

What is the *Exact* Security of the Signal Protocol?

Alexander Bienstock¹, Jaiden Fairuze², Sanjam Garg^{2,4}, Pratyay Mukherjee³, and Srinivasan Raghuraman³

¹New York University

²UC Berkeley

³Visa Research

⁴NTT Research

Abstract

In this work we develop comprehensive definitions in the Universal Composability framework to study the Signal Double Ratchet (**Signal** for short) protocol. Our definitions enable a more fine-grained and rigorous analysis of the *exact* security of **Signal** by explicitly capturing many new security guarantees, in addition to the ones that were already identified in the state-of-art work by Alwen, Coretti and Dodis [Eurocrypt 2019]. Moreover, our definitions provide the ability to more easily build on top of **Signal**, using the UC Composition Theorem. The **Signal** protocol, as it is described in the whitepaper, securely realizes our ideal functionality $\mathcal{F}_{\text{Signal}}$. However, as we interpret from the high-level description in the whitepaper, the guarantees of $\mathcal{F}_{\text{Signal}}$ seem slightly weaker than those one would expect **Signal** to satisfy. Therefore we provide a stronger, more natural definition, formalized through the ideal functionality $\mathcal{F}_{\text{Signal}^+}$. Based on our definitions we are able to make many important insights as follows:

- We observe *several shortcomings* of Alwen et al.’s game-based security notions. To demonstrate them, we construct *four* different modified versions of the **Signal** protocol, all of which are insecure according to our weaker $\mathcal{F}_{\text{Signal}}$ definition, but remain secure according to their definitions. Among them, one variant was suggested for use by Alwen et al.; another one was suggested for use by the **Signal** whitepaper itself.
- We identify the *exact assumptions* required for the full security of **Signal**. In particular, our security proofs use the gap-Diffie-Hellman assumption and the random oracle model, as opposed to the DDH assumption and standard model used in Alwen et al.
- We demonstrate the *shortcomings of Signal* with respect to our stronger functionality $\mathcal{F}_{\text{Signal}^+}$ by showing a non-trivial (albeit minor) weakness with respect to that definition.
- Finally, we complement the above weakness by providing a *minimalistic modification* to **Signal** (that we call the Triple Ratchet) and show that the resulting protocol securely realizes the stronger functionality $\mathcal{F}_{\text{Signal}^+}$. Remarkably, the modification incurs no additional communication cost and virtually no additional computational cost.

Contents

1	Introduction	3
1.1	Our work	4

1.2	Double Ratchet High-Level Summary	4
1.3	Technical Highlights of our Results	7
1.3.1	Shortcoming of ACD's Definition.	7
1.3.2	Exact Security Assumptions for Signal.	9
1.3.3	Shortcoming of Signal with Respect to $\mathcal{F}_{\text{Signal}^+}$.	9
1.3.4	Modification to Achieve $\mathcal{F}_{\text{Signal}^+}$: the Triple Ratchet protocol	10
1.4	Related Work	10
1.5	Roadmap for the Rest of the Paper	11
2	Defining Security of the Signal Protocol	12
2.1	Honest Execution	14
2.2	Execution with an Unrestricted Adversary	15
A	Preliminaries	23
A.1	Game-Based Security and Notation	23
A.2	Authenticated Encryption	24
B	Building Blocks	25
B.1	Key Derivation Function Chains	26
B.1.1	Differences from [ACD19]	26
B.2	Forward-Secure AEAD	27
B.2.1	Defining FS-AEAD	27
B.2.2	Differences from [ACD19]	29
B.2.3	Instantiating FS-AEAD	30
B.3	Continuous Key Agreement	31
B.3.1	Defining CKA	32
B.3.2	Differences from [ACD19]	34
B.3.3	Instantiating CKA	36
B.3.4	Instantiating CKA ⁺	36
B.3.5	Even stronger security	38
C	Composition	38
C.1	Constructions	39
C.2	Vulnerability of Signal with Respect to $\mathcal{F}_{\text{Signal}^+}$	41
D	Limitations of ACD's Secure Messaging Security Notion	42
D.1	Definitions from ACD [ACD19]	43
D.1.1	Secure Messaging (ACD)	43
D.1.2	ACD's CKA Definition	45
D.1.3	ACD's FS-AEAD Definition	46
D.1.4	ACD's PRF-PRNG Definition	46
D.1.5	ACD's Composition	47
D.2	The Transformations to Signal and Their (In)Security	48
D.2.1	T_1 : Postponed FS-AEAD Key Deletion	48
D.2.2	T_2 : Postponed CKA Key Deletion	49
D.2.3	T_3 : Eager CKA Randomness Sampling	51

D.2.4	T_4 : Malleable CKA	53
E	UC Security of Signal and Signal⁺	55
E.1	Hybrid Algorithms H_{t^*} and the Simulator \mathcal{S}	56
E.2	Type 1 Adversaries	58
E.3	Type 2 Adversaries	59
E.4	Type 3 Adversaries	62
E.5	Even Stronger Security of Signal ⁺ with CKA ⁺	64
F	The Model in Detail	64
F.1	UC Security: A Brief Overview	64
F.1.1	The Basic Model of Computation	65
F.1.2	Security of Protocols.	66

1 Introduction

Background. The Signal protocol is by far the most popular end-to-end secure messaging (SM) protocol, boasting of billions of users. The core underlying technique of the Signal protocol is commonly known as the *Double Ratchet*¹ algorithm and is beautifully explained in the whitepaper [MP16a] authored by the creators of Signal, Marlinspike and Perrin. The whitepaper also outlines the desired security properties, and provides intuitions on the design rationale for achieving them. Indeed, in addition to standard security against an eavesdropper who may modify ciphertexts, Signal attempts to achieve (i) *post-compromise security* (PCS) and *forward secrecy* (FS) with respect to leakages of secret state, (ii) *resilience to bad randomness*, and (iii) *immediate decryption*. PCS requires the conversation to naturally and quickly recover security after a leakage on one of the (or both) parties, as long as if the affected parties have good randomness. FS requires past messages to remain secure even after a leakage on one of the (or both) parties. Resilience to bad randomness requires that as long as if both parties’ secret states are secure (i.e., PCS has been achieved after any corruptions), then the conversation should remain secure, even if bad randomness is used in crafting messages. Finally, immediate decryption requires parties to — immediately upon reception of ciphertexts — obtain underlying plaintexts and place them in the correct order in the conversation, even if they arrive arbitrarily out of order and if some of them are completely lost by the network (the latter is also known as *message-loss resilience*).

However, despite the elegance and simplicity of the Double Ratchet, capturing its security turned out to be not so straightforward. The first formal analysis of the Signal protocol was provided by Cohn-Gordon et al. [CCD⁺20]. However, this analysis was complex and left open several question about the cryptographic security achieved by Signal. Following in the footsteps of Cohn-Gordon et al., a more generic and more comprehensive security definitions for Signal were provided by Alwen et al. in Eurocrypt 2019 [ACD19] (referred to as ACD henceforth), with close focus on immediate decryption property of the Signal protocol. They provided a modular analysis with respect to game-based definitions proposed therein. Indeed, they introduced new abstract primitives and composed them into SM protocols (including Signal itself) that capture the above properties: Continuous Key Agreement (CKA), Forward-Secure Authenticated Encryption with Associated Data (FS-AEAD), and PRF-PRNGs. While ACD’s work significantly improved our

¹Throughout the paper we use Double Ratchet and Signal interchangeably.

understanding of Signal’s Double Ratchet protocol, as we observe in this work, their definitional framework still falls short of capturing the *exact* security provided by the Signal protocol and the cryptographic strength that it requires.

1.1 Our work

In this work we provide a deeper insight on the security of Signal by developing more comprehensive definitions, in the Universal Composability [Can01] (UC) framework. We first note that providing a definition in the UC framework provides the ability to more easily build on top of Signal, using the UC Composition Theorem. Perhaps more importantly, our simulation-based definitional framework enables a more fine-grained and rigorous analysis of the *exact* security of Signal. In particular, we define two ideal functionalities, $\mathcal{F}_{\text{Signal}}$ and $\mathcal{F}_{\text{Signal}^+}$; as the name suggests $\mathcal{F}_{\text{Signal}^+}$ captures a strictly stronger notion than $\mathcal{F}_{\text{Signal}}$. The Signal protocol, as it is described in the whitepaper [MP16a] (in its strongest form), satisfies the functionality $\mathcal{F}_{\text{Signal}}$. However, as we interpret from the intuitive description in the whitepaper, the guarantees of $\mathcal{F}_{\text{Signal}}$ seem slightly weaker than those that one would expect Signal to satisfy. Therefore we provide a strengthened definition, formalized through the ideal functionality $\mathcal{F}_{\text{Signal}^+}$. Based on this framework we are able to make many important insights, as summarized below:

- We observe several shortcomings of the ACD notion. To demonstrate them we construct *four* different modified versions of the Signal protocol, all of which are insecure according to our weaker $\mathcal{F}_{\text{Signal}}$ definition, but remain secure according to ACD’s definition.
- We identify the exact assumptions that seem necessary for the full security of the Double Ratchet. In particular, our security proofs use the gap-Diffie-Hellman [OP01] (gap-DH) assumption and the random oracle model, as opposed to the decisional Diffie Hellman (DDH) assumption and standard model used in ACD.
- We find a non-trivial (albeit minor) weakness of Signal with respect to our stronger definition captured by $\mathcal{F}_{\text{Signal}^+}$. Note: this weakness is allowed in the definition in which ACD prove Signal’s security.
- Finally, we complement the above weakness by providing a minimalistic modification to Signal and prove the resulting protocol secure according to the stronger definition $\mathcal{F}_{\text{Signal}^+}$. We call this new protocol as the Triple Ratchet as it add another “mini ratchet” to the public-key ratchet in the Double Ratchet Protocol. Remarkably, the modification incurs no additional communication cost and virtually no additional computational cost.

We provide more details on these result in Section 1.3. Before explaining the details, below we first provide a brief summary of the Double Ratchet protocol, which will be helpful in understanding our results. For a detailed description we refer to the Double Ratchet whitepaper [MP16a]. Readers familiar with the Double Ratchet protocol could easily skip Section 1.2.

1.2 Double Ratchet High-Level Summary

Here we give a high-level overview of the Signal Double Ratchet before formally defining it in the main body. We note that although we here describe the double ratchet specifically in terms of its real-world implementation [MP16a], our paper still breaks it down into modular pieces which can

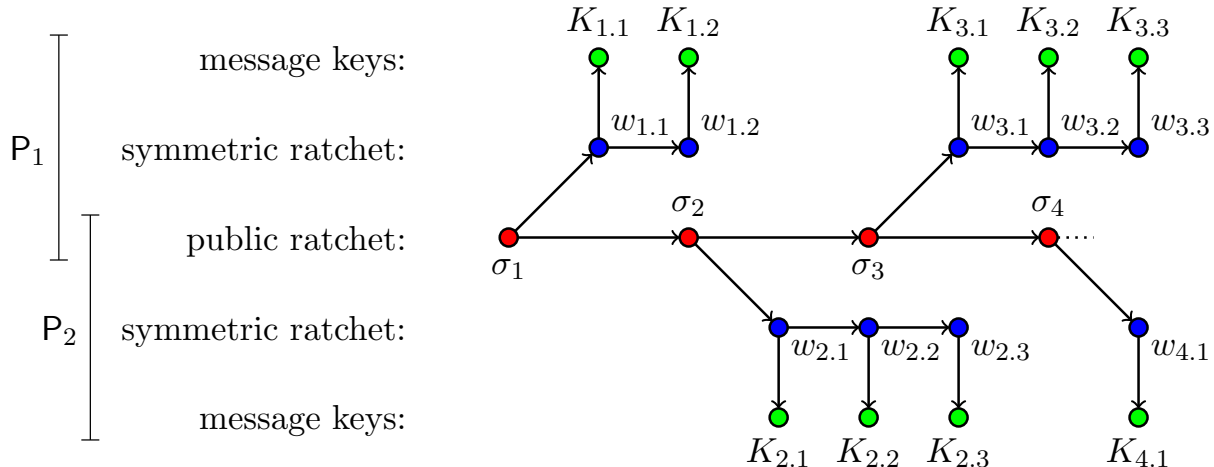


Figure 1: Sample Double Ratchet key evolution. In this depiction, P_1 sends and P_2 receives in epoch 1, followed by P_2 sending and P_1 receiving in epoch 2, and so on. As explained in the main body, initial symmetric chain keys $w_{i,1}$ for each epoch i are derived first by the sender, then also by the receiver, using the shared root keys σ_i and asynchronously exchanged shared secrets (via DDH). Then, updated symmetric chain keys $w_{i,j}$ and message keys $K_{i,j}$ are derived deterministically from $w_{i,1}$.

be instantiated in several different ways, as in [ACD19]. For the purpose of our paper, we assume that the two participants P_1 and P_2 share a common secret upon initialization. In Signal, this is achieved via the X3DH key exchange protocol [MP16b], but we consider this out of scope for our study of the double ratchet. Using their initial shared secret, P_1 and P_2 can derive the initial *root key* σ which seeds the public ratchet. Furthermore, upon initialization P_2 also holds some secret exponent x and P_1 holds the corresponding public value g^x . Once the initialization process completes, the ratcheting session can begin.

