Else $c'_1 \leftarrow \text{E}(1^\lambda, c'_2 \oplus k_2, m_2); h' \leftarrow \text{H}(p, c'_1)$</p> <p style="padding-left: 2em;">$\text{If } h = h' \text{ then}$</p> <p style="padding-left: 4em;">$\mathbf{C}_v[m] \leftarrow (c'_1, c'_2, c_3); \mathbf{C}_s[m] \leftarrow (c'_1, c'_2, c'_3); \mathbf{a} \leftarrow t_2$</p> <p style="padding-left: 2em;">Else $\mathbf{a} \leftarrow \perp; \text{Ret } (\mathbf{a}, \epsilon, \text{True})$</p> <p><u>Get</u>[2, 1](t)</p> <p>$(p, \mathbf{U}, \mathbf{fil}, \mathbf{own}) \leftarrow \sigma_S; (c_1, c_2) \leftarrow \mathbf{fil}[t]; c_3 \leftarrow \mathbf{own}[t, u]$</p> <p>$\text{If } c_3 = \perp \text{ then } (c'_1, c'_2) \leftarrow (\perp, \perp)$</p> <p>$\mathbf{M} \leftarrow (c_1, c_2, c_3); \sigma_S \leftarrow (p, \mathbf{U}, \mathbf{fil}, \mathbf{own}); \text{Ret } \sigma_S, \mathbf{M}, \text{True}$</p>	<p><u>Put</u>[1, 1](m, \mathbf{M})</p> <p>$\ell \leftarrow \text{s } \{0, 1\}^{\kappa(\lambda)}; c_1 \leftarrow \text{E}(1^\lambda, \ell, m); k_m \leftarrow \text{H}(p, m)$</p> <p>$t \leftarrow \text{H}(p, k_m); c_2 \leftarrow k_m \oplus \ell; c_3 \leftarrow \text{E}(1^\lambda, k, k_m)$</p> <p>$c \leftarrow (c_1, c_2, c_3); \mathbf{M} \leftarrow (u, c, t)$</p> <p>$\mathbf{a} \leftarrow (m, \ell, c_1, c_2, c_3, t); \text{Ret } \mathbf{a}, \mathbf{M}, \text{False}$</p> <p><u>Put</u>[2, 1](σ_S, \mathbf{M})</p> <p>$(p, \mathbf{U}, \mathbf{fil}, \mathbf{own}) \leftarrow \sigma_S; (u, c, t) \leftarrow \mathbf{M}; (c_1, c_2, c_3) \leftarrow c$</p> <p>$\text{If } \mathbf{fil}[t] = \perp \text{ then } \mathbf{fil}[t] = (c_1, c_2) \text{ Else}$</p> <p style="padding-left: 2em;">$(c_1, c_2) \leftarrow \mathbf{fil}[t]$</p> <p>$h \leftarrow \text{H}(p, p c_1)$</p> <p>$\text{If } \mathbf{own}[t, u] = \perp \text{ then } \mathbf{own}[t, u] = c_3$</p> <p>Else $c_3 \leftarrow \mathbf{own}[t, u]$</p> <p>$\mathbf{M} \leftarrow (h, c_2, c_3)$</p> <p>$\sigma_S \leftarrow (p, \mathbf{U}, \mathbf{fil}, \mathbf{own})$</p> <p>$\text{Ret } \mathbf{M}, \sigma_S, \text{True}$</p> <p><u>Put</u>[1, 2](\mathbf{a}, \mathbf{M})</p> <p>$(m, \ell, c_1, c_2, c_3, t) \leftarrow \mathbf{a}; (h, c'_1, c'_2, c'_3) \leftarrow \mathbf{M}$</p> <p>$\ell' \leftarrow c'_2 \oplus k_m; c'_1 \leftarrow \text{E}(1^\lambda, \ell', m)$</p> <p>$h' \leftarrow \text{H}(p, c'_1); \mathbf{a} \leftarrow \perp$</p> <p>$\text{If } h = h' \text{ then}$</p> <p style="padding-left: 2em;">$\mathbf{C}_v[m] \leftarrow (c'_1, c'_2, c_3); \mathbf{C}_s[m] \leftarrow (c'_1, c'_2, c'_3)$</p> <p>$\mathbf{a} \leftarrow t$</p> <p>$\text{Ret } (\mathbf{a}, \epsilon, \text{True})$</p> <p><u>WINCHECK</u>($j$) // G_2</p> <p>$\text{If } \mathbf{PS}[j] = \text{Put} \text{ then}$</p> <p style="padding-left: 2em;">$(\sigma_C, m) \leftarrow \text{first}(\mathbf{a}[j, 1]); f \leftarrow \text{last}(\mathbf{a}[j, 1])$</p> <p style="padding-left: 2em;">$T[f] \leftarrow m$</p> <p>$\text{If } \mathbf{PS}[j] = \text{Upd} \text{ then}$</p> <p style="padding-left: 2em;">$(\sigma_C, f, m_1, m_2) \leftarrow \text{first}(\mathbf{a}[j, 1])$</p> <p style="padding-left: 2em;">$f \leftarrow \text{last}(\mathbf{a}[j, 1])$</p> <p style="padding-left: 2em;">$T[f] \leftarrow m_2$</p> <p>$\text{If } \mathbf{PS}[j] = \text{Get} \text{ then}$</p> <p style="padding-left: 2em;">$(\sigma_C, f) \leftarrow \text{first}(\mathbf{a}[j, 1]); m' \leftarrow \text{last}(\mathbf{a}[j, 1])$</p> <p style="padding-left: 2em;">$m \leftarrow T[f]$</p> <p style="padding-left: 2em;">$\text{If } m' \neq m \text{ then}$</p> <p style="padding-left: 4em;">$\text{If } \mathbf{C}_r[m'] = \mathbf{C}_s[m] = \mathbf{C}_v[m] \text{ then win}_1 \leftarrow \text{True}$</p> <p style="padding-left: 4em;">$\text{Else If } \mathbf{C}_r[m'] \neq \mathbf{C}_s[m] \text{ then win}_2 \leftarrow \text{True} \text{ else}$</p> <p style="padding-left: 2em;">$\text{win}_3 \leftarrow \text{True}$</p> <p>$\text{win} \leftarrow \text{win}_1 \vee \text{win}_2 \vee \text{win}_3$</p>
---	--

Figure 10: Games G_1 and G_2 of Theorem 5.1. The INIT, STEP, STATE and MSG procedures for G_1 and G_2 , and WINCHECK for G_1 are the same as REC and hence omitted. Also omitted is the trivial Get[1, 1] procedure, which simply sets its outgoing message to its input f .

and Upd[1, 2], except for updating \mathbf{C}_v , \mathbf{C}_s , and \mathbf{C}_r . However, maintaining \mathbf{C}_v , \mathbf{C}_s , and \mathbf{C}_r itself has no effect on setting win and hence on the outcome of G_1 . Thus, we have $\Pr[\text{REC}_{\text{IRCE}}^A(1^\lambda)] = \Pr[G_1^A(1^\lambda)]$.

Game G_2 differs from G_1 only on how WINCHECK works. In G_1 , as in REC, the game maintains a table T through WINCHECK. When the adversary runs a put or an update protocol to completion through Step, the games set $T[m] = f$, where m is the put/updated plaintext and f is the returned identifier. When the adversary runs Get(f) to completion, the games check if the recovered plaintext m' matches $m = T[f]$, and if not, set win to true. However, in G_2 , when such a mismatch happens, the game goes through a series of steps before setting win. If $\mathbf{C}_v[m] = \mathbf{C}_s[m] = \mathbf{C}_r[m']$, the game sets win₁. If $\mathbf{C}_s[m] \neq \mathbf{C}_r[m']$, the game sets win₂. Otherwise, win₃ is set. It can be observed that if $m \neq m'$, one of win₁, win₂, win₃ will get set, and since the winning condition in G_2 is $\bigvee_{i=1}^3 \text{win}_i$, it follows that $\Pr[G_1^A(1^\lambda)] = \Pr[G_2^A(1^\lambda)]$.

Now, we show that the probabilities of setting win₁ and win₂ are zero, leaving win₃ as the only way for A to break soundness. If $(c_1^r, c_2^r, c_3^r) = (c_1^v, c_2^v, c_3^v)$, then, m' has to be m , by the correctness of SE, since the latter is a valid ciphertext for m . Hence, if $m' \neq m$, the three ciphertext triples cannot be equal, meaning that $\Pr[G_2^A(1^\lambda) \text{ sets win}_1] = 0$.

If $\mathbf{C}_s[m] \neq \mathbf{C}_r[m']$, the game sets win₂, but note that both these triples are **fil**[f] and **own**[u, f], the difference being that $\mathbf{C}_s[m]$ was derived during Put[2, 1], or Upd[2, 1] while $\mathbf{C}_r[m']$ was derived during a run of Get[2, 1]. However, given that **fil** and **own** are both immutable arrays, it follows that the two values have to be equal, meaning that $\Pr[G_2^A(1^\lambda) \text{ sets win}_2] = 0$.

In the setting of win₃, we have $\mathbf{C}_s[m] = \mathbf{C}_r[m']$, and $\mathbf{C}_s[m] \neq \mathbf{C}_v[m]$. Let $(c_1^s, c_2^s, c_3^s) \leftarrow \mathbf{C}_s[m]$ and $(c_1^v, c_2^v, c_3^v) \leftarrow \mathbf{C}_v[m]$. We know that $c_3^v = c_3^s$ because Put[1, 2] and Upd[1, 2] check for this condition, returning error on failure. Moreover, $c_2^v = c_2^s$ by construction, which means that $c_1^v \neq c_1^s$. But we know that $H(p, c_1^v) = H(p, c_1^s)$ because Put[1, 2] and Upd[1, 2] check for this condition as well, once again returning error on failure. This of course means that a collision has been found in HF. It is straightforward to describe an adversary B such that $\text{Adv}_{\text{HF}, \text{B}}^{\text{cr}} = \Pr[G_2^A(1^\lambda) \text{ sets win}_3]$. Adding the above equations, we have $\text{Adv}_{\text{IRCE}, \text{A}}^{\text{rec}}(\lambda) = \text{Adv}_{\text{HF}, \text{B}}^{\text{cr}}(\lambda)$. \square

C.2 Proof of Theorem 5.2

Proof. Let S be an unpredictable PT source which outputs $m(\lambda)$ plaintexts with length $\ell(\lambda, \cdot)$. and A be a PT adversary. Let $n : \mathbb{N} \rightarrow \mathbb{N}$ denote a bound on the total number of messages stored by the adversary (including both Put and Upd), and let $q_S(\lambda) : \mathbb{N} \rightarrow \mathbb{N}$ denote a bound on the number of RO₁ queries made by the source.