At its core, the double ratchet has two key components: the outer public-key ratchet, and the inner symmetric-key ratchet (often referred to as simply the public and symmetric ratchets, respectively). ACD abstract out the symmetric ratchet as their FS-AEAD primitive, the update mechanism of the public ratchet as their PRF-PRNG primitive, and the means by which shared secrets are produced to update the public ratchet as their CKA primitive. The goal of the double ratchet is to provide distinct *message keys* to encrypt/decrypt each new message. For each message the same message key is derived by both parties using a symmetric *chain key* which itself is derived from the aforementioned root key. Naturally, this results in a key *hierarchy* with the root key at the top, chain keys at an intermediate layer, and message keys at the bottom. Observe a graphical depiction of this hierarchy in Figure 1. In the Signal double ratchet, Diffie-Hellman key exchange is used to generate the root key for the next message, which can then be used to establish corresponding symmetric chain keys. Message keys are then derived from the current (newest) chain key, where chain keys are updated deterministically such that multiple messages can be sent in a row before a response, and no matter which of these messages is the first to arrive, the recipient can always compute its corresponding message key *immediately*. We now introduce the concept of asynchronous epochs before describing the two ratchets and the primary properties which they achieve:

Asynchronous Sending Epochs. In the double ratchet, the parties P_1 and P_2 asynchronously alternate sending messages in *epochs* (as termed in [ACD19]): Assume that P_1 starts the conversation, sending in epoch 1 at least one message. Then once P_2 receives one of these messages, she sends messages in epoch 2. Furthermore, once P_1 receives one of these message, she starts epoch 3, and so on. We emphasize that these sending epochs are *asynchronous* – for example, even if P_2 has started sending in epoch 2, if P_1 has not yet received any such epoch 2 messages and wants to send new messages, she will still send them in epoch 1. Not until she finally receives one of P_2 's epoch 2 messages will she send new messages in epoch 3.

Public Ratchet. The public ratchet forms the backbone of the double ratchet protocol. Parties update the root key using public-key cryptography (i.e. Diffie-Hellman secrets) every time a new epoch is initiated: if P_1 wishes to start a new epoch, she must first update the root key using the Diffie-Hellman public value from P_2 's latest epoch (or initialization). After deriving a new chain key from the root key, P_1 can send multiple separate messages in a row—this involves deriving a new message key for each message via the symmetric ratchet, as explained below.

We now describe the root key update process in more detail. To start a new epoch, P_1 samples a new private exponent y and corresponding public value g^y . Next, she uses the public value received from P_2 's latest epoch (or initialization), say g^x , to compute a shared secret $(g^x)^y = g^{xy}$. Then, P_1 uses a two-input Key Derivation Function (KDF) to update the current root key and derive a new chain key in one go. That is, she computes $(\sigma_{i+1}, w_{i+1}) \leftarrow \text{KDF}(\sigma_i, g^{xy})$.

P_1 includes in every message of the new epoch the fresh public share g^y to allow P_2 to compute the new shared secret g^{xy} that is used to update the root key, no matter which message of the epoch she receives first. So when P_2 receives a message in P_1 's new epoch, she recomputes the same above steps, i.e. she computes σ_{i+1} by first computing $(g^y)^x = g^{xy}$ where x is P_2 's own private share, followed by the same KDF computation. Once P_2 wishes to start her own new epoch, she generates another Diffie-Hellman pair (z, g^z) to ratchet the root key forward $(\sigma_{i+2}, w_{i+2}) \leftarrow \text{KDF}(\sigma_{i+1}, g^{yz})$. Essentially, P_2 has refreshed her component of the Diffie-Hellman shared secret while reusing P_1 's value from the previous epoch. Now, when P_1 receives a message for this epoch and again wishes to start a new one, she would similarly need to refresh her previous Diffie-Hellman share y . This process can continue asynchronously for as long as the session is active.

Symmetric Ratchet. The main purpose of the symmetric ratchet is to produce single-use symmetric keys for message encryption. When a party wishes to send (or receive) a new message, they derive a distinct message key K_i from the symmetric chain key w_i and simultaneously update the chain key. This is done by applying a KDF as follows: $(w_{i+1}, K_i) \leftarrow \text{KDF}(w_i)$ (if the KDF requires two inputs, a fixed value may be used to fill the other input). So, if P_1 just started a new epoch then she first computes initial symmetric chain key w_1 for the epoch as above. To derive a message key, P_1 puts this new chain key through the KDF to compute $(w_2, K_1) \leftarrow \text{KDF}(w_1)$. If P_1 wishes to send a second message, then she can derive $(w_3, K_2) \leftarrow \text{KDF}(w_2)$. When P_2 receives these messages from P_1 , she can repeat the key derivation in the same way as P_1 and use the subsequent message keys to decrypt the messages.

Post-Compromise Security and Forward Secrecy. Observe that even in the face of state leakages, as long as the parties have good sources of randomness, every time the root key is updated, fresh entropy is introduced via the use of Diffie-Hellman shares—this is the key to achieving Post-

Compromise Security (PCS). Symmetrically, even if parties use bad randomness when starting a new epoch, if the old root key was secure, then the new root key and chain key will still be secure as well. Furthermore, root keys are clearly forward secret, since the KDF is assumed to be one-way and new Diffie-Hellman secrets are sampled independently of past ones. Finally, the symmetric ratchet is also forward secure because the KDF is a one-way function. By updating the symmetric ratchet with every message, forward secrecy is provided at message-level granularity. Note however that the symmetric ratchet does not have PCS because of its deterministic nature.

Message Loss Resilience and Immediate Decryption. As explained above, whenever a new epoch is started by P_1 , no matter which message P_2 receives first for this epoch, she can update the root key appropriately since the public Diffie-Hellman share will be included in the message. Furthermore, since chain keys are computed deterministically using the KDF, P_2 can immediately derive the appropriate message key for decryption of this message from the derived root key. This property is called *immediate decryption*. Immediate decryption also implies a similarly property called *message loss resilience*, meaning that parties can always continue communication, even if some message is permanently lost by the network.

1.3 Technical Highlights of our Results

Given the above description we now provide the technical overview for our results here. Primarily, we first introduce ideal functionalities $\mathcal{F}_{\text{Signal}}$ and (stronger) $\mathcal{F}_{\text{Signal}^+}$ which provide more fine-grained guarantees than those previously in the literature. Indeed, we show several ways in which ACD’s security definition provides weaker guarantees than even $\mathcal{F}_{\text{Signal}}$ and thus allows for weaker protocols than **Signal**. Then, we show that **Signal** UC-realizes $\mathcal{F}_{\text{Signal}}$, and moreover (as in [ACD19]) show that swapping its building blocks with abstracted ones (CKA and FS-AEAD) that achieve the same level of security results in a generic SM protocol that UC-realizes $\mathcal{F}_{\text{Signal}}$. Finally, we make a minimalistic (and efficient) modification to **Signal** which allows it to achieve stronger functionality $\mathcal{F}_{\text{Signal}^+}$. As above, we also show changed generic requirements in one of the building blocks (CKA) which allow for this stronger security.

1.3.1 Shortcoming of ACD’s Definition.

We demonstrate the shortcoming of ACD’s definition by providing *four* distinct transformations to the original **Signal** protocol and showing their natural vulnerabilities here. In Section D.2 we formally show their insecurity with respect to our weaker functionality $\mathcal{F}_{\text{Signal}}$ and also their security with respect to ACD’s definition. Below we use the formalization of symmetric and public ratchets as done by ACD and also adapted by us – the symmetric ratchet is abstracted out as an *FS-AEAD* scheme and the public ratchet as a *CKA* scheme. We defer these definitions to Section B.

T_1 : Postponed FS-AEAD Key Deletion: This transformation slightly modifies the handling of symmetric ratchet secrets. In particular, when a party receives a new message for its counterpart’s next epoch, it does not immediately delete its (no longer needed) symmetric ratchet secrets from its previous sending epoch. Instead, it waits to delete these secrets until it starts its next sending epoch (i.e., sends its next message). In that case, an injection attack can be launched as follows: only focusing on the symmetric ratchet, suppose that for a sending epoch, P_1 derives $(w_2, K_1) \leftarrow \text{KDF}(w_1)$ and sends an encrypted message using K_1 , that is then received by P_2 . Then P_2 sends

a message, which is received by P_1 . Observe that unlike in (the strongest version of) **Signal**, this transformed protocol keeps w_2 in memory even after receiving this message from P_2 . Now if P_1 is compromised then the attacker obtains w_2 . Using this she can now launch an injection attack by encrypting any arbitrary message of her choice using the next message key $(\cdot, K_2) \leftarrow \text{KDF}(w_2)$ and sending that to P_2 . Notably, each time a sending epoch is started in the protocol, the information about how many messages were sent in the immediately past sending epoch is included. Nonetheless, that does not thwart this attack, because it is launched even before the next sending epoch starts at P_1 's end. Remarkably, both the primary description of the Double Ratchet in its white paper [MP16a] and ACD's description of **Signal** consider this version (however, the white paper does suggest to use the version which immediately deletes the symmetric ratchet secrets upon reception, for better security). Moreover, as evident by ACD's security proof, their definition does not require resistance against this attack; intuitively making their definition weaker than ours in this aspect.

T_2 : Postponed CKA Key Deletion: A similar problem arises if the keys from the public ratchet are kept for too long. The transformed protocol works as follows: suppose that in starting a new epoch, P_1 samples a secret exponent x_1 and combines that with the previous public value g^{x_2} from P_2 to compute $I = g^{x_1 x_2}$. At this point, instead of deleting I (as done in the actual **Signal** protocol), P_1 saves it. Now, it sends messages to P_2 encrypted with the message key derived from I and also sends g^{x_1} along. P_2 , who holds the secret x_2 can also compute I and decrypt the messages. Then, when P_1 again switches to a new sending epoch she generates a new I (deleting the old one). An attack can be executed on this protocol, by simply corrupting P_1 before the start of this newest sending epoch, and then using the leaked I to decrypt the earlier message sent by P_1 – thus breaking forward security. Note: this also requires another corruption of P_2 before the reception of the first epochs' messages to obtain the root key for the KDF. ACD's definition prevents such corruption of P_1 explicitly and thus does not require resistance against this attack.

T_3 : Eager CKA Randomness Sampling: If the secret-exponent of a public-ratchet is sampled too early, then that makes the protocol vulnerable. For example, consider the scenario when P_1 samples the exponent x for the next sending epoch when still in a receiving epoch. An attacker may compromise P_1 to obtain x (and the root key) at this stage and use that to decrypt the messages sent in the next epoch, thereby breaking PCS. ACD's game-based definition does not require resistance against this attack because it does not allow corruption right before the "challenge epoch". It is worth pointing out that as above, the Double Ratchet whitepaper [MP16a] presents this early sampling in its primary description (but again, later suggests deferring randomness sampling until actually sending, for better security). In this work we formally show the extra security that results from deferring randomness sampling.

T_4 : Malleable CKA: Finally we show that, if the public ratchet does not provide a strong non-malleability guarantee, then the **Signal** protocol could suffer from a malleability attack according to our weaker definition $\mathcal{F}_{\text{Signal}}$. More specifically, if the root key is leaked, without a strong enough key exchange mechanism to update the public ratchet, then there may exist attackers which, for example, can successfully maul **Signal** ciphertexts encrypting m into new ones that decrypt to $m+1$. This becomes evident when we prove the **Signal** protocol secure in our framework, as we need to rely on such a non-malleability property. ACD's definition however does not require resistance against

such an attack since injections are not allowed after corruptions. Thus, ACD’s security proof does not require such non-malleability guarantees.

1.3.2 Exact Security Assumptions for Signal.

Our security proofs assume that:

- The underlying CKA provides certain non-malleability guarantees. For example, in the specific (actual) implementation of Signal, we require that CDH is secure in a group even when given access to a DDH oracle (also known as gap-DH).
- The KDF function is modeled as a random oracle.

Both these assumptions are not only sufficient, but possibly also *necessary*. We do not formalize this, but provide some intuition. The first assumption is needed to ensure the aforementioned non-malleability of the underlying public ratchet (or CKA scheme). As shown above (attack on transformation T_4), the non-malleability property seems crucial to realize our definition; however, ACD’s security proof does not require such non-malleability properties.