Consider games G_1, G_2, G_3 and G_4 of Figure 11. Game G_1 is identical to the $\text{PRIV}_{\text{IRCE}}^{\text{S,A}}$ game. In G_1 , when the adversary imitates a Put protocol, and runs the first step through STEP, the game derives the ciphertext c_1, c_2, c_3 . Here, the key should be derived as $k_m \leftarrow \text{RO}(m)$, and the tag should be derived as $t \leftarrow \text{RO}(m)$. Instead, G_1 picks k_m and t as random $\kappa(\lambda)$ -bit strings. However, if $p \parallel k_m$ or $p \parallel m$ have already been queried at the RO by S or A, then G_1 ensures consistency by using the existing values. This event sets the bad flag in G_1 . When the adversary initiates an Update protocol, G_1 follows a similar procedure with k_2 and t_2 .

Moreover, G_1 ensures that subsequent queries to RO at $p \parallel m$ or $p \parallel k_m$ are replied with k_m or t respectively. Such events also set the bad flag in G_1 . Although Put[1, 1] and Upd[1, 1] are implemented differently in G_1 , this does not affect the outcome of the game. Game G_2 does away with maintaining consistency in the RO replies, but remains identical-until-bad to G_1 . We have

$$\Pr[\text{PRIV}_{\text{IRCE}}^{\text{S,A}}(1^\lambda)] = \Pr[G_1^{\text{S,A}}(1^\lambda)] \leq \Pr[G_2^{\text{S,A}}(1^\lambda)] + \Pr[G_2^{\text{S,A}}(1^\lambda) \text{ sets bad}]. \quad (1)$$

Game G_3 is identical to G_2 , and G_4 differs from G_3 in that the c_3 components of ciphertexts are encryptions of random strings, instead of encryptions of k_m and k_2 during Put[1, 1] and Upd[1, 1]. Here, parts of the Put[1, 1] procedures of G_3 (Left) and G_4 (Right) are provided for comparison. The Upd[1, 1] procedures are similarly related. There exists a CPA adversary B_1 , and a KR adversary B_2

<p>MAIN(1^λ) // $\boxed{G_1^{S,A}(1^\lambda), G_3^{S,A}(1^\lambda)}$, $G_2^{S,A}(1^\lambda), G_4^{S,A}(1^\lambda)$</p> <p>$\sigma_S \leftarrow \text{Init}(1^\lambda)$; $b \leftarrow \{0, 1\}$; $\vec{m}_0, \vec{m}_1 \leftarrow \text{S}^{\text{RO}}(1^\lambda, \epsilon)$ $b' \leftarrow \text{A}^{\text{REG,PUT,UPD,STEP,MSG,STATE,RO}}(1^\lambda)$; $\text{Ret } (b = b')$</p> <p>Put[1, 1](m, \mathbb{M})</p> <p>$\mathbf{n} \leftarrow \mathbf{n} + 1$; $k_m \leftarrow \{0, 1\}^{\kappa(\lambda)}$; $t \leftarrow \{0, 1\}^{\kappa(\lambda)}$ If $p \parallel m \in T$ then $\text{bad} \leftarrow \text{True}$; $k_m \leftarrow T[p \parallel m]$ If $p \parallel k_m \in T$ then $\text{bad} \leftarrow \text{True}$; $t \leftarrow T[p \parallel k_m]$ $\mathbf{k}[\mathbf{n}] \leftarrow k_m$; $\mathbf{m}[\mathbf{n}] \leftarrow m$; $\mathbf{t}[\mathbf{n}] \leftarrow t$ $\ell \leftarrow \{0, 1\}^{\kappa(\lambda)}$; $c_1 \leftarrow \text{E}(1^\lambda, \ell, m)$ $c_2 \leftarrow k_m \oplus \ell$; $c_3 \leftarrow \text{E}(1^\lambda, k, k_m)$; $c \leftarrow (c_1, c_2, c_3)$ $\mathbb{M} \leftarrow (u, c, t)$; $\mathbf{a} \leftarrow (m, \ell, c_1, c_2, c_3, t)$; $\text{Ret } \mathbf{a}, \mathbb{M}, 2, \text{False}$</p> <p>Put[2, 1](σ_S, \mathbb{M})</p> <p>$(p, \mathbf{U}, \mathbf{fil}, \mathbf{own}) \leftarrow \sigma_S$; $(u, c, t) \leftarrow \mathbb{M}$; $(c_1, c_2, c_3) \leftarrow c$ If $\mathbf{fil}[t] = \perp$ then $\mathbf{fil}[t] = (c_1, c_2)$ Else $(c_1, c_2) \leftarrow \mathbf{fil}[t]$ If $\mathbf{own}[t, u] = \perp$ then $\mathbf{own}[t, u] = c_3$ else $c_3 \leftarrow \mathbf{own}[t, u]$ $h \leftarrow \text{H}(p, p \parallel c_1)$; $\mathbb{M} \leftarrow (h, c_2, c_3)$; $\sigma_S \leftarrow (p, \mathbf{U}, \mathbf{fil}, \mathbf{own})$ $\text{Ret } \mathbb{M}, \sigma_S, 1, \text{True}$</p> <p>Put[1, 2](\mathbf{a}, \mathbb{M})</p> <p>$(m, \ell, c_1, c_2, c_3, t) \leftarrow \mathbf{a}$; $(h, c'_1, c'_2, c'_3) \leftarrow \mathbb{M}$ $\ell' \leftarrow c'_2 \oplus k_m$; $c'_1 \leftarrow \text{E}(1^\lambda, \ell', m)$; $h' \leftarrow \text{H}(p, c'_1)$ If $h = h'$ then $\mathbf{C}_v[m] \leftarrow (c'_1, c'_2, c_3)$; $\mathbf{C}_s[m] \leftarrow (c'_1, c'_2, c'_3)$; $\mathbf{a} \leftarrow t$ Else $\mathbf{a} \leftarrow \perp$ $\text{Ret } (\mathbf{a}, \epsilon, 2, \text{True})$</p> <p>RO($x$)</p> <p>For $i = 1$ to \mathbf{n} do If $x = \mathbf{k}[i]$ then $\text{bad} \leftarrow \text{True}$; $\text{Ret } \mathbf{t}[i]$ If $x = \mathbf{m}[i]$ then $\text{bad} \leftarrow \text{True}$; $\text{Ret } \mathbf{k}[i]$ If $x \notin T$ then $T[x] \leftarrow \{0, 1\}^{\kappa(\lambda)}$ $\text{Ret } T[x]$</p> <p>Put[1, 1](m, \mathbb{M}) // $\boxed{G_3^{S,A}(1^\lambda)}$, $G_4^{S,A}(1^\lambda)$</p> <p>$t \leftarrow \{0, 1\}^{\kappa(\lambda)}$; $\mathbf{k}[\mathbf{n}] \leftarrow k_m$; $\ell \leftarrow \{0, 1\}^{\kappa(\lambda)}$ $c_1 \leftarrow \text{E}(1^\lambda, \ell, m)$ $k_m \leftarrow \{0, 1\}^{\kappa(\lambda)}$; $c_2 \leftarrow \ell \oplus k_m$; $k' \leftarrow \{0, 1\}^{\kappa(\lambda)}$ $c_3 \leftarrow \text{E}(1^\lambda, k, k')$; $c_3 \leftarrow \text{E}(1^\lambda, k, k_m)$; $c \leftarrow (c_1, c_2, c_3)$ $\mathbb{M} \leftarrow (u, c, t)$; $\mathbf{a} \leftarrow (m, \ell, c_1, c_2, c_3, t)$; $\text{Ret } \mathbf{a}, \mathbb{M}, 2, \text{False}$</p>	<p>Upd[1, 1]($\mathbf{a} = (t, m_1, m_2), \mathbb{M}$)</p> <p>$\mathbf{n} \leftarrow \mathbf{n} + 1$; $k_2 \leftarrow \{0, 1\}^{\kappa(\lambda)}$; $t_2 \leftarrow \{0, 1\}^{\kappa(\lambda)}$ $\mathbf{k}[\mathbf{n}] \leftarrow k_2$; $\mathbf{m}[\mathbf{n}] \leftarrow m_2$ $\mathbf{t}[\mathbf{n}] \leftarrow t_2$; $\delta \leftarrow \text{diff}(m_1, m_2)$ If $p \parallel m_2 \in T$ then $\text{bad} \leftarrow \text{True}$; $k_2 \leftarrow T[p \parallel m_2]$ If $p \parallel k_2 \in T$ then $\text{bad} \leftarrow \text{True}$; $t_2 \leftarrow T[p \parallel k_2]$ $c_3 \leftarrow \text{E}(1^\lambda, k, k_2)$; $c' \leftarrow (k_1 \oplus k_2, c_3, \delta)$ $\mathbb{M} \leftarrow (u, t_1, t_2, c')$ $\mathbf{a} \leftarrow (m_1, m_2, k_1, k_2, t_1, t_2, c_1, c_2, c_3)$ $\text{Ret } (\mathbf{a}, \mathbb{M}, 2, \text{False})$</p> <p>Upd[2, 1](σ_S, \mathbb{M})</p> <p>$(u, t_1, t_2, c_d, c_3, \delta) \leftarrow \mathbb{M}$; $(p, \mathbf{U}, \mathbf{fil}, \mathbf{own}) \leftarrow \sigma_S$ If $\mathbf{own}[u, t_1] = \perp$ then $(c'_1, c'_2) \leftarrow (\perp, \perp)$ Else If $\mathbf{fil}[t_2] = \perp$ then $(c_1, c_2) \leftarrow \mathbf{fil}[t_1]$; $c'_1 \leftarrow \text{patch}(c_1, \delta)$ $c'_2 \leftarrow c_2 \oplus c_d$; $\mathbf{fil}[t_2] \leftarrow (c'_1, c'_2)$ Else $(c'_1, c'_2) \leftarrow \mathbf{fil}[t_2]$ If $\mathbf{own}[t_2, u] = \perp$ then $\mathbf{own}[t_2, u] = c_3$ Else $c_3 \leftarrow \mathbf{own}[t_2, u]$ $h \leftarrow \text{H}(p, c'_1)$; $\mathbb{M} \leftarrow (c'_1, h, c'_2, c_3)$ $\sigma_S \leftarrow (p, \mathbf{U}, \mathbf{fil}, \mathbf{own})$; $\text{Ret } \sigma_S, \mathbb{M}, 1, \text{True}$</p> <p>Upd[1, 2](\mathbf{a}, \mathbb{M})</p> <p>$(m_1, m_2, k_1, k_2, t_1, t_2, c_1, c_2, c_3) \leftarrow \mathbf{a}$ $(h, c'_1, c'_2, c'_3) \leftarrow \mathbb{M}$ If $c_3 \neq c'_3$ then $\mathbf{a} \leftarrow \perp$ Else $c'_1 \leftarrow \text{E}(1^\lambda, c'_2 \oplus k_2, m_2)$; $h' \leftarrow \text{H}(p, c'_1)$ If $h = h'$ then $\mathbf{C}_v[m] \leftarrow (c'_1, c'_2, c_3)$ $\mathbf{C}_s[m] \leftarrow (c'_1, c'_2, c'_3)$; $\mathbf{a} \leftarrow t_2$ Else $\mathbf{a} \leftarrow \perp$ $\text{Ret } (\mathbf{a}, \epsilon, 2, \text{True})$</p> <p>Upd[1, 1]($\mathbf{a} = (t, m_1, m_2), \mathbb{M}$) // $\boxed{G_3^{S,A}(1^\lambda)}$, $G_4^{S,A}(1^\lambda)$</p> <p>$k_2 \leftarrow \{0, 1\}^{\kappa(\lambda)}$; $t_2 \leftarrow \{0, 1\}^{\kappa(\lambda)}$ $\delta \leftarrow \text{diff}(m_1, m_2)$ $k' \leftarrow \{0, 1\}^{\kappa(\lambda)}$; $c_3 \leftarrow \text{E}(1^\lambda, k, k')$ $c_3 \leftarrow \text{E}(1^\lambda, k, k_m)$ $c' \leftarrow (k_1 \oplus k_2, c_3, \delta)$; $\mathbb{M} \leftarrow (u, t_1, t_2, c')$ $\mathbf{a} \leftarrow (m_1, m_2, k_1, k_2, t_1, t_2, c_1, c_2, c_3)$ $\text{Ret } (\mathbf{a}, \mathbb{M}, 2, \text{False})$</p>
--	---