The necessity of the random oracle arises for two reasons. First, to achieve non-malleability guarantees with such fantastic efficiency, Signal requires the random oracle (as in CCA-security for hashed ElGamal [ABR01, CS03, KM04]). Also, the fact that we do not restrict the adversary’s ability to leak on parties *at all* in our simulation-based definition requires the random oracle. This is in contrast with ACD’s game-based approach which does not allow the adversary to leak on parties when challenged-messages are in-transit. This is similar to the necessity of (programmable) random oracles for non-committing encryptions, as shown by a separation by Nielsen [Nie02] as opposed to just semantically secure encryption. Intuitively, to simulate a corruption that happens after “committing” the ciphertext one must use the random oracle programmability at its full extent – this common simulation paradigm also comes up naturally in our security analysis.

1.3.3 Shortcoming of Signal with Respect to $\mathcal{F}_{\text{Signal}^+}$.

We show that the Signal Double Ratchet protocol fails to realize our stronger functionality $\mathcal{F}_{\text{Signal}^+}$. We demonstrate this by providing an attack against Signal. The attack stems from the fact that a party needs to hold on to the secret exponent for the public ratchet that it generated in her newest sending epoch until it receives a message in her counterpart’s next sending epoch—this is important because the corresponding public component will be used by the other party to encrypt messages in the next epoch. For example, consider a setting in that party P_1 is about to start a sending epoch – at this point P_1 ’s state has g^{y_0} and P_2 has y_0 . Now when the sending epoch commences, P_1 samples a fresh secret (random) exponent x_1 and combines that with g^{y_0} to derive the CKA key $I_{1,0} = g^{x_1 y_0}$ which it then combines with the root key σ_1 to derive first the symmetric chain key w_1 , followed by message key K_1 . In this epoch P_1 sends g^{x_1} to P_2 , who then derives the same key $I_{1,0}$ by computing $(g^{x_1})^{y_0}$, and subsequently K_1 . In the next epoch, P_2 becomes a sender. Then P_2 samples a fresh y_1 to derive a new CKA key $I_{1,1} = (g^{x_1})^{y_1}$ and sends g^{y_1} to P_1 . Now, P_1 needs to compute $I_{1,1}$ as $(g^{y_1})^{x_1}$. To execute this step P_1 must have stored x_1 throughout her sending epoch. The attack exploits this by compromising P_1 twice in a short interval:

- first compromise P_1 before starting the sending epoch to obtain the root key σ_1 ;

- then compromise P_1 at any time after it sends a few messages (at least one) to obtain x_1 ;

and then combine σ_1 and x_1 to derive K_1 , given which *all* messages within P_1 's sending epoch are vulnerable including the ones that were sent between two corruptions. Intuitively, this breaks PCS with respect to the first corruption, as well as FS with respect to the second corruption. For more details we refer to Section C.2.

1.3.4 Modification to Achieve $\mathcal{F}_{\text{Signal}^+}$: the Triple Ratchet protocol

We provide a minimalistic modification of Signal, which we call the Triple Ratchet protocol, or simply Signal^+ , with virtually no overhead over Signal. This protocol realizes our stronger ideal functionality, $\mathcal{F}_{\text{Signal}^+}$. The Signal^+ protocol modifies the underlying public ratchet in a way that the sampled secret exponent is deterministically updated after starting a sending epoch; thus, adding a “mini ratchet” on top of Signal’s public ratchet. In particular, in the modified public ratchet, a party (say P_1) after sampling x_1 , and deriving $I_{1,0} = (g_0^y)^{x_1}$, stores $x'_1 = x_1 \cdot \mathbf{H}(I)$ instead of x_1 . Once P_2 receives g^{x_1} , it also derives I and computes $g^{x'_1} = g^{x_1 \cdot \mathbf{H}(I)}$ that she uses as the next public ratchet. In particular, in the next epoch when P_2 becomes the sender, she samples a fresh y_1 , and uses the key $I_{1,1} = g^{x'_1 y_1}$ where $x'_1 = x_1 \cdot \mathbf{H}(I)$. P_2 sends g^{y_1} , on receiving which P_1 can compute $I_{1,1}$ as $(g^{y_1})^{x'_1}$ and so on. Assuming \mathbf{H} to be a random oracle, or instead, circular-security of ElGamal encryption, we can show that given x'_1 , $I_{1,0}$ is completely hidden, rendering the above attack useless. Note that the communication cost remains the same for the modified protocol, that is one group element. The computation cost increases only slightly, specifically exactly once per epoch. We also note that, the alternate CKA scheme based on generic KEMs proposed by [ACD19] seems to achieve this security too, albeit with doubling the communication cost.

1.4 Related Work

The Off-The-Record (OTR) protocol introduced the concept of *ratcheting* [BGB04]. This scheme uses DH public keys to seed the generation of symmetric keys for message encryption, giving rise to a “double ratchet” structure similar to that of Signal. However, the simpler OTR protocol was not designed to work in an asynchronous environment: early versions of the protocol could not tolerate out-of-order messages and thus could not achieve immediate decryption. Signal’s double ratchet algorithm extended OTR into a robust, fully asynchronous messaging protocol.

Following the first formal analysis of Signal by Cohn-Gordon et al. [CCD⁺20], researchers proposed a number of protocols that provided stronger security than the Signal protocol [BSJ⁺17, PR18, JS18, DV19, JMM19a]. In the process of strengthening security, Alwen et al. [ACD19] observed that all such protocols suffer from steep efficiency costs and loss of *immediate decryption*, rendering these protocols impractical for real-world use.

In 2017, Bellare et al. [BSJ⁺17] studied the unidirectional ratcheting case. In this model, the receiver is never corrupted and never needs to update its state. These properties, in addition to sacrificing immediate decryption, give way to a simple protocol: fix the receiver’s public key as g^y . For any message m_i that the sender wishes to send, the sender performs DH key exchange using a fresh public key g^{x_i} and the receiver’s fixed key. This key is then used to encrypt m_i . Unfortunately, Bellare et al.’s protocol does not naturally extend to the bidirectional case, which is necessary for real-world messaging.

Jaeger and Stepanovs [JS18] and Poettering and Rösler [PR18] addressed the bidirectional case in 2018. Both papers model a bidirectional channel satisfying FS and PCS without considering

immediate decryption. The papers explore a stronger “fine-grained” PCS notion that guarantees partial security even during the healing period after a compromise. In summary, outgoing messages must still be private, and incoming messages must still be authentic; Signal’s double ratchet does not achieve this security. This additional security requires the use of expensive public-key primitives constructed from hierarchical identity-based encryption (HIBE) [GS02]. Recent work by Balli et al. [BRV20] showed that HIBE-like primitives are necessary to realize strong ratcheting security. Additionally, user state is growing and potentially unbounded in such schemes. These issues prevent stronger protocols from real-world deployment.

In 2019, Durak and Vaudenay [DV19] and Jost, Maurer and Mularczyk [JMM19a] introduced new notions of ratcheting security and presented new protocols. Like the works discussed above, the protocols therein do not tolerate out-of-order messages. Moreover, they make heavy use of standard public-key operations. For these reasons, they suffer from the same core issues as the schemes with fine-grained security, albeit to a lesser extent. Additionally, Jost, Maurer and Mularczyk [JMM19b] analyzed ratcheting with the Constructive Cryptography framework [Mau11]. They aimed to capture the security of various sub-protocols used in the construction of larger ratcheting protocols. An example of a sub-protocol is FS-AEAD, used by ACD to construct the Signal protocol. In contrast with our work, Jost, Maurer and Mularczyk develop a framework for analyzing ratcheting components with composability guarantees, whereas we focus on a holistic security definition for the full ratcheting protocol.

In the same year, Alwen, Coretti and Dodis (ACD) [ACD19] took a modular approach in analyzing the Signal double ratchet. ACD were the first to present a clear approximation of the level of security achieved by Signal. By breaking the protocol down into its core components, they were able to define security in a cleaner way than all other bidirectional ratcheting works. ACD were also the first to point out the impracticality of schemes from academia due to their lack of immediate decryption. Their modular protocol allowed for post-quantum secure ratcheting by instantiating the CKA module with a post-quantum secure KEM.

More recently, there has been work on the X3DH key exchange protocol used in Signal. Following Cohn-Gordon et al.’s initial analysis [CCD+20], Brendel et al. [BFG+20] gave a generalized framework for constructing X3DH key exchange protocols. Their formalism allows for a post-quantum secure version of X3DH, though the framework only supports DH-style assumptions. Hashimoto et al. [HKKP21] present a new class of authenticated key exchange (AKE) protocols, denoted *Signal-conforming* AKE protocols, and cast X3DH in this framework. They give a post-quantum Signal-conforming AKE from well-established assumptions. Many recent works have studied the offline deniability property believed to be present in X3DH [VGIK20, UG15, UG18, HKKP21].

1.5 Roadmap for the Rest of the Paper

In the next section we provide our UC-based ideal functionalities in Figure 2. We put a lot of discussions around it for reader’s convenience.

Due to space limitation we defer the rest of the technical sections to the supplementary body. Here we provide a brief summary of the content of the supplementary body and a roadmap to that.

In Appendix A we provide the preliminaries containing mostly definitions borrowed from literature.

In Appendix B we provide the building blocks required for Signal (and Signal⁺), such as CKA (capturing public ratchet) and FS-AEAD (formalizing the private ratchet part). While the concepts were borrowed from ACD [ACD19], our definitions differ significantly from theirs. In particular, we

propose a stronger CKA definition, called CKA+, which is achieved by our modified CKA protocol (this formalizes our Triple Ratchet protocol). Within this section we first describe our key-derivation chains and the differences from ACD (c.f. Section B.1) Then in the next part (Section B.2) we give our FS-AEAD definition (Definition 3, Definition 4 and Figure 4) – this differs from ACD significantly, which is also discussed in this section. The instantiation of FS-AEAD scheme, which is essentially the same as that of ACD, is provided in Figure 5. The security analysis with respect to our definition is formalized in Theorem 1. In the next subsection (Section B.3) we provide technical details of our CKA+ notion. The definitions for CKA and CKA+ are found in Definition 6, Definition 7, Definition 8 and Figure 6. The two definitions are merged into a security game (Figure 6). The definition significantly differs from ACD in terms of fine-grained security guarantees and the differences are discussed within this subsection. The CKA instantiation is provided next – this is the same as ACD’s instantiation and we prove that it satisfies security according to our (weaker) CKA definition in Theorem 2. Next we provide our new scheme that satisfies the CKA+ definition (Theorem 3) – looking ahead this scheme forms the core of our Triple Ratchet protocol. At the end of this section we also discuss how our new scheme actually provides even stronger security guarantees beyond what is captured by our CKA+ definition.

In Appendix C we provide the technical details of protocols Signal (Double Ratchet) and Signal+ (Triple Ratchet). In Subsection C.1 we provide our constructions, which are merged into one figure (Figure 7). Finally in Subsection C.2 we formally show the vulnerability of Signal with respect to our stronger functionality $\mathcal{F}_{\text{Signal}^+}$.

In Appendix D we provide the full technical details of our transformations to Signal, their insecurity with respect to our functionality $\mathcal{F}_{\text{Signal}}$ and their security with respect to ACD’s notion. To formalize, we require to provide ACD’s formal definitions – those are given in Subsection D.1. Specifically we provide ACD’s secure-messaging definition in Definition 11 and Figure 8. Their CKA definition is provided in Definition 13 and Figure 9; ACD’s FS-AEAD definition is given in Definition 14 and Figure 10. We also provide (a simplified version of) ACD’s composition theorem in Theorem 4. In Subsection D.2 we provide the full details of the transformation and their (in)security.

In Appendix E we provide the security analyses of the Double Ratchet Signal and Triple Ratchet Signal+. It is formalized in Theorem 5. The simulator is described in Figure 13 and the corresponding hybrids are detailed in the next Subsection E.1. The rest of this section contains the full technical details of the analysis.

Appendix F contains technical descriptions of our model the UC framework, mostly borrowed from the literature, but adapted to our setting.