Figure 11: Games G_1, G_2, G_3 and G_4 . The boxed code is part of G_1 and G_3 .

<p>MAIN(1^λ) // $\boxed{H_1^{S,A}(1^\lambda)}$, $\boxed{H_2^{S,A}(1^\lambda)}$, $H_3^{S,A}(1^\lambda)$</p> <p>$s_S \leftarrow_s \text{Init}(1^\lambda)$; $b \leftarrow_s \{0, 1\}$; $\vec{m}_0, \vec{m}_1 \leftarrow_s \mathbf{S}^{\text{RO}_1}(1^\lambda, \epsilon)$</p> <p>$b' \leftarrow_s \mathbf{A}^{\text{REG,PUT,UPD,STEP,MSG,STATE,RO}_2}(1^\lambda)$</p> <p>Ret bad</p> <p>Put[1, 1](m, \mathbf{M})</p> <p>$\mathbf{n} \leftarrow \mathbf{n} + 1$; $k_m \leftarrow_s \{0, 1\}^{\kappa(\lambda)}$; $t \leftarrow_s \{0, 1\}^{\kappa(\lambda)}$</p> <p>$\mathbf{k}[\mathbf{n}] \leftarrow k_m$; $\mathbf{m}[\mathbf{n}] \leftarrow m$; $\mathbf{t}[\mathbf{n}] \leftarrow t$</p> <p>$\boxed{\text{If } p \parallel m \in T_1 \text{ or } p \parallel k_m \in T_1 \text{ then bad} \leftarrow \text{True}}$</p> <p>$\boxed{\text{If } p \parallel m \in T_2 \text{ or } p \parallel k_m \in T_2 \text{ then bad} \leftarrow \text{True}}$</p> <p>$\ell \leftarrow_s \{0, 1\}^{\kappa(\lambda)}$; $c_1 \leftarrow \mathbf{E}(1^\lambda, \ell, m)$</p> <p>$c_2 \leftarrow k_m \oplus \ell$; $c_3 \leftarrow \mathbf{E}(1^\lambda, k, k_m)$</p> <p>$c \leftarrow (c_1, c_2, c_3)$</p> <p>$\mathbf{M} \leftarrow (u, c, t)$; $\mathbf{a} \leftarrow (m, \ell, c_1, c_2, c_3, t)$</p> <p>Ret $\mathbf{a}, \mathbf{M}, 2, \text{False}$</p>	<p>Upd[1, 1]($\mathbf{a} = (t, m_1, m_2), \mathbf{M}$)</p> <p>$\mathbf{n} \leftarrow \mathbf{n} + 1$; $k_2 \leftarrow_s \{0, 1\}^{\kappa(\lambda)}$; $t_2 \leftarrow_s \{0, 1\}^{\kappa(\lambda)}$</p> <p>$\mathbf{k}[\mathbf{n}] \leftarrow k_2$; $\mathbf{m}[\mathbf{n}] \leftarrow m_2$</p> <p>$\mathbf{t}[\mathbf{n}] \leftarrow t_2$; $\delta \leftarrow \text{diff}(m_1, m_2)$</p> <p>$\boxed{\text{If } p \parallel m_2 \in T_1 \text{ or } p \parallel k_2 \in T_1 \text{ then bad} \leftarrow \text{True}}$</p> <p>$\boxed{\text{If } p \parallel m_2 \in T_2 \text{ or } p \parallel k_2 \in T_2 \text{ then bad} \leftarrow \text{True}}$</p> <p>$c_3 \leftarrow \mathbf{E}(1^\lambda, k, k_2)$; $c' \leftarrow (k_1 \oplus k_2, c_3, \delta)$</p> <p>$\mathbf{M} \leftarrow (u, t_1, t_2, c')$</p> <p>$\mathbf{a} \leftarrow (m_1, m_2, k_1, k_2, t_1, t_2, c_1, c_2, c_3)$</p> <p>Ret $(\mathbf{a}, \mathbf{M}, 2, \text{False})$</p> <p><u>RO₁($x$)</u></p> <p>If $x \notin T$ then $T[x] \leftarrow_s \{0, 1\}^{\kappa(\lambda)}$</p> <p>$T_1[x] \leftarrow T[x]$; Ret $T[x]$</p> <p><u>RO₂(x)</u></p> <p>If $x \in \mathbf{k} \cup \mathbf{m}$ then bad $\leftarrow \text{True}$</p> <p>If $x \notin T$ then $T[x] \leftarrow_s \{0, 1\}^{\kappa(\lambda)}$</p> <p>$T_2[x] \leftarrow T[x]$; Ret $T[x]$</p>
--	--

Figure 12: Games H_1, H_2 and H_3 .

such that

$$\text{Adv}_{\text{SE}, \mathbf{B}_1}^{\text{cpa}}(\lambda) = \Pr[G_4^{S,A}(1^\lambda)] - \Pr[G_3^{S,A}(1^\lambda)] \quad , \quad \Pr[G_4^{S,A}(1^\lambda)] \leq m(\lambda) \text{Adv}_{\text{SE}, \mathbf{B}_2}^{\text{cpa}}(\lambda) + m(\lambda) \text{Adv}_{\text{SE}, \mathbf{B}_3}^{\text{kr}}(\lambda) \quad , \quad (2)$$

where m is the bound on the size of \mathbf{S} output along with the total number of procedure calls made by \mathbf{A} .

Consider games H_1, H_2 and H_3 of Figure 12, where, H_1 is G_2 , except that setting **bad** wins H_1 . Both the source \mathbf{S} and the adversary \mathbf{A} can set **bad**. Given that the source is not provided p , the probability that \mathbf{S} sets **bad** in H_1 is bounded by $q_S(\lambda)n(\lambda)/2^{\kappa(\lambda)}$. In H_2 , the source cannot set **bad**, as the RO_1 query points of \mathbf{S} are not taken into account while testing for **bad**. We have

$$\Pr[H_2^{S,A}(1^\lambda) \text{ sets bad}] - \Pr[H_1^{S,A}(1^\lambda) \text{ sets bad}] \leq \frac{q_S(\lambda)(\lambda)n(\lambda)}{2^{\kappa(\lambda)}}. \quad (3)$$

Consider an adversary \mathbf{A}' which runs \mathbf{A} , keeps track of all the RO queries of \mathbf{A} and when \mathbf{A} finishes, repeats all the queries. By running \mathbf{A}' , testing for **bad** can be localized to RO_2 and dropped in $\text{Put}[1, 1]$ and $\text{Upd}[1, 1]$. This change is implemented in H_3 . We have

$$\Pr[H_2^{S,A}(1^\lambda) \text{ sets bad}] \leq \Pr[H_3^{S,A'}(1^\lambda) \text{ sets bad}]. \quad (4)$$

In H_4 , as in G_4 , the c_3 components of the ciphertexts are replaced by encryptions of random strings. In H_5 , the ciphertexts c_1 are derived as encryptions of random strings, and not as encryptions of messages output by \mathbf{S} . There exist adversaries $\mathbf{B}_3, \mathbf{B}_4$ and \mathbf{B}_5 such that

$$|\Pr[H_5^{S,A'}(1^\lambda) \text{ sets bad}] - \Pr[H_3^{S,A'}(1^\lambda) \text{ sets bad}]| \leq \text{Adv}_{\text{SE}, \mathbf{B}_4}^{\text{cpa}}(\lambda) + n(\lambda)(\text{Adv}_{\text{SE}, \mathbf{B}_4}^{\text{cpa}}(\lambda) + \text{Adv}_{\text{SE}, \mathbf{B}_5}^{\text{kr}}(\lambda)). \quad (5)$$