2 Defining Security of the Signal Protocol

In this section, we focus on obtaining an ideal functionality $\mathcal{F}_{\text{Signal}}$ that captures the *exact* security provided by the Signal protocol. We emphasize that we study the security provided by the *strongest* implementation of Signal of which we are aware. For more on this, see Section D.2.1. We also provide an ideal functionality $\mathcal{F}_{\text{Signal}^+}$ that captures the security of our stronger Signal+ protocol. Both functionalities are provided in Figure 2.

$$\mathcal{F}_{\text{Signal}} \text{ and } \boxed{\mathcal{F}_{\text{Signal}^+}}$$

Notation: The ideal functionality interacts with two parties P_1, P_2 , and an ideal adversary \mathcal{S} . The ideal functionality initializes lists of *used message-ids* $P_1.M$, *in-transit* messages $P_1.T$, *adversarially injected* message-ids $P_1.I$,

and *vulnerable messages* $P_1.V$ sent by P_1 to P_2 to ϕ . Analogously, lists $P_2.M, P_2.T, P_2.I$ and $P_2.V$ are also initialized to ϕ . The ideal functionality also initializes leakage flags of both P_1 and P_2 for their corresponding (i) public ratchet secrets: $P_1.PLEK, P_2.PLEK$, (ii) current sending epoch symmetric secrets: $P_1.CUR_SLEK, P_2.CUR_SLEK$, and (iii) previous sending epoch symmetric secrets: $P_1.PREV_SLEK, P_2.PREV_SLEK$, all to 0. Further, it initializes bad-randomness flags $P_1.BAD, P_2.BAD$ and takeover possible flags $P_1.TAKEOVER_POSS, P_2.TAKEOVER_POSS$ to 0. Finally, it initializes the turn flag $TURN$ as \perp .

- **On input** (sid, **SETUP**) **from** P where $P \in \{P_1, P_2\}$: Send (sid, **SETUP**, P) to \mathcal{S} . When \mathcal{S} returns (sid, **SETUP**) then set $TURN \leftarrow P$, and send (sid, **INITIATED**) to both P_1 and P_2 . Ignore all future messages until this step is completed for sid. Once this happens P can send the first message.
 - **On input** (sid, mid, **SEND**, m) **from** $P \in \{P_1, P_2\}$:
 1. Ignore if $mid \in P.M$.
 2. If $\bar{P}.CUR_SLEK \vee (P.V \neq \emptyset)$ then $P.V \cup \{(sid, mid, m)\}$.
 3. If $\text{New}(P, TURN, P.T)^a$ then set (i) $P.PLEK \leftarrow P.BAD$, (ii) $P.CUR_SLEK \leftarrow \bar{P}.CUR_SLEK \wedge (P.PLEK \vee \bar{P}.PLEK)$, and (iii) $\bar{P}.TAKEOVER_POSS \leftarrow P.CUR_SLEK$.
 4. Add mid to $P.M$; if $mid \notin P.I$ then add (sid, mid, **IN_TRANSIT**, $m, P.CUR_SLEK, TURN$) to $P.T$; and pass (sid, mid, **IN_TRANSIT**, $P, |m|, m'$) to \mathcal{S} where $m' \leftarrow m$ if $P.CUR_SLEK$ and \perp otherwise.
 - **On input** (sid, mid, **DELIVER**, P, m') **from** \mathcal{S} where $P \in \{P_1, P_2\}$:
 1. Find (sid, mid, **IN_TRANSIT**, m, β, γ) $\in P.T$ and remove it from $P.T$. Skip rest of the steps if no such entry is found.
 2. If $\gamma = P$ then set (i) $TURN \leftarrow \bar{P}$, (ii) $P.T \leftarrow \text{Flip}(P, P.T)^b$, (iii) $P.PREV_SLEK \leftarrow 0$, (iv) $\bar{P}.PREV_SLEK \leftarrow \bar{P}.CUR_SLEK$, (v) $\bar{P}.CUR_SLEK \leftarrow 0$, (vi) $\bar{P}.PLEK \leftarrow 0$, (vii) $P.TAKEOVER_POSS \leftarrow 0$, and (viii) $\bar{P}.V \leftarrow \emptyset$.
 3. If $\beta = 1$ then set $m \leftarrow m'$. Send (sid, mid, **DELIVER**, m) to \bar{P} .
-
- **On input** (sid, **LEAK**, P) **from** \mathcal{S} where $P \in \{P_1, P_2\}$:
 1. If $\neg \text{New}(P, TURN, P.T)$ then set $P.CUR_SLEK \leftarrow 1$, $P.PLEK \leftarrow 1$, and $\bar{P}.TAKEOVER_POSS \leftarrow 1$.
 2. If $\neg \text{New}(\bar{P}, TURN, \bar{P}.T)$ then set $\bar{P}.CUR_SLEK \leftarrow 1$.
 3. If $TURN = \bar{P}$ then set $\bar{P}.PREV_SLEK \leftarrow 1$.
 4. If $\text{New}(P, TURN, P.T) \vee (\neg \text{New}(\bar{P}, TURN, \bar{P}.T) \wedge TURN = \bar{P})$ then set $P.TAKEOVER_POSS \leftarrow 1$.
 5. Execute $\bar{P}.T \leftarrow \text{Unsafe}(\bar{P}.T)^c$ and $P.T \leftarrow \text{Unsafe}'(P.T, P.V)^d$ then send $\bar{P}.T$ and $P.V$ to \mathcal{S} .
 - **On input** (sid, **BAD_RANDOMNESS**, P, ρ) **from** \mathcal{S} where $\rho \in \{0, 1\}$ and $P \in \{P_1, P_2\}$: Set $P.BAD \leftarrow \rho$.
 - **On input** (sid, mid, **INJECT**, P, m, δ, γ) **from** \mathcal{S} :
 1. Skip if $(mid \in P.I \cup P.M) \vee \neg (P.TAKEOVER_POSS \vee P.PREV_SLEK \vee P.CUR_SLEK)$.
 2. If $P.TAKEOVER_POSS \wedge \delta$ then forward all subsequent incoming messages from \bar{P} to \mathcal{S} and from \mathcal{S} for \bar{P} directly to \bar{P} . Also, remove from $P.T$ all elements of the form $(\cdot, \cdot, \text{IN_TRANSIT}, \cdot, \cdot, P)$ and drop all subsequent incoming messages for \bar{P} generated by the ideal functionality (i.e., do not send them to \bar{P}), except the ones generated according to the **DELIVER** command.
 3. Otherwise, add mid to $P.I$ and add to $P.T$ (i) if $TURN = \bar{P} \vee \neg P.CUR_SLEK$ then (sid, mid, **IN_TRANSIT**, $m, 1, \perp$), (ii) if $TURN = P \wedge \neg P.PREV_SLEK$ then (sid, mid, **IN_TRANSIT**, $m, 1, P$), and (iii) else (sid, mid, **IN_TRANSIT**, $m, 1, \gamma$).

^a $\text{New}(\mathbf{P}, \text{TURN}, \mathbf{P.T})$ outputs 1 if we have $\text{TURN} = \mathbf{P}$ and for all $(\text{sid}, \text{mid}, \text{IN_TRANSIT}, m, \beta, \gamma) \in \mathbf{P.T}$ we have that $\gamma \neq \mathbf{P}$; otherwise output 0.

^b $\text{Flip}(\mathbf{P}, \mathbf{P.T})$ for each $(\text{sid}, \text{mid}, \text{IN_TRANSIT}, m, \beta, \mathbf{P}) \in \mathbf{P.T}$ replaces it with $(\text{sid}, \text{mid}, \text{IN_TRANSIT}, m, \beta, \perp)$.

^c $\text{Unsafe}(\mathbf{P.T})$ for each $(\text{sid}, \text{mid}, \text{IN_TRANSIT}, m, \beta, \gamma) \in \mathbf{P.T}$ replaces it with $(\text{sid}, \text{mid}, \text{IN_TRANSIT}, m, 1, \gamma)$.

^d $\text{Unsafe}'(\mathbf{P.T}, \mathbf{P.V})$ for each $(\text{sid}, \text{mid}, m) \in \mathbf{P.V}$, if there is a corresponding $(\text{sid}, \text{mid}, \text{IN_TRANSIT}, m, \beta, \gamma) \in \mathbf{P.T}$, replaces it with $(\text{sid}, \text{mid}, \text{IN_TRANSIT}, m, 1, \gamma)$.

Figure 2: The ideal functionalities $\mathcal{F}_{\text{Signal}}$ and $\mathcal{F}_{\text{Signal}+}$, respectively.

2.1 Honest Execution

We start with a simplified view of the functionality where only the first three commands, namely **SETUP**, **SEND**, and **DELIVER** are executed. In other words, we consider a restricted view of the ideal functionality where leakage, bad randomness and injection attacks are not allowed. The adversary is still allowed to delay, reorder, and drop messages at will.

SETUP Command. This command can be initiated by either $\mathbf{P} = \mathbf{P}_1$ or $\mathbf{P} = \mathbf{P}_2$, and allows for initializing the communication channel between \mathbf{P} and $\bar{\mathbf{P}}$. Looking ahead, in the real-world protocol, this initialization will involve sharing cryptographic secrets between the real-world \mathbf{P} and real-world $\bar{\mathbf{P}}$, then properly initializing their states using these secrets. While the actual Signal protocol uses the X3DH key exchange [MP16b] for this, the focus of our work is to analyze the security and functionality of the double ratchet algorithm, and not X3DH. Therefore, we present a simple description for the **SETUP** command, that may be stronger than what X3DH achieves, but nonetheless suffices for our purposes.

We note that both \mathbf{P}_1 and \mathbf{P}_2 must receive $(\text{sid}, \text{INITIATED})$ before the communication between them can proceed. Turn status flag **TURN** is set to the initiator \mathbf{P} to denote that \mathbf{P} will be the first party to send a message.

SEND Command. This command allows $\mathbf{P} \in \{\mathbf{P}_1, \mathbf{P}_2\}$ to send a message m , under a unique assigned message id mid , to $\bar{\mathbf{P}}$. Naturally, the ideal functionality only allows \mathbf{P} to send one message under each such mid , which it ensures by aborting in Step 1 if mid is already in list $\mathbf{P.M}$, and subsequently adding mid to $\mathbf{P.M}$ in Step 4 otherwise. Now, this message might be dropped or delayed while in transit. Thus, at this point, the message is only added to the in-transit list $\mathbf{P.T}$ (Step 4) and the ideal-functionality waits for the instruction from the ideal-world adversary on when this message is to be delivered (if at all).

Observe that the last element of each tuple in $\mathbf{P.T}$ is **TURN**: the turn status at the time when \mathbf{P} attempted to send this message (i.e., when it was added to $\mathbf{P.T}$). Looking ahead, this element is used by the functionality in helper function $\text{New}(\mathbf{P}, \text{TURN}, \mathbf{P.T})$ within **SEND** (Step 3) to determine whether \mathbf{P} is initiating a new epoch when sending a message and, if so, the (in)security of the new epoch. When discussing the **DELIVER** command below, we will explain the role the last element of $\mathbf{P.T}$ plays in the logic of $\text{New}(\mathbf{P}, \text{TURN}, \mathbf{P.T})$ and further understand its role elsewhere in the functionality.

Finally, as is typical with encryption methods, in the real-world the length of the encrypted message is often leaked by the ciphertext. Thus, the ideal functionality leaks the length of the sent messages to the ideal adversary.

DELIVER Command. This command allows the ideal adversary to instruct the ideal functionality that a certain message, with unique message id mid , is no longer in-transit, and should be delivered to the recipient. The ideal functionality restricts the ideal adversary to delivering the message associated with this mid only once, which reflects the forward security of Signal – once a message is delivered, the recipient should no longer be able to decrypt it (in case she is leaked on afterwards). This is done by removing the entry for mid from P.T when it is delivered, so that subsequent deliveries cannot occur (Step 1).

As part of the delivery process (Step 2), the ideal functionality also checks if TURN was set to P when this message was sent. If so, the message was indeed the first of P 's newest epoch that is delivered to $\bar{\text{P}}$ (out of possibly many messages that can be the first delivered in the epoch). Thus, subsequently, it will next be $\bar{\text{P}}$'s turn to start a new epoch. So, if this is the case, then TURN is flipped to $\bar{\text{P}}$. Additionally, helper function $\text{Flip}(\text{P}, \text{P.T})$ flips the last entry of each message from P to $\bar{\text{P}}$ in P.T to \perp . This is done so that subsequently, when P starts its next sending epoch, $\text{New}(\text{P}, \text{TURN}, \text{P.T})$ will return 1: TURN will flip back to P once a message of $\bar{\text{P}}$'s next sending epoch is delivered to P for the first time, and there will be no element in P.T whose last entry is P . (Note: before P receives a message for $\bar{\text{P}}$'s next sending epoch, P 's sent messages will not initiate a new epoch, which is captured by $\text{New}(\text{P}, \text{TURN}, \text{P.T})$, since TURN will be set to $\bar{\text{P}}$.)

2.2 Execution with an Unrestricted Adversary

In addition to delaying, reordering, and dropping messages, we assume that the real-world adversary can: (i) provide bad randomness for both parties, (ii) leak the secret states of both parties; possibly multiple times at various points in the execution, (iii) tamper with in-transit messages between the parties, and (iv) attempt to inject messages on behalf of both parties. Here, we explain how the ideal functionality captures this behavior.

The Ideal Functionality's Flags. The ideal functionality uses several binary flags to properly capture adversarial behavior. The functionality initializes all of them to 0. Binary flag P.BAD captures bad randomness for party $\text{P} \in \{\text{P}_1, \text{P}_2\}$. Naturally, P.BAD is set to 0 or 1 when the ideal-world adversary issues a $(\text{sid}, \text{BAD_RANDOMNESS}, \text{P}, \rho)$ command to the ideal functionality, depending on the value of ρ . If P.BAD is set to 1 then P is provided with bad randomness (by the adversary, c.f. Appendix D.2.1) when she tries to sample some. Otherwise, P samples fresh randomness.

The ideal functionality further utilizes the following binary flags for each party $\text{P} \in \{\text{P}_1, \text{P}_2\}$ to capture the rest of the possible adversarial behavior. We first introduce their real-world semantic meaning here before explaining (i) their evolution within the ideal functionality as a result of the ideal world adversary's behavior, then (ii) how they thus allow the ideal functionality to determine security for the session.

- P.PLEK (Public Ratchet Secrets Leakage): If P.PLEK is set to 1 then P 's public ratchet secrets are leaked to the real-world adversary. Otherwise, they should be hidden from the real-world adversary.
- P.CUR_SLEK (Current Sending Symmetric Ratchet Secrets Leakage): If P.CUR_SLEK is set to 1 then the symmetric ratchet secrets of P 's *current* sending epoch are leaked to the real-world adversary. Otherwise, they should be hidden from the real-world adversary.

- **P.PREV_SLEK** (Previous Sending Symmetric Ratchet Secrets Leakage): If **P.PREV_SLEK** is set to 1 then the symmetric ratchet secrets of the *previous* sending epoch of P are leaked to the real-world adversary. Otherwise, they should be hidden from the real-world adversary.
- **P.TAKEOVER_POSS** (Takeover Possible): If **P.TAKEOVER_POSS** is set to 1 then the real-world adversary has the option to take over the role of P in the conversation with \bar{P} . Otherwise, the real-world adversary should not have this option.

How the Flags are Affected by Leakages. We first describe how a leakage on one of the parties $P \in \{P_1, P_2\}$ effects the above flags. For **P.PLEK**, when $\text{New}(P, \text{TURN}, P.T) = 1$, it is P 's turn to start her next sending epoch, but she has not yet started it. Thus she does not currently have any public ratchet secret state (just \bar{P} 's public value), so there is no effect on **P.PLEK** if leakage on P occurs in this case. If $\text{New}(P, \text{TURN}, P.T) = 0$ when leakage on P occurs, P of course does have a public ratchet secret state, as she needs to be able to receive a message for \bar{P} 's next sending turn; thus in command **LEAK**, the ideal functionality sets **P.PLEK** to 1 (Step 1). Since P never stores \bar{P} 's public ratchet secrets, there is never any effect on **\bar{P} .PLEK** when P 's state is leaked.

For **P.CUR_SLEK**, the functionality has similar behavior. If $\text{New}(P, \text{TURN}, P.T) = 1$ when leakage on P occurs, P has not yet generated the secrets for her next sending epoch, so **P.CUR_SLEK** is not modified. Otherwise, P has started the epoch, and so she stores the corresponding secrets in order to send new messages for the epoch; thus in command **LEAK**, we set **P.CUR_SLEK** to 1 (Step 1). Additionally, if $\text{New}(\bar{P}, \text{TURN}, \bar{P}.T) = 1$ then \bar{P} has not yet generated the secrets for her next sending epoch, so **\bar{P} .CUR_SLEK** is not modified. Otherwise, \bar{P} has indeed started the epoch, in which case P must be able to derive the epoch's symmetric secrets (possibly using in-transit messages, which the adversary has), and thus in command **LEAK** we set **\bar{P} .CUR_SLEK** to 1 (Step 2).

For **P.PREV_SLEK**, since in the most secure version of Signal (see Section D.2.1 for more on this), P only ever stores the secrets for her current sending epoch (if she has indeed started it), leakage on P has no effect on **P.PREV_SLEK**. However, once it is \bar{P} 's turn to start a new sending epoch, P still stores the secrets of \bar{P} 's previous sending epoch (in case she needs to receive messages for it; she does not yet know \bar{P} will never again send a message for that epoch), until she receives a message in \bar{P} 's new epoch for the first time. Therefore, if $\text{TURN} = \bar{P}$ then in command **LEAK**, we set **\bar{P} .PREV_SLEK** to 1 (Step 3); otherwise, if $\text{TURN} = P$, **\bar{P} .PREV_SLEK** is not modified.

Finally, for **P.TAKEOVER_POSS**, if it is P 's turn to start a new sending epoch, then of course a leakage on P will enable the adversary to forge the first message of this new epoch and thus influence the subsequent state of \bar{P} upon delivery such that the adversary can take over P 's role in the conversation (if it wishes). This is because the adversary will obtain the double ratchet root key, and can thus send such a message herself. Also, note that this key is derived from (i) P 's previous state before she received any message for \bar{P} 's newest epoch and (ii) any message of \bar{P} 's newest sending epoch. Thus, additionally, if P is leaked while any message from \bar{P} 's newest epoch is in-transit, but before P receives any such message, then the adversary can obtain the root key as above, and so will have the ability to forge the first message of P 's next sending epoch. Therefore, in command **LEAK**, if $\text{New}(P, \text{TURN}, P.T) = 1$, or $\text{New}(\bar{P}, \text{TURN}, \bar{P}.T) = 0$ and $\text{TURN} = \bar{P}$, then we set **P.TAKEOVER_POSS** to 1 (Step 4). Otherwise, if P has already sent the first message of the epoch, and \bar{P} has not yet started her next sending epoch, leakage on P does not reveal the root key, so **P.TAKEOVER_POSS** is not modified. In the former case, this is because P deletes the key after sending the message, and in the latter case, this is because the key does not yet exist. Furthermore, if P has indeed sent a message for her current sending epoch, then a leakage on P will provide the

adversary with the new root key. The adversary will therefore be able to forge the first message for \bar{P} 's next sending epoch. So, if $\text{New}(P, \text{TURN}, P.T) = 0$ then in command **LEAK**, we additionally set $\bar{P}.\text{TAKEOVER_POSS}$ to 1 (Step 1). Otherwise, if it is P 's turn to start a new epoch, and she has not yet started it, then the new root key has not yet been generated, so $\bar{P}.\text{TAKEOVER_POSS}$ is not modified.

How the Flags are Affected by Epoch Initialization. The effects on the ideal functionality's flags of epoch initialization via a **SEND** command are determined in Step 3 of the command. First, if $P.\text{BAD} = 1$ when starting a new epoch (i.e. P uses bad randomness to start it), then we of course set $P.\text{PLEK}$ to 1 (In Signal^+ we may still have security of P 's public ratchet secret state in this circumstance, but we choose to capture slightly weaker security for simplicity); otherwise we set $P.\text{PLEK}$ to 0. Now, consider the privacy of the root key when $\bar{P}.\text{CUR_SLEK}$ is 1 and P is initializing a new epoch:

- If $\bar{P}.\text{CUR_SLEK}$ was set to 1 when \bar{P} initialized her newest epoch (as we explain below), then the root key must have been leaked in addition to the corresponding symmetric ratchet secrets, since they are both part of the same KDF output.
- If $\bar{P}.\text{CUR_SLEK}$ was set to 1 as a result of a leakage on P , then the root key must have been also leaked, since P needs it to start her new sending epoch.
- Finally, if $\bar{P}.\text{CUR_SLEK}$ was set to 1 as a result of a leakage on \bar{P} , then the root key must have been also leaked, since \bar{P} needs it to receive a message for P 's new sending epoch.

So, if $\bar{P}.\text{CUR_SLEK}$ is 1 when P initializes her new sending epoch, then it must be that the root key is leaked. Thus, only if P and \bar{P} have a secure key exchange can security for Signal be recovered, which only happens if both $P.\text{PLEK}$ and $\bar{P}.\text{PLEK}$ are 0, i.e., their public ratchet secrets are both hidden from the adversary. In this case, we set $P.\text{CUR_SLEK}$ to 0; otherwise, we set it to 1. If $\bar{P}.\text{CUR_SLEK}$ is 0 at the time of initialization, then the root key must be hidden. This is because if not, then the current symmetric ratchet secrets of \bar{P} would also not be hidden, since they were part of the same KDF output when \bar{P} started her latest sending epoch, and there were no subsequent leakages on either party. So we set $P.\text{CUR_SLEK}$ to 0 upon initialization, in this case.

Finally, if we do indeed set $P.\text{CUR_SLEK}$ to 1 at this time, as we noted above, this means that the new root key is known by the adversary, and thus the adversary could forge the first message for \bar{P} 's next turn; otherwise the root key is hidden, and so the adversary does not have this ability. So, we set $\bar{P}.\text{TAKEOVER_POSS} \leftarrow P.\text{CUR_SLEK}$.

How the Flags are Affected by Epoch Termination. When the ideal adversary issues a **DELIVER** command for the first message of P 's newest sending epoch, the ideal functionality needs to properly evolve the flags it uses to capture adversarial behavior (Step 2). First, when such a delivery occurs, \bar{P} 's latest sending epoch terminates, as her next message will be sent in a new epoch. To reflect this, upon such a delivery, the ideal functionality sets $\bar{P}.\text{PREV_SLEK} \leftarrow \bar{P}.\text{CUR_SLEK}$. Also, since \bar{P} deletes her public ratchet secrets upon reception of such a message, and her newest epoch has not actually started at this point, the functionality sets $\bar{P}.\text{CUR_SLEK} \leftarrow 0$ and $\bar{P}.\text{PLEK} \leftarrow 0$.

Furthermore, in Signal , P includes in each message of an epoch the number of messages she sent in her previous epoch (see Section C). Thus, once \bar{P} receives such a message in Signal , she knows exactly how many messages P sent in her previous epoch. So, the adversary can no longer inject

messages in P’s previous epoch (just modify them) and there is no more adversarial action possible for that epoch, so the functionality sets $P.PREV_SLEK \leftarrow 0$. Finally, since a message for P’s newest sending epoch has indeed been delivered at this point, the secrets needed to start her next sending epoch are yet to be determined. Thus, the adversary cannot yet forge a message to start her next sending epoch, so the functionality sets $P.TAKEOVER_POSS \leftarrow 0$.

Determining New Messages’ Privacy and Authenticity. We know from above that if $P.CUR_SLEK = 1$, then P’s current symmetric ratchet secrets are leaked to the adversary. Thus, if P issues a **SEND** command for message m with id mid , and $P.CUR_SLEK = 1$, then the ideal functionality leaks the corresponding message to the ideal adversary (Step 4). Additionally, the ideal functionality sets the penultimate element of mid ’s entry in $P.T$ to 1. This will allow the ideal adversary to modify the message associated with mid upon delivery: the adversary will issue a **DELIVER** command for mid to the functionality with input modified message m' , which will then be delivered P, instead of m (Step 3).

Otherwise, if $P.CUR_SLEK = 0$ when P issues the **SEND** command, then the ideal functionality only leaks the message length to the adversary and sets the penultimate element of the corresponding entry of $P.T$ to 0, ensuring (for now) privacy and authenticity of m .

The Consequences of Leakages. When the adversary leaks on P in the real-world, the privacy of in-transit messages from \bar{P} to P is no longer guaranteed, since P must preserve all keys that will be necessary for authenticating and decrypting them. Therefore, when the ideal adversary issues a **LEAK** command on P, the ideal functionality leaks the in-transit messages from \bar{P} to P, $P.T$, to the ideal adversary, and allows the ideal adversary to modify them in the future (Step 5). The ideal functionality accomplishes the latter using helper function $Unsafe(\bar{P}.T)$ which sets the penultimate element of each in-transit message of $\bar{P}.T$ to 1. As a result, the ideal adversary can modify these in-transit messages in the **DELIVER** command, as described above.

Vulnerable Messages in *Signal*. As explained in Section D.2.1, if in *Signal*, the root key is leaked when it is P’s turn to start a new sending epoch, but she has not yet started it, then the messages of that epoch are *vulnerable*. This means that if P is leaked on before she receives a message of \bar{P} ’s next sending epoch for the next time, the messages that she sent in her epoch become insecure.

To capture this, the ideal functionality in the **SEND** command adds messages to list $P.V$ if they are indeed vulnerable (Step 2). At the start of the epoch, this is the case if $\bar{P}.CUR_SLEK = 1$ (as explained above); in the middle of the epoch, this is the case if $P.V$ is non-empty. Hence, if the adversary issues a **LEAK** command on P, in addition to the consequences of the above paragraph, the ideal functionality *also* leaks $P.V$ and allows for future modification of its elements that are still in-transit (Step 5). The latter is accomplished via helper function $Unsafe'(P.T, P.V)$, similarly as in $Unsafe(\bar{P}.T)$. Finally, if the adversary issues a **DELIVER** command for the first message of P’s next sending epoch, the ideal functionality sets $P.V = \emptyset$: P properly deletes the secrets which make those messages vulnerable at this time.