Finally, in H_5 , the adversary receives no information about the messages output by \mathbf{S} , and it follows that

$$\Pr[H_5^{S,A}(1^\lambda) \text{ sets bad}] \leq \frac{q_A(\lambda)n(\lambda)}{2^{\kappa(\lambda)}} + \leq q_S(\lambda)n(\lambda)\mathbf{GP}_S(\lambda). \quad (6)$$

where $q_A(\lambda) : \mathbb{N} \rightarrow \mathbb{N}$ is a bound on the number of RO_2 queries made by \mathbf{A} . Adding the above, we

<p><u>$\text{dist}_e(s_1, s_2)$</u> For $i = 1$ to s_1 do $D[i, 0] \leftarrow i$ For $j = 1$ to s_2 do $D[0, j] \leftarrow j$ For $i = 1$ to s_1 do For $j = 1$ to s_2 do If $s_1[i] = s_2[j]$ then $D[i, j] \leftarrow D[i - 1, j - 1]$ Else $D[i, j] \leftarrow \min(D[i - 1, j] + 1,$ $D[i, j - 1] + 1, D[i - 1, j - 1] + 1)$ Ret $D[s_1 , s_2]$ <u>$\text{patch}_e(S, s_1)$</u> For $(\alpha, \beta, \gamma) \in S$ If $\alpha = r$ then $s_1[\beta] \leftarrow \gamma$ If $\alpha = i$ then $s_1 \leftarrow s_1[1, \beta] + \gamma + s_1[\beta + 1, s_1]$ If $\alpha = d$ then $s_1 \leftarrow s_1[1, \beta - 1] + s_1[\beta + 1, s_1]$ Ret s_1</p>	<p><u>$\text{diff}_e(D, s_1, s_2)$</u> $i \leftarrow s_1 ; j \leftarrow s_2 ; S \leftarrow \epsilon$ While $i > 0$ and $j > 0$ $d \leftarrow D[i - 1, j - 1]; t \leftarrow D[i - 1, j]$ $l \leftarrow D[i, j - 1]$ If $d < t$ and $d < l$ then If $d < D[i, j]$ then $S \leftarrow S \cup \{(r, i, s_2[j])\}$ $i \leftarrow i - 1; j \leftarrow j - 1$; continue If $l < t$ then and $l \leq D[i, j]$ $S \leftarrow S \cup \{(i, i, s_2[j])\}; j \leftarrow j - 1$; continue $S \leftarrow S \cup \{(d, i)\}; i \leftarrow i - 1$ Ret S</p>
---	--

Figure 13: The $\text{dist}_e, \text{diff}_e,$ and patch_e algorithms for edit distance.

have

$$\text{Adv}_{\text{IRCE}[\text{SE}], S, A}^{\text{priv}}(\lambda) \leq 2(n(\lambda) + 1)\text{Adv}_{\text{SE}, C_1}^{\text{cpa}}(\lambda) + 2n(\lambda)\text{Adv}_{\text{SE}, C_2}^{\text{kr}}(\lambda) + \frac{(q_S(\lambda) + q_A(\lambda))n(\lambda)}{2^{\kappa(\lambda)}} + q_A(\lambda)n(\lambda)\text{GP}_S(\lambda).$$

where C_1 and C_2 are the ones among the CPA and KR adversaries with highest advantage. \square

D The IRCE2 scheme

EDIT DISTANCE. The edit distance between two strings s_1 and s_2 over Σ is the minimum number of single character modifications, including insertion, deletion, and substitution that need to be performed to convert s_1 to s_2 . We define edit distance dist_e and associated algorithms diff_e and patch_e in Figure 13.

IVT. Let E be a blockcipher with blocksize $w(\lambda)$ and $\kappa(\lambda)$ -bit keys. The $\text{IVT}[E] = (E_{\text{IVT}}, D_{\text{IVT}})$ SE scheme is described in Figure 14. We assume that plaintext lengths are exact multiples of $w(\lambda)$, a restriction that can be circumvented via an appropriate padding.

At a high level, IVT is like CTR mode of operation with E , but instead of having a single starting point for the counter, it contains a table accompanying the ciphertext that says what counter value to use for each block of data. This table can be compressed down when the counter values are increasing incrementally, and encryption the scheme works just like CTR. But having this table enables inserting, deleting and changing blocks. For example, given a ciphertext of ℓ -blocks, if a block has to be inserted in the middle, it is XORed with the output of E on counter $\ell + 1$. The table is modified to indicate this aberration in the middle, but can still be compressed efficiently.

The diff_{IVT} algorithm (Figure 14), given two plaintexts m_1, m_2 computes the information S that needs to be applied to a ciphertext c_1 of m_1 under k to change it to a ciphertext of m_2 . The $\text{patch}_{\text{IVT}}$ algorithm (Figure 14) takes S and c_1 and returns c_2 , a ciphertext of m_2 . Here, c_1 does not have to be an output of E_{IVT} ; it could be a result of previous patching efforts and hence contain a modified IV table.

We state the following proposition, which is easy to verify from the pseudocode of IVT.

Proposition D.1. *For all $c = (c_1, c_2, \text{iv}_l)$, for all $k \in \{0, 1\}^{\kappa(\lambda)}$, if $D_{\text{IVT}}(1^\lambda, k, c) = m$ then for all $m_2 \in \{0, 1\}^*$, it holds that $D_{\text{IVT}}(1^\lambda, k, \text{patch}_{\text{IVT}}(1^\lambda, c, \text{diff}_{\text{IVT}}(1^\lambda, k, m_1, m_2))) = m_2$.*

<p>$\underline{E_{IVT}(1^\lambda, k, m)}$ For $i = 1$ to $m /w(\lambda)$ do $iv[i] = i$ $c_1 \leftarrow \text{CTR}[E](k, m)$; $c_2 \leftarrow \text{compress}_{IVT}(iv)$ $iv_l \leftarrow i$; Ret (c_1, c_2, iv_l)</p> <p>$\underline{\text{compress}_{IVT}(iv)}$ While $i < iv$ $s \leftarrow iv[i]$; $j \leftarrow 0$; While $iv[i+j] = s+j$; $j \leftarrow j+1$ $iv_c \leftarrow iv_c \cup \{s, j\}$ Ret iv_c</p> <p>$\underline{\text{patch}_{IVT}(1^\lambda, c, (\delta, iv_l))}$ $(c_1, c_2, iv_l) \leftarrow c$; $iv \leftarrow \text{expand}_{IVT}(c_2)$ For $(\alpha, \beta, \gamma) \in \delta$ If $\alpha = r$ then $c_1[\beta] \leftarrow c_1[\beta] \oplus \gamma$ If $\alpha = i$ then $c_1 \leftarrow c_1[1, \beta] \parallel \gamma[1] \parallel c_1[\beta+1, c_1]$ $iv \leftarrow iv[1, \beta] \parallel \gamma[2] \parallel iv[\beta+1, iv]$ If $\alpha = d$ then $c_1 \leftarrow c_1[1, \beta-1] \parallel c_1[\beta+1, c_1]$; $iv \leftarrow iv[1, \beta] \parallel iv[\beta+1, iv]$ $c_2 \leftarrow \text{compress}_{IVT}(iv)$; Ret (c_1, c_2, iv_l)</p>	<p>$\underline{D_{IVT}(1^\lambda, k, c)}$ $(c_1, c_2, iv_l) \leftarrow c$; $iv \leftarrow \text{expand}_{IVT}(c_2)$ For $i = 1$ to $c_1 /w(\lambda)$ do $m[i] \leftarrow E(1^\lambda, k, iv[i]) \oplus c_1[i]$ Ret m</p> <p>$\underline{\text{expand}_{IVT}(iv_c)}$ $k \leftarrow 1$ For (a, b) in iv_c do For $i = 1$ to b do $iv[k] \leftarrow a+i$; $k \leftarrow k+1$ Ret iv</p> <p>$\underline{\text{diff}_{IVT}(1^\lambda, k, m_1, m_2, iv_l)}$ $D \leftarrow \text{dist}_e(m_1, m_2)$; $\delta \leftarrow \text{diff}_e(D, m_1, m_2)$ $n \leftarrow m_1 /w(\lambda)$ For $(\alpha, \beta, \gamma) \in \delta$ If $\alpha = r$ then $\alpha' \leftarrow r$; $\beta' \leftarrow \beta$; $\gamma' \leftarrow m_1[\beta] \oplus m_2[\beta]$ If $\alpha = d$ then $\alpha' \leftarrow d$; $\beta' \leftarrow \beta$ If $\alpha = i$ then $\alpha' \leftarrow i$; $\beta' \leftarrow \beta$; $iv_l \leftarrow iv_l + 1$ $\gamma'[1] \leftarrow E(1^\lambda, k, iv_l) \oplus \gamma$; $\gamma'[2] \leftarrow iv_l$ $S \leftarrow S \cup \{(\alpha', \beta', \gamma')\}$ Ret S, iv_l</p>
---	--

Figure 14: **Bottom:** The IVT SE scheme with a blockcipher E of blocksize $w(\lambda)$.

We do not explicitly prove or require incremental encryption security of IVT, but we will reason about its security as part of the IRCE2 construction. However, we do observe that the SE scheme $IVT[E] = (E_{IVT}, D_{IVT})$ is deterministic CPA secure, as encryption here simply runs the CTR mode with E and adds a few other elements which can be generated knowing only the length of the ciphertext.

THE IRCE2 SCHEME. Let H denote a hash function H with $\kappa(\lambda)$ -bit keys and $\kappa(\lambda)$ -bit outputs. Let E be a blockcipher with blocksize E and $\kappa(\lambda)$ -bit keys. The $IRCE2[H, E] = (\text{Init}, \text{Reg}, \text{Put}, \text{Get}, \text{Upd})$ iMLE scheme resembles IRCE, but uses $IVT[E] = (E_{IVT}, D_{IVT})$ as the SE scheme. The Reg protocol and the Init algorithm are the same as in IRCE (Figure 7). The Put , Get and Upd protocols are described in Figure 15.

Proposition C.1 can be extended in a simple manner to argue that IRCE2 supports deduplication, as the Put protocols in the two schemes are essentially the same. The difference in the Upd protocols does not play any role.

In Section 5, we define incremental updates w.r.t Hamming distance. That definition can be modified for edit distance by simply replacing the $u(\text{HAMM}(m_1, m_2))$ bound with $u(\text{dist}_e(m_1, m_2))$. Specifically, we consider edit distance with alphabet $\Sigma = \{0, 1\}^{w(\lambda)}$.