Injections and Takeovers. If $P.CUR_SLEK = 1$ or $P.PREV_SLEK = 1$, then the adversary has the secrets required to inject its own messages into P’s current or previous sending epoch, respectively. Also, if $P.TAKEOVER_POSS = 1$, then the adversary can forge the first message to be delivered in

P's next sending epoch to \bar{P} . In either case, the ideal adversary issues the **INJECT** command to inject message m under unique message id mid on behalf of P. Of course, the ideal functionality only allows the adversary to inject one message under each such mid , which it ensures by aborting in Step 1 if mid is already in **P.I**, and adding it to **P.I** in Step 3 if not. The ideal functionality also aborts if a message with message id mid was already sent by P, i.e., it is in **P.M**, in which case injection of mid is not allowed, only modification.

First, if **P.TAKEOVER.POSS** = 1, and the ideal adversary inputs $\delta = 1$ to the **INJECT** command, indicating that it wishes to takeover for P, then the ideal functionality thereafter directly forwards messages sent from P to the ideal adversary, and vice versa (Step 2).

If the ideal adversary injects a message with id mid that is not a takeover forgery, then before actual delivery of the injection occurs, a corresponding entry is added to **P.T**. However, the ideal functionality has to be careful to set the last element of this entry correctly:

- If **TURN** = \bar{P} , then the first message of P's current sending epoch has already been delivered to \bar{P} . Thus, the last element of the entry is set to \perp , so that if **TURN** is flipped to P, the entry's subsequent delivery does not prematurely flip **TURN** back. Moreover, if **P.CUR_SLEK** = 0, then the adversary must be injecting into P's previous sending epoch, so for the same reason as above, we set its last entry to \perp .
- If **P.PREV_SLEK** = 0 and **TURN** = P, then the adversary must be injecting into P's current sending epoch, and moreover, it might be that the injected message could be the first of the epoch delivered to \bar{P} . Therefore, we set **TURN** to P.
- If neither of the above are true, i.e., **TURN** = P, **P.PREV_SLEK** = 1, and **P.CUR_SLEK** = 1, then it could be that the adversary is injecting into *either* P's previous or current sending epoch. Therefore, the ideal adversary specifies its choice of the last element with the last input γ to the **INJECT** command.

Actual delivery of injections is then handled in the **DELIVER** command, in the same simple manner as specified in the Honest Execution Section (Section 2.1). Namely, delivery of injected message with message id mid is done by removing it from **P.T** (if such an entry exists), and sending it to \bar{P} . The functionality works this way in order to capture the scenario in which the real-world adversary modifies the first message of a new sending epoch for P to inform \bar{P} that P's last sending epoch contains more messages than it actually does. Therefore, the real-world adversary will be able to in the future inject such additional messages whenever it wants. The ideal-world adversary thus issues an **INJECT** command for all of these message ids at the time of the first modification, so that later it can actually send them to \bar{P} using **DELIVER** commands (regardless of the status of the functionality's flags at that time).

If an injected message with id mid is indeed added to **P.T**, then the ideal functionality needs to also make sure that P can send a message with the same mid (since it does not know about the injection), but not allow the ideal adversary to deliver two messages with the same mid (since \bar{P} will only accept one such message in Signal). Therefore, in the **SEND** command, the ideal functionality checks if $\text{mid} \notin \text{P.I}$ and if so adds the corresponding message to **P.T** as in the honest execution. However, if mid is in **P.I**, the ideal functionality does not add the corresponding message to **P.T**, but still sends the length of the message (and the message itself if **P.CUR_SLEK** = 1) to the ideal adversary, mirroring that a ciphertext is still created in the real-world.

References

- [ABR01] Michel Abdalla, Mihir Bellare, and Phillip Rogaway. The oracle Diffie-Hellman assumptions and an analysis of DHIES. In David Naccache, editor, *Topics in Cryptology – CT-RSA 2001*, volume 2020 of *Lecture Notes in Computer Science*, pages 143–158, San Francisco, CA, USA, April 8–12, 2001. Springer, Heidelberg, Germany.
- [ACD18] Joël Alwen, Sandro Coretti, and Yevgeniy Dodis. The double ratchet: Security notions, proofs, and modularization for the Signal protocol. Cryptology ePrint Archive, Report 2018/1037, 2018. <https://eprint.iacr.org/2018/1037>.
- [ACD19] Joël Alwen, Sandro Coretti, and Yevgeniy Dodis. The double ratchet: Security notions, proofs, and modularization for the Signal protocol. In Yuval Ishai and Vincent Rijmen, editors, *Advances in Cryptology – EUROCRYPT 2019, Part I*, volume 11476 of *Lecture Notes in Computer Science*, pages 129–158, Darmstadt, Germany, May 19–23, 2019. Springer, Heidelberg, Germany.
- [BCH12] Nir Bitansky, Ran Canetti, and Shai Halevi. Leakage-tolerant interactive protocols. In Ronald Cramer, editor, *TCC 2012: 9th Theory of Cryptography Conference*, volume 7194 of *Lecture Notes in Computer Science*, pages 266–284, Taormina, Sicily, Italy, March 19–21, 2012. Springer, Heidelberg, Germany.
- [BFG⁺20] Jacqueline Brendel, Marc Fischlin, Felix Günther, Christian Janson, and Douglas Stebila. Towards post-quantum security for signal’s x3dh handshake. In *Selected Areas in Cryptography–SAC 2020*, 2020.
- [BGB04] Nikita Borisov, Ian Goldberg, and Eric Brewer. Off-the-record communication, or, why not to use pgp. In *Proceedings of the 2004 ACM workshop on Privacy in the electronic society*, pages 77–84, 2004.
- [BRV20] Fatih Balli, Paul Rösler, and Serge Vaudenay. Determining the core primitive for optimally secure ratcheting. In Shiho Moriai and Huaxiong Wang, editors, *Advances in Cryptology – ASIACRYPT 2020, Part III*, volume 12493 of *Lecture Notes in Computer Science*, pages 621–650, Daejeon, South Korea, December 7–11, 2020. Springer, Heidelberg, Germany.
- [BSJ⁺17] Mihir Bellare, Asha Camper Singh, Joseph Jaeger, Maya Nyayapati, and Igors Stepanovs. Ratcheted encryption and key exchange: The security of messaging. In Jonathan Katz and Hovav Shacham, editors, *Advances in Cryptology – CRYPTO 2017, Part III*, volume 10403 of *Lecture Notes in Computer Science*, pages 619–650, Santa Barbara, CA, USA, August 20–24, 2017. Springer, Heidelberg, Germany.
- [Can01] Ran Canetti. Universally composable security: A new paradigm for cryptographic protocols. In *42nd Annual Symposium on Foundations of Computer Science*, pages 136–145, Las Vegas, NV, USA, October 14–17, 2001. IEEE Computer Society Press.
- [CCD⁺20] Katriel Cohn-Gordon, Cas Cremers, Benjamin Dowling, Luke Garratt, and Douglas Stebila. A formal security analysis of the signal messaging protocol. *J. Cryptol.*, 33(4):1914–1983, 2020.

- [CS03] Ronald Cramer and Victor Shoup. Design and analysis of practical public-key encryption schemes secure against adaptive chosen ciphertext attack. *SIAM Journal on Computing*, 33(1):167–226, November 2003.
- [DV19] F. Betül Durak and Serge Vaudenay. Bidirectional asynchronous ratcheted key agreement with linear complexity. In Nuttapong Attrapadung and Takeshi Yagi, editors, *IWSEC 19: 14th International Workshop on Security, Advances in Information and Computer Security*, volume 11689 of *Lecture Notes in Computer Science*, pages 343–362, Tokyo, Japan, August 28–30, 2019. Springer, Heidelberg, Germany.
- [FIP95] PUB FIPS. 180-1. secure hash standard. *National Institute of Standards and Technology*, 17:45, 1995.
- [GMR89] Shafi Goldwasser, Silvio Micali, and Charles Rackoff. The knowledge complexity of interactive proof systems. *SIAM J. Comput.*, 18(1):186–208, 1989.
- [Gol04] Oded Goldreich. *The Foundations of Cryptography - Volume 2: Basic Applications*. Cambridge University Press, 2004.
- [GS02] Craig Gentry and Alice Silverberg. Hierarchical ID-based cryptography. In Yuliang Zheng, editor, *Advances in Cryptology – ASIACRYPT 2002*, volume 2501 of *Lecture Notes in Computer Science*, pages 548–566, Queenstown, New Zealand, December 1–5, 2002. Springer, Heidelberg, Germany.
- [HKKP21] Keitaro Hashimoto, Shuichi Katsumata, Kris Kwiatkowski, and Thomas Prest. An efficient and generic construction for signal’s handshake (x3dh): Post-quantum, state leakage secure, and deniable. In *Public Key Cryptography (2)*, pages 410–440, 2021.
- [JMM19a] Daniel Jost, Ueli Maurer, and Marta Mularczyk. Efficient ratcheting: Almost-optimal guarantees for secure messaging. In Yuval Ishai and Vincent Rijmen, editors, *Advances in Cryptology – EUROCRYPT 2019, Part I*, volume 11476 of *Lecture Notes in Computer Science*, pages 159–188, Darmstadt, Germany, May 19–23, 2019. Springer, Heidelberg, Germany.
- [JMM19b] Daniel Jost, Ueli Maurer, and Marta Mularczyk. A unified and composable take on ratcheting. In Dennis Hofheinz and Alon Rosen, editors, *TCC 2019: 17th Theory of Cryptography Conference, Part II*, volume 11892 of *Lecture Notes in Computer Science*, pages 180–210, Nuremberg, Germany, December 1–5, 2019. Springer, Heidelberg, Germany.
- [JS18] Joseph Jaeger and Igors Stepanovs. Optimal channel security against fine-grained state compromise: The safety of messaging. In Hovav Shacham and Alexandra Boldyreva, editors, *Advances in Cryptology – CRYPTO 2018, Part I*, volume 10991 of *Lecture Notes in Computer Science*, pages 33–62, Santa Barbara, CA, USA, August 19–23, 2018. Springer, Heidelberg, Germany.
- [KE10] Hugo Krawczyk and Pasi Eronen. Hmac-based extract-and-expand key derivation function (hkdf). Technical report, RFC 5869, May, 2010.

- [KM04] Kaoru Kurosawa and Toshihiko Matsuo. How to remove MAC from DHIES. In Huaxiong Wang, Josef Pieprzyk, and Vijay Varadharajan, editors, *ACISP 04: 9th Australasian Conference on Information Security and Privacy*, volume 3108 of *Lecture Notes in Computer Science*, pages 236–247, Sydney, NSW, Australia, July 13–15, 2004. Springer, Heidelberg, Germany.
- [Mau11] Ueli Maurer. Constructive cryptography—a new paradigm for security definitions and proofs. In *Joint Workshop on Theory of Security and Applications*, pages 33–56. Springer, 2011.
- [MP16a] M. Marlinspike and T. Perrin. The Double Ratchet Algorithm, 11 2016. <https://whispersystems.org/docs/specifications/doubleratchet/doubleratchet.pdf>.
- [MP16b] M. Marlinspike and T. Perrin. The X3DH Key Agreement Protocol, 11 2016. <https://signal.org/docs/specifications/x3dh/x3dh.pdf>.
- [Nie02] Jesper Buus Nielsen. Separating random oracle proofs from complexity theoretic proofs: The non-committing encryption case. In Moti Yung, editor, *Advances in Cryptology – CRYPTO 2002*, volume 2442 of *Lecture Notes in Computer Science*, pages 111–126, Santa Barbara, CA, USA, August 18–22, 2002. Springer, Heidelberg, Germany.
- [OP01] Tatsuaki Okamoto and David Pointcheval. The gap-problems: A new class of problems for the security of cryptographic schemes. In Kwangjo Kim, editor, *PKC 2001: 4th International Workshop on Theory and Practice in Public Key Cryptography*, volume 1992 of *Lecture Notes in Computer Science*, pages 104–118, Cheju Island, South Korea, February 13–15, 2001. Springer, Heidelberg, Germany.
- [PR18] Bertram Poettering and Paul Rösler. Towards bidirectional ratcheted key exchange. In Hovav Shacham and Alexandra Boldyreva, editors, *Advances in Cryptology – CRYPTO 2018, Part I*, volume 10991 of *Lecture Notes in Computer Science*, pages 3–32, Santa Barbara, CA, USA, August 19–23, 2018. Springer, Heidelberg, Germany.
- [Sip97] Michael Sipser. *Introduction to the theory of computation*. PWS Publishing Company, 1997.
- [UG15] Nik Unger and Ian Goldberg. Deniable key exchanges for secure messaging. In Indrajit Ray, Ninghui Li, and Christopher Kruegel, editors, *ACM CCS 2015: 22nd Conference on Computer and Communications Security*, pages 1211–1223, Denver, CO, USA, October 12–16, 2015. ACM Press.
- [UG18] Nik Unger and Ian Goldberg. Improved strongly deniable authenticated key exchanges for secure messaging. *Proc. Priv. Enhancing Technol.*, 2018(1):21–66, 2018.
- [VGIK20] Nihal Vatandas, Rosario Gennaro, Bertrand Ithurburn, and Hugo Krawczyk. On the cryptographic deniability of the Signal protocol. In Mauro Conti, Jianying Zhou, Emiliano Casalicchio, and Angelo Spognardi, editors, *ACNS 20: 18th International Conference on Applied Cryptography and Network Security, Part II*, volume 12147 of *Lecture Notes in Computer Science*, pages 188–209, Rome, Italy, October 19–22, 2020. Springer, Heidelberg, Germany.

Supplementary Material

A Preliminaries

A.1 Game-Based Security and Notation

In addition to our use of the UC framework to capture the security and functionality of `Signal`, we also consider some game-based security definitions for the primitives that are used within `Signal`, i.e., games executed between a challenger and an adversary. The games have one of the following formats:

- **Unpredictability games:** First, the challenger executes the special `init` procedure, which sets up the game. Subsequently, the attacker is given access to a set of oracles that allow it to interact with the scheme in question. The goal of the adversary is to provoke a particular, game-specific *winning* condition. The *advantage* of an adversary \mathcal{A} against construction C in an unpredictability game Γ^C is

$$\mathcal{A}_\Gamma^C(\mathcal{A}) := \Pr[\mathcal{A} \text{ wins } \Gamma^C] .$$

- **Indistinguishability games:** In addition to setting up the game, the `init` procedure samples a secret bit $b \in \{0, 1\}$. The goal of the adversary is to determine the value of b . Once more, upon completion of `init`, the attacker interacts arbitrarily with all available oracles up to the point where it outputs a guess bit b' . The adversary *wins* the game if $b = b'$. The *advantage* of an adversary \mathcal{A} against construction C in an indistinguishability game Γ is

$$\mathcal{A}_\Gamma^C(\mathcal{A}) := 2 \cdot |\Pr[\mathcal{A} \text{ wins } \Gamma^C] - 1/2| .$$

- **Recoverability games:** In such games, the attacker is once more given access to a set of oracles that allow it to interact with the scheme in question after the initial `init` procedure. In this case, the goal of the adversary is to recover some secret value S (usually security parameter-many bits long) that is used within the scheme. The adversary *wins* the game if they guess that the secret value is $S' = S$. The *advantage* of an adversary \mathcal{A} against construction C in a recoverability game Γ is

$$\mathcal{A}_\Gamma^C(\mathcal{A}) := \Pr[\mathcal{A} \text{ wins } \Gamma^C] .$$

With the above in mind, to describe any security (or correctness) notion, one need only specify the `init` oracle and the oracles available to \mathcal{A} . The following special keywords are used to simplify the exposition of the security games:

- **req** is followed by a condition; if the condition is not satisfied, the oracle/procedure containing the keyword is exited and all actions by it are undone.
- **win** is used to declare that the attacker has won the game; it can be used for all types of games above.
- **end** disables all oracles and returns all values following it to the attacker.

Moreover, the descriptions of some games/schemes involve *dictionaries*. For ease of notation, these dictionaries are described with the *array-notation* described next, but it is important to note that they are to be implemented by a data structure whose size grows (linearly) with the number of elements *in* the dictionary (unlike arrays):

- **Initialization:** The statement $D[\cdot] \leftarrow \perp$ initializes an *empty* dictionary D .
- **Adding elements:** The statement $D[i] \leftarrow v$ adds a value v to dictionary D with key i , overriding the value previously stored with key i if necessary.
- **Retrieval:** The expression $D[i]$ returns the value v with key i in the dictionary; if there are no values with key i , the value \perp is returned.
- **Deletion:** The statement $D[i] \leftarrow \perp$ *deletes* the value v corresponding to key i .

Additionally, sometimes the random coins of certain probabilistic algorithms are made explicit. For example, $y \leftarrow A(x; r)$ means that A , on input x and with random tape r , produces output y . If r is not explicitly stated, it is assumed to be chosen uniformly at random; in this case, the notation $y \stackrel{\$}{\leftarrow} A(x)$ is used.

Finally, many of the protocols in this work, including **Signal** itself, may consist of algorithms which take in some party's state. In this case, some such algorithms, upon failing, may throw an exception (**error**), which causes the calling party's state to be rolled back to where it was before the algorithm was invoked.

A.2 Authenticated Encryption

Definition 1. An authenticated encryption with associated data (AEAD) scheme is a pair of algorithms $AE = (\text{Enc}, \text{Dec})$ with the following syntax:

- **Encryption:** Enc takes a key K , associated data a , and a message m and produces a ciphertext $e \leftarrow \text{Enc}(K, a, m)$.
- **Decryption:** Dec takes a key K , associated data a , and a ciphertext e and produces a message $m \leftarrow \text{Dec}(K, a, e)$.

All AEAD schemes in this paper are assumed to be deterministic, i.e., all randomness stems from the key K .

Correctness. An AEAD scheme is *correct* if for all keys K and all pairs (K, a) ,

$$\text{Dec}(K, a, \text{Enc}(K, a, m)) = m.$$

Security. In order to be used in the constructions in this paper, AEAD schemes need to satisfy one-time IND-CCA security. This is captured by the game depicted in Figure 3. It provides access to a one-time encryption oracle that on input associated data a and messages m_0, m_1 , returns an encryption of message m_b , depending on a randomly chosen bit b . Moreover, the attacker may query a decryption oracle arbitrarily many times (except on the challenge ciphertext), which, however, always returns \perp if $b = 1$.

Oracles for AEAD One-Time IND-CCA Game

init $K \leftarrow \mathcal{K}$ $e^* \leftarrow \perp$ $b \leftarrow \{0, 1\}$	encrypt (a, m_0, m_1) $e^* \leftarrow \text{Enc}(K, a, m_b)$ return e^*	decrypt (a, e) if $e = e^*$ or $b = 1$ \quad return \perp return $\text{Dec}(K, a, e)$
corr if $e^* \neq \perp$ and $b = 1$ \quad $\text{AEAD-Expl-Ct}(K, a, m_0, e^*)$ end K		

Figure 3: Oracles of the *one-time* IND-CCA security game for an AEAD scheme (Enc, Dec) , where **encrypt** is a one-time oracle.

An additional property that hash-based AEAD schemes satisfy and which we require (and which [ACD19] do not require in their weaker AEAD notion) is the ability to “explain” ciphertexts. That is, the ability to program any ciphertext so that given a chosen key, message and associated data, the ciphertext decrypts with the key consistently to the message after the fact, and furthermore so that encrypting the message under the key indeed results in the same ciphertext (similar to non-committing encryption).² It will later be seen that such explainability is required by our ideal functionality for Signal.

For instance, by appropriately programming a random oracle after the fact, we can ensure that a ciphertext generated ahead of time correctly decrypts to the chosen message with the chosen key, and that the message encrypts to that ciphertext with the key. We model this “explainability” aspect via an additional algorithm called AEAD-Expl-Ct that takes an AEAD key, associated data, message, and ciphertext, and programs them to be consistent. In the security game, the adversary can query oracle **corr** at which point if $b = 1$ the challenger explains the ciphertext as an encryption of m_0 under K and returns K , or if $b = 0$, simply returns K . No matter the setting of bit b , after corruption the game ends without loss of generality as the adversary can encrypt and decrypt on its own. All AEAD schemes considered in this work are required to provide a AEAD-Expl-Ct algorithm.

The advantage of an adversary \mathcal{A} attacking an AEAD scheme AE is denoted by $\mathcal{A}_{\text{ot-cca}}^{\text{AE}}(\mathcal{A})$; the attacker is parametrized by its running time t .

Definition 2. An AEAD scheme AE is (t, ε) -one-time-CCA-secure if for all t -attackers \mathcal{A} ,

$$\mathcal{A}_{\text{ot-cca}}^{\text{AE}}(\mathcal{A}) \leq \varepsilon.$$

B Building Blocks

In this work, we present **Signal** and our stronger **Signal⁺** protocols as modular constructions that use three components (following [ACD19]): Key Derivation Function Chains (KDF Chains)—used to advance forward the public ratchet; continuous key-agreement (CKA)—used to generate the secrets

²Constructions used by Signal, such as AEAD based on CBC+HMAC and SIV, will satisfy this property.

that advance forward the public ratchet; and forward-secure authenticated encryption with associated data (FS-AEAD)—the symmetric ratchet itself. These components are presented in isolation in this section before combining them into the Signal and Signal⁺ schemes in Section C. For CKA and FS-AEAD we first provide security definitions for the two primitives and then constructions achieving them. However, the Signal and Signal⁺ schemes we show later can use any construction of the two which satisfy their corresponding definitions. For KDF chains, we do not provide any definitions but rather achieve their security and functionality by modelling the underlying KDF as a Random Oracle in Section C. For a concrete instantiation of the underlying KDF, we point to that which the Signal protocol uses: HKDF [KE10] with SHA-256 or SHA-512 [FIP95]. This section often follows the work of [ACD19] verbatim.³

B.1 Key Derivation Function Chains

The Signal protocol makes use of Key Derivation Function (KDF) Chains for the public ratchet. We take (almost) verbatim from its whitepaper by Marlinspike and Perrin [MP16a] the desired syntax and security properties of KDFs and KDF Chains: A KDF is a cryptographic function that takes as input a secret and random KDF key σ and some input data I and returns output data R (i.e., $R \leftarrow \text{KDF}(\sigma, I)$). The output data R is indistinguishable from random provided the key isn't known (i.e. a KDF satisfies the requirements of a cryptographic “PRF”). If the key is not secret and random, the KDF should still provide a secure cryptographic hash of its key and input data.

To build a KDF Chain, Marlinspike and Perrin use iterated computations of a KDF, where one part of the output is used as a separate output key k and the other part is used to replace the KDF key σ , which can then be used with another input (i.e., $(\sigma', k) \leftarrow \text{KDF}(\sigma, I)$). KDF chains have the following intuitive properties:

- **Resilience:** The output keys appear random to an adversary without knowledge of the KDF keys. This is true even if the adversary can control the KDF inputs.
- **FS:** Output keys from the past appear random to an adversary who learns the KDF key at some point in time.
- **PCS:** Future output keys appear random to an adversary who learns the KDF key at some point in time, provided that future inputs have added sufficient entropy.

As stated earlier, we simply model the KDF that underlies the KDF chain in Signal and Signal⁺ as a Random Oracle H . Thus, provided that either the KDF key or the input data is random and unknown to the adversary, the output will be random and unknown to the adversary, and thus both resilience as well as PCS are achieved. Likewise, FS is easily seen to be achieved.

B.1.1 Differences from [ACD19]

Alwen *et al.* capture the above security properties of a KDF Chain in a primitive which they call *PRF-PRNGs* [ACD19]. They also provide a corresponding (deterministic) construction in the standard model. We refer the reader to Section 4.3 of their work for more details. However, we emphasize that a standard model PRF-PRNG construction suffices in their work only because in

³ [ACD19] use A and B to refer to the two parties of CKA and FS-AEAD. To remain consistent with the notation of our ideal functionality for Secure Messaging, we instead use P_1 and P_2 .

their SM security definition, leakages are not allowed during challenge epochs (see Section D.2.1). On the other hand, our functionality requires the simulator to generate fake ciphertexts for an epoch, even if the adversary may later leak on the parties to reveal the keys for these ciphertexts. So, while these fake ciphertexts need to be at first indistinguishable from the real ciphertexts, and thus the keys used to encrypt them need to be sampled randomly, the adversary can later obtain the root KDF chain key and CKA shared secret for that epoch, from which the FS-AEAD initialization key is derived using the function underlying the KDF chain. Thus, we need to be able to properly program this function to output the correct (random) key. Furthermore, we also need the ability to program this function in order to prove that Signal achieves certain non-malleability properties (similar to the proof of CCA-security for hashed ElGamal [ABR01, CS03, KM04]; c.f. Section D.2.4).

B.2 Forward-Secure AEAD

B.2.1 Defining FS-AEAD

Forward-secure authenticated encryption with associated data is a stateful primitive between a sender P_1 and a receiver P_2 and can be considered a single-epoch variant of Signal, a fact that is also evident from its security definition.