Proposition D.2. *Then $IRCE2[H, E]$ supports incremental updates w.r.t edit distance.*

Proof. It is easy to observe that the IVT SE scheme supports incremental updates w.r.t edit distance: the diff_{IVT} algorithm (Figure 14), given two plaintexts m_1, m_2 and a key k produces S such that $|S| \leq 2|\delta|$, where δ in turn consists of $\text{dist}_e(m_1, m_2)$ elements, each element no more in size than $\log(|m_1| + |m_2|)$. The patch_{IVT} algorithm, given such an S , can convert a ciphertext for m_1 (under k) to a ciphertext for m_2 . Now, by inspecting the Upd protocol, it can be checked that the total

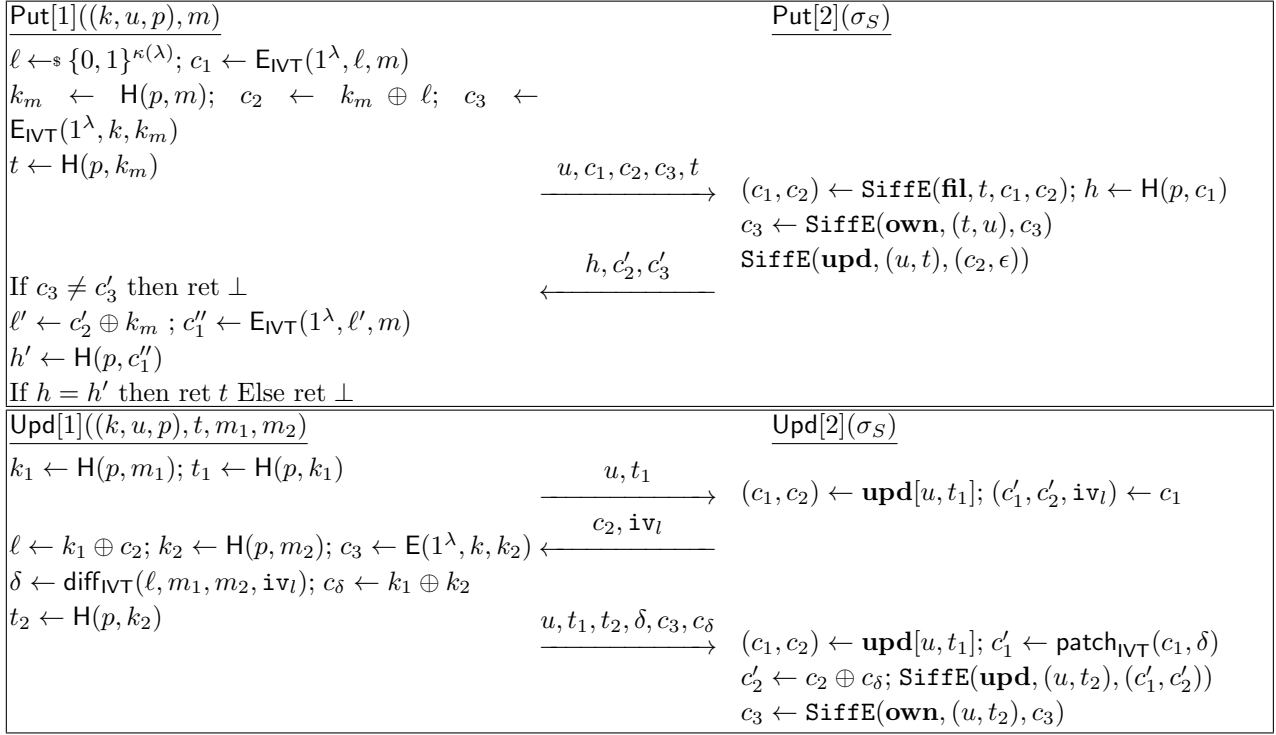


Figure 15: The Put and Upd protocols of IRCE2. The **fil** and **own** tables are immutable, and support the set-iff-empty operation (**SiffE**) explained in text.

length of the transmitted messages is $2\lambda + 6\kappa(\lambda) + |S|$, where $S = \mathbf{diff}_{IVT}(m_1, m_2)$ as above. Letting $\kappa(\lambda) = \lambda$, this in turn is bounded by $a \log(|m_1| + |m_2|) + b\lambda$ for small constants a and b , proving the proposition. \square

The proofs closely follow those from IRCE, so we reuse parts where applicable, and highlight the differences.

Theorem D.3. *If H is collision resistant, then IRCE2[H, E] is REC-secure.*

Proof. As with IRCE, the adversary **A** cannot hope to get a non-error output by furnishing a new t for **Get**, as the server checks that the tag was previously handled by the client. Now, if t was returned during a previous **Put** instance, then the same argument we used with IRCE carries over, as **Upd** does not play a role. Specifically, using the immutability of **fil** and **own**, we argue that the ciphertext stored in the server cannot change between the failed **Get** request and the last time the file was put. However, **Put** ensures that the hash of the ciphertext stored on the server matches with the hash of a correctly formed ciphertext of the plaintext. Consequently, if t does correspond to a **Put** instance, then **A** is in effect finding a pair of colliding inputs, namely the hash inputs involved in the comparison. We can argue that a CR adversary **B** can be built which has the same advantage as **A** in this case. The other case is when t is due to a ciphertext formed from an update. Consider such an **Upd** instance, and let $m_1, c_1, t_1, m_2, c_2, t_2$ denote the original and updated plaintext-ciphertext-tags. If c_2 is a valid ciphertext for m_2 , then by immutability and correctness of decryption, **Get** with t_2 will return m_2 . If c_2 is not a valid ciphertext for m_2 , then either c_1 is not a valid ciphertext for m_1 , or the update process must have introduced an error. The latter is ruled out by Proposition D.1. Now, if c_1 was placed on the server via a **Put** instance, we fall back to the first case (meaning we have an invalid ciphertext from a **Put** instance), and can build a CR adversary with matching advantage. If c_1 was placed via **Upd**, then the same argument can be applied recursively, until we reach a ciphertext from a **Put** instance. \square

PRIV SECURITY WITH EDIT DISTANCE. The following theorem shows PRIV-security. We consider a variant of the PRIV game of Figure 4, where the δ supplied to UPD is interpreted in terms of edit distance. Here, supporting delete is tricky, as the adversary could potentially start with an unpredictable plaintext, but then delete a sufficient number of characters to make the result predictable, and then break security. To prevent this, we require that the messages formed as result of the updates should also be unpredictable. We enforce this condition by making the source S output a list of allowed diff_e -outputs along with two tuples of messages. The adversary can only pick updates from this list. An edit source S is an algorithm that on input 1^λ returns $(\mathbf{m}_0, \mathbf{m}_1, \mathbf{d})$. Here \mathbf{d} is a list of diff_e -values. Consider source S' defined as below, where \parallel denotes adding an element to a tuple.

```


$$\begin{array}{l} \underline{S'(1^\lambda)} \\ (\mathbf{m}_0, \mathbf{m}_1, \mathbf{d}) \leftarrow S(1^\lambda) \\ \text{For } b \in \{0, 1\} \text{ and } m \in |\mathbf{m}_b| \text{ and } \delta \in \mathbf{d} \text{ do } \mathbf{m}_b \leftarrow \mathbf{m}_b \parallel \text{patch}_e(1^\lambda, m, \delta) \\ \text{Ret } (\mathbf{m}_0, \mathbf{m}_1) \end{array}$$


```

We require that S' should satisfy the formatting restrictions of regular sources, including length and equality. We say that S is unpredictable if S' is unpredictable. We now sketch the PRIV game for edit distances, with an edit source. Only the main, and UPD procedures differ. In main, the source S is executed to get $\mathbf{m}_0, \mathbf{m}_1, \delta$. In $\text{UPD}(i, \delta)$, before running the operations in Figure 4, a check that $\delta \in \mathbf{d}$ is performed.

PRF SECURITY. Let E be a blockcipher with blocksize $w(\lambda)$ and $\kappa(\lambda)$ -bit keys. The PRF game with E and adversary A starts by picking $b \leftarrow_s \{0, 1\}$ and $k \leftarrow_s \{0, 1\}^{\kappa(\lambda)}$ and runs A with access to a FN procedure, which A can query at points x . Upon such a query, if $b = 1$, the game returns $E(1^\lambda, k, x)$; otherwise, it returns a random but consistent value. Finally A exits with output b' and wins the game if $b = b'$. We define advantage $\text{Adv}_{E,A}^{\text{PRF}}(\lambda) = 2 \Pr[\text{PRF}_E^A(1^\lambda)] - 1$ and say that E is PRF-secure if no PT A has non-negligible advantage.

Theorem D.4. *If E is PRF-secure, then $\text{IRCE}_{\text{RO}}[E]$ is PRIV-secure.*

Proof. Let G_1 denote the $\text{PRIV}_{\text{IRCE}}$ game included with an unpredictable edit source S . Let A denote a PT adversary. Without loss of generality, assume that A makes only permitted UPD queries. Let n denote a bound on the number of plaintexts put on the server by A and $q_S(\lambda) : \mathbb{N} \rightarrow \mathbb{N}$ denote a bound on the number of RO queries made by S . Let G_2 denote the game similar to G_1 , where the c_3 , instead of being encryptions of the message-derived keys k_m , are replaced with encryptions of random strings. All the c_3 encryptions are performed under the legitimate client's secret key, which is never revealed to the server. There exists B such that

$$\text{Adv}_{\text{IVT}[E],B}^{\text{cpa}}(\lambda) = \Pr[G_1^A(1^\lambda)] - \Pr[G_2^A(1^\lambda)].$$

It follows that there exists another adversary B' such that $\text{Adv}_{\text{IVT}[E],B}^{\text{cpa}}(\lambda) = \text{Adv}_{E,B'}^{\text{PRF}}(\lambda)$. Let G_3 denote the game where in $\text{Put}[1, 1]$, the key and tag, which should be derived as $k_m \leftarrow \text{RO}(m)$, and $t \leftarrow \text{RO}(m)$ are picked instead as random $\kappa(\lambda)$ -bit strings, but if $p \parallel k_m$ or $p \parallel m$ have already been queried at the RO by S or A , then G_3 ensures consistency by using the existing values, but sets **bad**. When the adversary initiates an Update protocol, G_3 follows a similar procedure with k_2 and t_2 . Subsequent queries to RO at $p \parallel m$ or $p \parallel k_m$ are replied with k_m or t respectively, but sets **bad**. In G_4 , all the consistency measures are done away with, and the k_m and t values have no relations with the RO outputs on the associated \mathbf{m} points. We have $\Pr[G_3^A(1^\lambda)] \leq \Pr[G_4^A(1^\lambda)] + \Pr[G_4^A(1^\lambda) \text{ sets bad}]$.

Consider the PRF-game with n keys $\text{PRF}[n]$. We build a $\text{PRF}[n]$ -adversary C which runs A on G_4 . The n -keys of the game form the n -keys of A 's plaintexts. Adversary C starts by running S to get $\mathbf{m}_0, \mathbf{m}_1, \mathbf{d}$. Then it picks a random bit b and use \mathbf{m}_b in the rest of the simulation. When A calls $\text{PUT}(i)$, adversary C prepares a ciphertext for $\mathbf{m}_b[i]$. Here, c_2 is a random string, and c_1 , supposed to be the output of E_{IVT} , is formed by C making queries to its FN oracle. When A makes a (permitted) update query, then C forms the update plaintext m_2 and runs diff_{IVT} , with m, m_2 , and iv_l which it

can find from the server state it maintains using its FN oracle in lieu of the key. When C's FN oracle is implemented by E, it simulates A on G_4 ,

Consider G_5 , where E is replaced with a different random function for each key. If A makes no UPD queries, only the c_1 components depend on b . But these are CTR mode encryptions with a random function. No queries to the random function are repeated, and hence the c_1 values can be picked as random strings, independent of b . When A does make update queries, deleting blocks and modifying blocks does not help towards finding b ; inserting blocks could help, but in the IVT construction, in each ciphertext, a value iv_l is maintained, which keeps track of the last value of the counter used. When a new block is to be inserted, $\text{patch}_{\text{IVT}}$ increments IVT and uses a fresh counter value each time, meaning that these FN outputs can also be picked at random, independent of b and A cannot tell the difference. Overall, A learns no information about b in G_5 , and hence, $\Pr[G_5^A(1^\lambda)] = 1/2$. From a simple hybridization argument on PRF[n], we have $\Pr[G_4^A(1^\lambda)] - \Pr[G_5^A(1^\lambda)] = n(\lambda) \cdot \text{Adv}_{\text{E,C}}^{\text{PRF}}(\lambda)$.

Consider games H_1, H_2 and H_3 described as follows. Game H_1 is the same as G_4 , except that the winning condition in H_1 is setting **bad**. In game H_1 , both the source S and the adversary A can set **bad**, by querying the random oracle at a $\mathbf{m}[i]$ or $\mathbf{k}[i]$ point. The probability it sets **bad** is bounded by $q_S(\lambda)n(\lambda)/2^{\kappa(\lambda)}$. Game H_2 changes from H_1 only A query points are taken into account when testing for **bad**. In H_3 , the ciphertexts c_1 are derived as encryptions of random strings. There exist adversaries C' and D such that

$$\Pr[H_3^A(1^\lambda) \text{ sets bad}] - \Pr[H_2^A(1^\lambda) \text{ sets bad}] \leq n(\lambda)\text{Adv}_{\text{E,C}}^{\text{PRF}}(\lambda) + n(\lambda)\text{Adv}_{\text{E,D}}^{\text{kr}}(\lambda).$$

Here, C' works like C, except that it keeps track of **bad** queries, and D checks if any of the RO queries are keys in its KR game. Finally, in H_3 , the adversary learns nothing about the messages output by S, and we have

$$\Pr[H_3^{\text{S,A}}(1^\lambda) \text{ sets bad}] \leq q_S(\lambda)q(\lambda)2^{-\kappa(\lambda)} + q_A(\lambda)q(\lambda)\mathbf{GP}_S(\lambda).$$

where $q(\lambda)_A : \mathbb{N} \rightarrow \mathbb{N}$ is a bound on the number of RO queries made by A. Adding the equations so far, we have

$$\begin{aligned} \text{Adv}_{\text{IRCE2[E],S,A}}^{\text{priv}}(\lambda) &\leq 2(n(\lambda) + 1)\text{Adv}_{\text{SE,C}_1}^{\text{PRF}}(\lambda) + 2n(\lambda)\text{Adv}_{\text{E,C}_2}^{\text{kr}}(\lambda) + (q_S(\lambda) + q_A(\lambda))n(\lambda)2^{-\kappa(\lambda)} \\ &\quad + q_A(\lambda)n(\lambda)\mathbf{GP}_S(\lambda), \end{aligned}$$

where C_1 and C_2 are the ones among the PRF and KR adversaries with highest advantage. \square

E Parameter-dependent security: Proofs and extensions

Proposition E.1. *If FHE has evaluation correctness, then FCHECK[FHE, MLEWC] scheme supports deduplication.*

Proof. We need to show that 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) σ_C, σ'_C , for all $m \in \{0, 1\}^*$, the expected increase in size of σ''_S over σ'_S when $(f', \sigma'_S) \leftarrow_s \text{Run}(\text{Put}, (\sigma_C, m), \sigma_S)$ and $(f', \sigma''_S) \leftarrow_s \text{Run}(\text{Put}, (\sigma'_C, m), \sigma'_S)$ is bounded by $\ell(\lambda)$. In FCHECK, if a client with params σ_C runs Put with m , then, at the end, a parameter-ciphertext pair p, c (where $c = \text{E}(1^\lambda, \text{K}(1^\lambda, p, m), m)$) is stored on **fil**. Now, when the second client with parameters σ'_C tries to put m , the search in Put[2, 1] should find p, c . If the search happened over plaintexts, it is easy to see that the match will be detected, and as a result the client will only store an encryption of $\text{K}(1^\lambda, p, m), m$, which by assumption on MLEWC is independent of the size of m . But the search happens over FHE ciphertexts. We invoke evaluation correctness: as that the coins involved in K_f to generate the pk, sk in σ'_C , the coins in E_f to encrypt m , and the coins in Ev_f are all picked uniformly at random, except with negligible probability, the match is detected and the client decrypts to p , which leads the client to download p, c and hence stops the client from putting another ciphertext for m . \square

E.1 Proof of Theorem 4.2

We now that for all PT A , for all unpredictable PT sources S , there exists a PT unpredictable source S' and adversaries B and C such that

$$\text{Adv}_{\text{FCHECK}[\text{FHE}, \text{MLEWC}], S, A}^{\text{ldpriv}}(\lambda) \leq \text{Adv}_{\text{MLEWC}, S', C}^{\text{wpriv}}(\lambda) + \text{Adv}_{\text{FHE}, B}^{\text{cpa}}(\lambda) + 2m(\lambda)q(\lambda)\mathbf{GP}_S(\lambda). \quad (7)$$

where $q: \mathbb{N} \rightarrow \mathbb{N}$ is a bound on the total number of procedure queries made by A . Then, the theorem follows from the assumed CPA security of FHE, and from the WPRIV security of MLEWC. Consider games G_1 through G_4 of Figure 16. Without loss of generality, we assume that A does not repeat PUT queries. Here G_1 is the PDPRIV game with the code of S and $\text{FCHECK}[\text{FHE}, \text{MLEWC}]$. We have

$$\Pr[G_1^A(1^\lambda)] = \Pr[\text{PDPRIV}^{S, A}(1^\lambda)].$$

Game G_2 , as in the proof of Theorem 4.1, performs the search not over ciphertexts, but instead over plaintexts. However, following the argument from the proof of Theorem 4.1, the correctness of FHE ensures that this does not affect the outcome of the game. We have $\Pr[G_1^{S, A}(1^\lambda)] = \Pr[G_2^{S, A}(1^\lambda)]$.

In game G_2 , the c_2 components stored on the server correspond to the FHE encryptions of the MLEWC keys of the messages. These are replaced with random strings in G_3 . Moreover, in $\text{Put}[2, 1]$, when searching for a match for the put m , if a match is found, and the p, c pair for the match came as a result of a MSG instance (i.e. not from PUT), then the game sets **bad**. However, setting **bad** has not effect on the outcome of the game. Finally, in G_3 , the algorithms in the Put protocol do not send the plaintext m in messages, but instead use the index of the plaintext in \mathbf{m}_b . This change does not affect the outcome of the game either, and only serves to simplify the code. There exists an adversary B such that

$$\Pr[G_2^{S, A}(1^\lambda)] - \Pr[G_3^{S, A}(1^\lambda)] \leq \frac{1}{2}\text{Adv}_{\text{FHE}, B}^{\text{cpa}}(\lambda).$$

The description of B is straightforward, and we omit it here. Game G_4 is identical-until-**bad** to G_3 . From the fundamental lemma of game-playing [15], we have

$$|\Pr[G_3^{S, A}(1^\lambda)] - \Pr[G_4^{S, A}(1^\lambda)]| \leq \Pr[G_4^{S, A}(1^\lambda) \text{ sets bad}].$$

To set **bad** in G_4 , adversary A must produce a p, c pair that is a valid encryption of some $\mathbf{m}_b[i]$, put it on the server with MSG, and subsequently run $\text{PUT}(i)$ followed by STEP to make the search at $\text{Put}[2, 1]$ find a match at p, c . However, A receives no information about the output of S , not even the ciphertexts. To see the ciphertexts, A must query STATE, but it can no longer query STEP after doing so. Thus, setting **bad** in G_4 can be bounded by $m(\lambda)q(\lambda)\mathbf{GP}_S(\lambda)$.

Consider source S' and adversary C which work as follows. S' picks coins r at random and starts running G_4^A (except for picking the random bit b) with r as the only source of randomness up until the point when A makes a PTXT query with input d . At this point, S' runs $S(1^\lambda, d)$ with fresh coins to get $\mathbf{m}_0, \mathbf{m}_1$ and exits with output $\mathbf{m}_0, \mathbf{m}_1, r$. Adversary C , when invoked with $\mathbf{p}, \mathbf{c}, r$ runs G_4^A with r , and when a PTXT(d) query is made, it lets \mathbf{p}, \mathbf{c} play the role of the corresponding variables in G_4 . When A finishes with output b' , then C also exits with output b' . Together, S' and C simulate A in G_4 , and hence it follows that $\Pr[\text{WPRIV}_{\text{MLEWC}}^{S', C}(1^\lambda)] = \Pr[G_4^{S, A}(1^\lambda)]$. Moreover, since S' runs S on fresh coins, which are not provided to C , it follows that S' is unpredictable if S is unpredictable. Adding up the above equations leads to Equation (7), completing the proof.

F MLEWC: Proofs and extensions

F.1 Proof of Theorem 4.3

Proof. Let S be an unpredictable auxiliary source and A be an adversary. Let $m: \mathbb{N} \rightarrow \mathbb{N}$ denote a bound on the length of message tuples output by S . Consider the constructions of PF source S' and B described in Figure 18. It can be checked that $\text{Adv}_{\text{HtO}[\text{HF}, \text{os}], S, A}^{\text{wpriv}}(\lambda) \leq \text{Adv}_{\text{OS}, S', B}^{\text{cdipfo}}(\lambda)$.

<p><u>MAIN(1^λ)</u> // $G_1^A(1^\lambda) - G_4^A(1^\lambda)$ $b \leftarrow_s \{0, 1\}$; $\sigma_S \leftarrow_s \text{Init}(1^\lambda)$ $b' \leftarrow_s \mathbf{A}^{\text{PUT, STEP, PTXT, MSG, REG, STATE}}(1^\lambda)$; Ret ($b = b'$)</p> <p><u>Put[1, 2]($m, \mathbf{M}$)</u> $p, i \leftarrow \text{D}_f(1^\lambda, sk, c_r)$ If $p = 0^{\kappa(\lambda)}$ then $p \leftarrow_s \mathbf{P}(1^\lambda)$; $k \leftarrow_s \mathbf{K}(1^\lambda, p, m)$; $c_1 \leftarrow \text{E}(1^\lambda, m)$ Else $k \leftarrow_s \mathbf{K}(1^\lambda, p, m)$ $c_2 \leftarrow \mathbf{E}_{pk}(k)$ Ret $m, (c_1, c_2, p, u, i), 2, \text{False}$</p> <p><u>Put[2, 2]($\sigma_S, \mathbf{M}$)</u> $(n_c, \mathbf{U}, \mathbf{fil}, \mathbf{own}) \leftarrow \sigma_S$; $(u, p, c_1, c_2, i) \leftarrow \mathbf{M}$ If $c_1 \neq \epsilon$ then $n_f \leftarrow n_f + 1$; $i \leftarrow n_f$; $\mathbf{fil}[i] \leftarrow (p, c_1)$ $\mathbf{own}[u, i] \leftarrow c_2$; $\sigma_S \leftarrow (p, \mathbf{U}, \mathbf{fil}, \mathbf{own})$ $\mathbf{M} \leftarrow i$; Ret $\sigma_S, \mathbf{M}, 1, \text{True}$</p> <p><u>PTXT($d$)</u> // $G_3^A(1^\lambda), G_4^A(1^\lambda)$ $\vec{m}_0, \vec{m}_1 \leftarrow_s \mathbf{S}(1^\lambda, d)$ For $i = 1$ to \mathbf{m}_b $\mathbf{p}[i] \leftarrow_s \mathbf{P}(1^\lambda)$; $\mathbf{k}[i] \leftarrow \mathbf{K}(1^\lambda, \mathbf{p}[i], \mathbf{m}_b[i])$ $\mathbf{c}[i] \leftarrow \text{E}(1^\lambda, \mathbf{k}[i], \mathbf{m}_b[i])$</p> <p><u>Put[1, 2]($i, \mathbf{M}$)</u> // $G_3^A(1^\lambda), G_4^A(1^\lambda)$ $p, i \leftarrow \text{D}_f(1^\lambda, sk, c_r)$ If $p = 0^{\kappa(\lambda)}$ then $p \leftarrow \mathbf{p}[i]$; $k \leftarrow_s \mathbf{k}[i]$; $c_1 \leftarrow \mathbf{c}[i]$ Else $k \leftarrow \mathbf{K}(1^\lambda, p, m)$ $k' \leftarrow_s \{0, 1\}^{ k }$; $c_2 \leftarrow \mathbf{E}_{pk}(k')$ Ret $m, (c_1, c_2, p, u, i), \text{False}$</p> <p><u>Put[1, 1]($m, \mathbf{M}$)</u> // $G_2^A(1^\lambda)$ Ret m, m, False</p>	<p><u>Put[1, 1](m, \mathbf{M})</u> // $G_1^A(1^\lambda)$ $c_f \leftarrow_s \mathbf{E}_{pk}(m)$; Ret m, c_f, False</p> <p><u>Put[2, 1](σ_S, \mathbf{M})</u> // $G_1^A(1^\lambda)$ $(n_c, \mathbf{U}, \mathbf{fil}, \mathbf{own}) \leftarrow \sigma_S$ $c_r \leftarrow_s \text{E}_f(1^\lambda, pk, 0^{\kappa(\lambda)})$; $c_i \leftarrow_s \text{E}_f(1^\lambda, pk, 0)$; $c_n \leftarrow c_i$ For $(p, c) \in \mathbf{fil}$ do $c_p \leftarrow_s \text{E}_f(1^\lambda, pk, p)$; $c_c \leftarrow_s \text{E}_f(1^\lambda, pk, c)$ $c_r, c_n, c_i \leftarrow_s \text{E}_{vf}(1^\lambda, pk, \mathbf{cmp}, c_f, c_p, c_c, c_r, c_n, c_i)$ $\sigma_S \leftarrow (n_c, \mathbf{U}, \mathbf{fil}, \mathbf{own})$; $\mathbf{M} \leftarrow c_r, c_n$; Ret $\sigma_S, \mathbf{M}, \text{False}$</p> <p><u>Put[2, 1]($\sigma_S, m$)</u> // $G_2^A(1^\lambda)$ $(n_c, \mathbf{U}, \mathbf{fil}, \mathbf{own}) \leftarrow \sigma_S$; $r \leftarrow 0^{\kappa(\lambda)}$; $i \leftarrow 0$; $n \leftarrow 0$ For $(k, c) \in \mathbf{fil}$ do $(r, n, i) \leftarrow \text{Eval}(\mathbf{cmp}, m, p, c, r, n, i)$ $\mathbf{M} \leftarrow (\text{E}_f(1^\lambda, pk, r), \text{E}_f(1^\lambda, pk, n))$; Ret $\sigma_S, \mathbf{M}, \text{False}$</p> <p><u>PUT($i$)</u> // $G_3^A(1^\lambda), G_4^A(1^\lambda)$ $\mathbf{p} \leftarrow \mathbf{p} + 1$; $\mathbf{PS}[\mathbf{p}] = \text{Put}$; $\mathbf{a}[\mathbf{p}, 1] \leftarrow i$ $\mathbf{N}[\mathbf{p}] \leftarrow 1$; $\mathbf{M}[\mathbf{p}] \leftarrow \epsilon$; Ret \mathbf{p}</p> <p><u>Put[1, 1](i, \mathbf{M})</u> // $G_3^A(1^\lambda), G_4^A(1^\lambda)$ Ret i, i, False</p> <p><u>Put[2, 1](σ_S, j)</u> // $G_3^A(1^\lambda), G_4^A(1^\lambda)$ $m \leftarrow \mathbf{m}_b[j]$; $(n_c, \mathbf{U}, \mathbf{fil}, \mathbf{own}) \leftarrow \sigma_S$; $r \leftarrow 0^{\kappa(\lambda)}$ $i \leftarrow 0$; $n \leftarrow 0$ For $(p, c) \in \mathbf{fil}$ do $k \leftarrow \mathbf{K}(1^\lambda, p, m)$; $c' \leftarrow \text{E}(1^\lambda, k, m)$; $i \leftarrow i + 1$ If $c = c'$ and $p \neq \mathbf{p}[j]$ then $\text{bad} \leftarrow \text{True}$; $i_f \leftarrow i$; $p_f \leftarrow p$ $\mathbf{M} \leftarrow (\text{E}_f(1^\lambda, pk, p_f), \text{E}_f(1^\lambda, pk, i_f))$; Ret $\sigma_S, \mathbf{M}, \text{False}$</p>
--	--

Figure 16: Games G_1 through G_4 of Theorem 4.2. Procedures with code unchanged from PDP_{PRIV} are omitted.

<p><u>MAIN(1^λ)</u> // WPRIV^{S,A}(1^λ) $(\mathbf{m}_0, \mathbf{m}_1, z) \leftarrow_s \mathbf{S}(1^\lambda, \epsilon)$; $b \leftarrow_s \{0, 1\}$ For $i \in [\mathbf{m}_b]$ do $\mathbf{p}[i] \leftarrow_s \mathbf{P}(1^\lambda)$; $\mathbf{k}[i] \leftarrow_s \mathbf{K}(1^\lambda, \mathbf{p}[i], \mathbf{m}_b[i])$ $\mathbf{c}[i] \leftarrow_s \text{E}(1^\lambda, \mathbf{k}[i], \mathbf{m}_b[i])$ $b' \leftarrow_s \mathbf{A}_2(1^\lambda, \mathbf{p}, \mathbf{c}, z)$; Ret ($b = b'$)</p>	<p><u>MAIN(1^λ)</u> // CDIPFO^{S,A}(1^λ) $(\mathbf{p}, z) \leftarrow_s \mathbf{S}(1^\lambda)$; $b \leftarrow_s \{0, 1\}$ For $i \in [\mathbf{p}]$ do If $b = 1$ then $(\alpha, \beta) \leftarrow \mathbf{p}[i]$; $\mathbf{F}[i] \leftarrow_s \text{Obf}(1^\lambda, (\alpha, \beta))$ Else $(\alpha', \beta') \leftarrow \mathbf{p}[i]$; $\alpha \leftarrow_s \{0, 1\}^{ \alpha' }$; $\beta \leftarrow_s \{0, 1\}^{ \beta' }$ $\mathbf{F}[i] \leftarrow_s \text{Obf}(1^\lambda, (\alpha, \beta))$ $b' \leftarrow_s \mathbf{A}(1^\lambda, \mathbf{F}[i], z)$; Ret ($b = b'$)</p>
---	--

Figure 17: The WPRIV game on the left, and the and CDIPFO game on the right.

It remains to show that S' is unpredictable. Consider source S_2 works by running S to get $(\mathbf{m}_0, \mathbf{m}_1, z)$, then picks keys $\mathbf{k}_H[i] \leftarrow_s \mathbf{K}_h(1^\lambda)$ and a bit b , and for $i \in [|\mathbf{m}_b|]$, computes $\mathbf{k}_H[i] \leftarrow_s \mathbf{K}_h(1^\lambda)$

$\underline{S'(1^\lambda)}$ $b \leftarrow_s \{0, 1\}; (\mathbf{m}_0, \mathbf{m}_1, z) \leftarrow_s S(1^\lambda)$ For $i \in [\mathbf{m}_b]$ do $\mathbf{k}_H[i] \leftarrow_s K_h(1^\lambda); \mathbf{k}[i] \leftarrow H(1^\lambda, \mathbf{k}_H[i], \mathbf{m}_b[i])$ For $j \in [\mathbf{m}_b[i]]$ do $\mathbf{m}'[i\ell(\lambda) + j] \leftarrow \mathbf{k}[i][\langle \ell, i \rangle] \mathbf{m}_b[i, j]$ Ret $\mathbf{m}', (b, \mathbf{k}_H, z)$	$\underline{B(1^\lambda, F, (b, \mathbf{k}_H, z))}$ For $i \in [m(\lambda)]$ do $\mathbf{c}[i] \leftarrow (F[(i-1)\ell(\lambda) + 1], \dots, F[i\ell(\lambda)])$ $b' \leftarrow_s A(1^\lambda, \mathbf{k}_H, \mathbf{c}, z)$ If $b' = b$ return 1 else return 0
---	---

Figure 18: Source S' and adversary B of Theorem 4.3.

$\underline{\text{MAIN}(1^\lambda)} \quad // \text{MUCE}_{\text{HF}}^{\text{S,D}}(1^\lambda)$ $(1^n, t) \leftarrow_s S(1^\lambda, \epsilon); \text{For } i = 1 \text{ to } n \text{ do } \mathbf{k}_H[i] \leftarrow_s K(1^\lambda)(1^\lambda)$ $b \leftarrow_s \{0, 1\}; L \leftarrow_s S^{\text{HASH}}(1^n, t); b' \leftarrow_s D(1^\lambda, \mathbf{k}_H, L)$ Return $(b' = b)$ $\underline{\text{HASH}(x, 1^\ell, i)}$ If $T[x, \ell, i] = \perp$ then If $b = 1$ then $T[x, \ell, i] \leftarrow H(1^\lambda, \mathbf{k}_H[i], x, 1^\ell)$ Else $T[x, \ell, i] \leftarrow_s \{0, 1\}^\ell$ Return $T[x, \ell, i]$	$\underline{\text{MAIN}(1^\lambda)} \quad // \text{MPRED}_S^{\text{P}}(1^\lambda)$ $(1^n, t) \leftarrow_s S(1^\lambda, \epsilon); \text{done} \leftarrow \text{False}$ $\mathbf{P} \leftarrow \emptyset; L \leftarrow_s S^{\text{HASH}}(1^n, t)$ $\text{done} \leftarrow \text{True}; \mathbf{P}' \leftarrow_s \text{P}^{\text{HASH}}(1^\lambda, 1^n, L)$ Return $(\mathbf{P} \cap \mathbf{P}' \neq \emptyset)$ $\underline{\text{HASH}(x, 1^\ell, i)}$ If $\text{done} = \text{False}$ then $\mathbf{P} \leftarrow \mathbf{P} \cup \{x\}$ If $T[x, \ell, i] = \perp$ then $T[x, \ell, i] \leftarrow_s \{0, 1\}^\ell$ Return $T[x, \ell, i]$
--	--

Figure 19: The MUCE and MPRED games.

and $\mathbf{k}[i] \leftarrow H(1^\lambda, \mathbf{k}_H[i], \mathbf{m}_b[i])$. Then it outputs $\mathbf{k}, (\mathbf{k}_H, b, z)$. Clearly, unpredictability of S' follows from that of S_2 . Further, consider $S_2[i]$ which operates as above, but only outputs $\mathbf{k}[i], (\mathbf{k}_H[i], b, z)$. Once again, if $S_2[i]$ is unpredictable for all $i \in [m(\lambda)]$, then S_2 is also unpredictable. Note that $S_2[i]$ runs S , an unpredictable source, picks one of its outputs $\mathbf{m}_b[i]$, generates $k_H \leftarrow_s K_h(1^\lambda)$ and outputs k_H and the hash $H(1^\lambda, k_H, \mathbf{m}_b[i])$. Knowing that CR hash functions are randomness condensers [29], it is easy to show that $S_2[i]$ is unpredictable for $i \in [m(\lambda)]$. We omit the details. \square

REMARKS ON HtO. The HtO construction can be made more efficient by using a CDIPFO where the special output can depend on the special input, so that the plaintext can be obfuscated in one shot, instead of bit-by-bit. However, such obfuscators in the standard model come only from UCes. The HtO construction can be modified so that ciphertext length is only additive, by changing E to encrypt m under an SE scheme with a fresh random key, and encrypting the key as above. Now OS should be secure against CDIPFO-secure against computationally unpredictable sources and this can be achieved from the construction in [16], by extending the t -Strong Vector Decision Diffie Hellman assumption to computationally unpredictable distributions.

F.2 PRV\$-CDA secure MLE from UCE

A family of functions $\text{HF} = (K_h, H)$ is a pair of deterministic algorithms. Key generation $K(1^\lambda)$ returns a key $k \in \{0, 1\}^{\kappa(\lambda)}$ on input 1^λ , and evaluation H takes 1^λ , a key k , an input $m \in \{0, 1\}^*$, and a unary encoding 1^ℓ of an output length to return an output $H(1^\lambda, k, x, 1^\ell) \in \{0, 1\}^\ell$.

We now recall the definition of statistical multi-key security UCE for hash functions. Consider the game MUCE of Figure 19. A source is an algorithm which begins by keys indicating n the number of instances, along with state t . The game creates n independent keys. The source gets a procedure HASH which is either implemented by HF with n -independent keys, or via n random functions depending on the bit b chosen in the game. Then S returns with output leakage $L \in \{0, 1\}^*$

and the distinguisher D on input L and all keys should guess b to win. We associate advantage $\text{Adv}_{\text{HF}, \mathcal{S}, D}^{\text{m-uce}}(\lambda) = 2 \Pr[\text{MUCE}_{\text{HF}}^{\mathcal{S}, D}(\lambda)] - 1$.

A statistical predictor P is an algorithm (not necessarily PT) if there exist polynomials q, s such that for all $\lambda \in \mathbb{N}$, in the MPRED game predictor P makes at most $q(\lambda)$ oracle queries and outputs a set \mathbf{P}' of size at most $s(\lambda)$. A source \mathcal{S} is multi-key statistically unpredictable if for all statistical predictors P , it holds that $\text{Adv}_{\mathcal{S}, P}^{\text{pred}}(\lambda) = \Pr[\text{MPRED}_{\mathcal{S}}^P(\lambda)]$ is negligible. We say that HF is statistical multi-key UCE secure ($\text{UCE}[\mathcal{S}^{\text{sup-m}}]$ -secure), if it holds that for all $\text{Adv}_{\text{HF}, \mathcal{S}, D}^{\text{m-uce}}(\lambda)$ is negligible for all PT statistical unpredictable \mathcal{S} , for all PT D .

Consider the MLE scheme $\text{CE}[\text{HF}] = (P, K, E, D, T)$ described below, where $T(1^\lambda, p, c) = c$.

$$\begin{array}{l|l|l|l} \underline{P(1^\lambda)} & \underline{K(1^\lambda, k_H, m)} & \underline{E(1^\lambda, k_H, k, m)} & \underline{D(1^\lambda, k_H, k, c)} \\ k_H \leftarrow_s K_h(1^\lambda) & k \leftarrow H(1^\lambda, k_H, m, 1^\lambda) & c \leftarrow m \oplus H(1^\lambda, k_H, k, 1^{|m|}) & m \leftarrow c \oplus H(1^\lambda, k_H, k, 1^{|c|}) \\ \text{Return } k_H & \text{Return } k & \text{Return } c & \text{Return } m \end{array}$$

Correctness of the scheme is easy to check. The following theorem shows that $\text{CE}[\text{HF}]$ is WPRIV -secure if HF is statistical mUCE secure.

Theorem F.1. If HF is $\text{UCE}[\mathcal{S}^{\text{sup-m}}]$ -secure, then $\text{CE}[\text{HF}]$ is WPRIV -secure.

Proof. Let A be a PT unpredictable WPRIV adversary with functions m, ℓ . Consider \mathcal{S}, D described below.

$$\begin{array}{l|l} \underline{S^{\text{HASH}}(1^\lambda)} & \underline{D(1^\lambda, \mathbf{k}_H, L)} \\ b \leftarrow_s \{0, 1\}; \mathbf{m}_0, \mathbf{m}_1, z \leftarrow_s A_1(1^\lambda) & \mathbf{c}, b, z \leftarrow L \\ \text{For } i = 1 \text{ to } |\mathbf{m}_b| \text{ do} & b' \leftarrow A_2(1^\lambda, \mathbf{k}_H, \mathbf{c}, z) \\ \quad \mathbf{k}[i] \leftarrow \text{HASH}(\mathbf{m}[i], 1^\lambda, i) & \text{If } b = b' \text{ then return 1 else return 0} \\ \quad \mathbf{c}[i] \leftarrow \mathbf{m}[i] \oplus \text{HASH}(\mathbf{k}[i], 1^{|\mathbf{m}[i]|}, i) & \\ L \leftarrow \mathbf{c}, b, z; \text{Return } L & \end{array}$$

It can be seen that $\text{Adv}_{\text{CE}[\text{HF}], A}^{\text{wpriv}}(\lambda) = 2 \text{Adv}_{\text{HF}, \mathcal{S}, D}^{\text{m-uce}}(\lambda)$. It remains to show that \mathcal{S} is simple statistically unpredictable by Lemma 4.7 of [10]. If P is a simple predictor, it follows that $\text{Adv}_{P, \mathcal{S}}^{\text{m-spred}}(\lambda) \leq m \cdot (\mathbf{GP}_A(\lambda) + 2^{-\lambda})$. Simple statistical unpredictability of \mathcal{S} follows from unpredictability of A . \square