Definition 3. Forward-secure authenticated encryption with associated data (FS-AEAD) is a tuple of algorithms $\text{FS-AEAD} = (\text{FS-Init-S}, \text{FS-Init-R}, \text{FS-Send}, \text{FS-Rcv})$, where

- FS-Init-S (and similarly FS-Init-R) takes a key k and outputs a state $v_{P_1} \leftarrow \text{FS-Init-S}(k)$,
- FS-Send takes a state v , associated data a , and a message m and produces a new state and a ciphertext $(v', e) \leftarrow \text{FS-Send}(v, a, m)$, and
- FS-Rcv takes a state v , associated data a , and a ciphertext e and produces a new state, an index, and a message $(v', i, m) \leftarrow \text{FS-Rcv}(v, a, e)$.

Observe that all algorithms of an FS-AEAD scheme are deterministic.

Memory management. In addition to the basic syntax above, it is useful to define the following two functions FS-Stop (called by the sender) and FS-Max (called by the receiver) for memory management:

- FS-Stop , given an FS-AEAD state v , outputs how many messages have been sent (sic) and then “erases” the FS-AEAD session corresponding to v from memory; and
- FS-Max , given a state v and an integer ℓ , remembers ℓ internally such that the session corresponding to v is erased from memory as soon as ℓ messages have been received.

These features will be useful in the full protocol (cf. Section C) to be able to terminate individual FS-AEAD sessions when they are no longer needed. Providing a formal requirement for these additional functions is omitted. Moreover, since an attacker can infer the value of the message counter from the behavior of the protocol anyway, there is no dedicated oracle included in the security game below.

Security Game for FS-AEAD

<p>init</p> <pre> $k \xleftarrow{\\$} \mathcal{K}$ $v_{P_1} \leftarrow \text{FS-Init-S}(k)$ $v_{P_2} \leftarrow \text{FS-Init-R}(k)$ $i_{P_1} \leftarrow 0$ $\text{corr} \leftarrow \text{false}$ $\text{trans, comp, sent} \leftarrow \emptyset$ $\text{RCVD} \leftarrow \emptyset$ $b \xleftarrow{\\$} \{0, 1\}$ </pre> <p>corr-P₁</p> <pre> $\text{corr} \leftarrow \text{true}$ if $\text{comp} = \emptyset$ $\text{comp} \stackrel{\pm}{\leftarrow} i_{P_1}$ return v_{P_1} </pre> <p>corr-P₂</p> <pre> $\text{corr} \leftarrow \text{true}$ $\text{chall} = \{(i, a, m, e) : (i, \text{chall}, a, *, m, e) \in \text{trans}\}$ if $b = 1$ $\text{FS-Expl-In-Trans-Cts}(v_{P_2}, \text{chall})$ set-comp() return v_{P_2} </pre>	<p>transmit-P₁ (a, m)</p> <pre> $i_{P_1} ++$ $(v_{P_1}, e) \leftarrow \text{FS-Send}(v_{P_1}, a, m)$ record(non-chall, a, m, \perp, e) return e </pre> <p>chall-P₁ (a, m_0, m_1)</p> <pre> req $\neg \text{corr}$ and $m_0 = m_1$ $i_{P_1} ++$ $(v_{P_1}, e) \leftarrow \text{FS-Send}(v_{P_1}, a, m_b)$ record(chall, a, m_b, m_0, e) return e </pre> <p>corr-init-key</p> <pre> $\text{corr} \leftarrow \text{true}$ $\text{chall} = \{(i, a, m, e) : (i, \text{chall}, a, *, m, e) \in \text{sent}\}$ if $b = 1$ $\text{FS-Expl-Vul-Cts}(k, \text{chall})$ $\text{comp} = \{1\}$ return k </pre>	<p>deliver-P₂ (a, e)</p> <pre> req $(i, \text{flag}, a, m, *, e) \in \text{trans}$ for some i, m $(v_{P_2}, i', m') \leftarrow \text{FS-Rcv}(v_{P_2}, a, e)$ if $(i', m') \neq (i, m)$ win if $\text{flag} = \text{chall}$ $m' \leftarrow \perp$ $\text{RCVD} \stackrel{\pm}{\leftarrow} i$ delete(i) return (i', m') </pre> <p>inject-P₂ (a, e)</p> <pre> req $(*, *, a, *, *, e) \notin \text{trans}$ $(v_{P_2}, i', m') \leftarrow \text{FS-Rcv}(v_{P_2}, a, e)$ if $m' \neq \perp$ and $(i' \in \text{RCVD}$ or $(i' \notin \text{comp}$ and $i' < \max\{\text{comp}\}))$ win $\text{RCVD} \stackrel{\pm}{\leftarrow} i'$ delete(i') return (i', m') </pre>
<p>set-comp ()</p> <pre> $\text{mr} \leftarrow \max\{\text{RCVD}\}^a$ if $\text{comp} = \emptyset$ $\text{comp} \stackrel{\pm}{\leftarrow} \text{mr}$ else if $\max\{\text{comp}\} > \text{mr}$ $\text{comp} \stackrel{-}{\leftarrow} \max\{\text{comp}\}$ $\text{comp} \stackrel{\pm}{\leftarrow} \text{mr}$ for $i : ((i, *, *, *, *, *) \in \text{trans})$ and $(i < \max\{\text{comp}\})$ $\text{comp} \stackrel{\pm}{\leftarrow} i$ </pre>	<p>record (flag, a, m, e)</p> <pre> $\text{rec} \leftarrow (i_{P_1}, \text{flag}, a, m, m', e)$ $\text{trans, sent} \stackrel{\pm}{\leftarrow} \text{rec}$ </pre>	<p>delete (i)</p> <pre> $\text{rec} \leftarrow (i, \text{flag}, a, m, m', e)$ for m, m', a, e s.t. $(i, \text{flag}, a, m, m', e) \in \text{trans}$ $\text{trans} \stackrel{-}{\leftarrow} \text{rec}$ </pre>

^a $\max\{\text{RCVD}\} = 0$ if $\text{RCVD} = \emptyset$

Figure 4: Oracles corresponding to party P₁ of the FS-AEAD security game for a scheme FS-AEAD = (FS-Init-S, FS-Init-R, FS-Send, FS-Rcv).

Correctness and security. Both correctness and security are built into the security game depicted in Figure 4. One can observe that it corresponds to a single epoch of Signal and it thus inherits those properties. When messages are sent, they are identified by a simple counter. Thus, for correctness, all honestly generated ciphertexts should be decrypted by the recipient to their

corresponding messages with their corresponding indices.

For security, the $\text{chall-P}_1(a, m_0, m_1)$ oracle on input messages m_0 and m_1 of the same length returns an encryption of m_b , depending on a randomly chosen bit b . FS with respect to corruptions of either party is required: If P_1 is corrupted then the adversary still should not know the underlying plaintext of previously generated ciphertexts. Moreover, if P_1 has sent i_{P_1} messages thus far, no messages with lower indices can be modified or injected. If P_2 is corrupted then the adversary still should not know the underlying plaintext of previously generated and delivered ciphertexts. Moreover, for challenge ciphertexts that are in-transit at the time of such a corruption, we require the FS-AEAD scheme to “explain” them. That is, we require the ability to program ciphertexts so that they decrypt consistently to any chosen message after the fact, and furthermore so that encrypting the message under the corrupted state in the same manner as P_1 would have honestly done indeed results in the same ciphertext (similar to non-committing encryption).⁴ We model this explicitly as an algorithm $\text{FS-Expl-In-Trans-Cts}$ that takes the corrupted state, as well as the list of in-transit challenge messages and their corresponding associated data and ciphertexts, and programs the messages and ciphertexts to be consistent. Furthermore, we allow the initialization key k to be corrupted, in which case we require *all* challenge ciphertexts to be “explained” in the same way. We model this explicitly as an algorithm FS-Expl-Vul-Cts that takes the initial key k , as well as all sent challenge messages and their corresponding associated data and ciphertexts, and programs the messages and ciphertexts to be consistent. All FS-AEAD schemes considered in this work are required to provide both algorithms (really FS-Expl-Vul-Cts is only required for the proof of Signal , but it can easily be achieved with a random oracle, anyway).

The advantage of an attacker \mathcal{A} against an FS-AEAD scheme FS-AEAD is denoted by the expression $\mathcal{A}_{\text{fs-aead}}^{\text{FS-AEAD}}(\mathcal{A})$. The attacker is parameterized by its running time t and the total number of queries q it makes.

Definition 4. An FS-AEAD scheme FS-AEAD is (t, q, ε) -secure if for all (t, q) -attackers \mathcal{A} ,

$$\mathcal{A}_{\text{fs-aead}}^{\text{FS-AEAD}}(\mathcal{A}) \leq \varepsilon .$$

Definition 5. An FS-AEAD scheme is simply called ε -secure if for every $t, q \in \text{poly}(\kappa)$, $\mathcal{A}_{\text{fs-aead}}^{\text{FS-AEAD}}(\mathcal{A}) \leq \varepsilon$ and $\varepsilon \in \text{negl}(\kappa)$, where κ is the security parameter.

B.2.2 Differences from [ACD19]

Explaining ciphertexts and init key corruption. Our ideal functionalities $\mathcal{F}_{\text{Signal}}$ and $\mathcal{F}_{\text{Signal}^+}$ of course allow the adversary to leak on either party at all times (unlike the Secure Messaging definition of [ACD19]). Thus, although the ideal adversary may at first only get the length of a message when it is sent and still need to simulate the corresponding ciphertext, if a leak occurs on the recipient while the ciphertext is still in-transit, the ideal adversary needs to be able to *explain* the ciphertext as the real message. Thus, the underlying FS-AEAD security game must also require such explanation of in-transit challenge ciphertexts for our security proofs of Signal and Signal^+ . Furthermore, in (only) Signal , state leakages on the sender of an epoch may leak all *vulnerable* messages of the epoch to the adversary, and thus the FS-AEAD must explain *all* challenge ciphertexts sent.

⁴Constructions used by Signal , which use an underlying AEAD based on CBC+HMAC and SIV, will satisfy this property.

Continuing the game if P_2 is corrupted. In the FS-AEAD definition of [ACD19], the game ends once P_2 is corrupted. However, even if such a corruption occurs, if the adversary chooses to still deliver honest ciphertexts, then we must still require correctness of the FS-AEAD.

Fixing comp. Intuitively, if an attacker corrupts P_1 then it can successfully inject messages for a later index $i > i_{P_1}$. However, the [ACD19] definition declares that the attacker wins if this happens (since for their differently defined dictionary **comp** in this situation, $i \notin \text{comp}$ if the attacker did not make any **transmit- P_1** or **chall- P_1** queries, which triggers **win** in the **inject- P_2** oracle of their game). Thus, we use our own logic for defining **comp**. (We also explicitly require FS-Rcv to only accept one message per index and output \perp if it does receive more than one, in order to be properly used for our ideal functionalities. The security game uses set RCVD for this purpose.)

B.2.3 Instantiating FS-AEAD

An FS-AEAD scheme can be easily constructed from two components:

- an AEAD scheme $AE = (\text{Enc}, \text{Dec})$, and
- a KDF $H : \mathcal{W} \rightarrow \mathcal{W} \times \mathcal{K}$, where \mathcal{K} is the key space of the AEAD scheme.

The scheme is described in Figure 5. For simplicity the states of sender P_1 and receiver P_2 are not made explicit; it consists of the variables set during initialization. The main idea of the scheme, is that P_1 and P_2 share KDF key w . KDF key w is initialized with a pre-shared key $k \in \mathcal{W}$, which is assumed to be chosen uniformly at random. Both parties keep local counters i_{P_1} and i_{P_2} , respectively.⁵ P_1 , when sending the i^{th} message m with associated data (AD) a , uses H to expand the current state to a new state and an AEAD key $(w, K) \leftarrow G(w)$ and computes an AEAD encryption under K of m with AD $h = (i, a)$.

Since P_2 may receive ciphertexts out of order, whenever he receives a ciphertext, he first checks whether the key is already stored in a dictionary \mathcal{D} . If the index of the message is higher than expected (i.e., larger than $i_{P_2} + 1$), P_2 skips the KDF chain ahead and stores the skipped keys in \mathcal{D} . In either case, once the key is obtained, it is used to decrypt. If decryption fails, FS-Rcv throws an exception (**error**), which causes the state to be rolled back to where it was before the call to FS-Rcv.

Algorithms FS-Expl-In-Trans-Cts and FS-Expl-Vul-Cts. Lastly, we explain the algorithms FS-Expl-In-Trans-Cts(v_{P_2}, chall) and FS-Expl-Vul-Cts(k, chall) provided by our FS-AEAD. In order to explain ciphertexts correctly, they rely on the explanation algorithm AEAD-Expl-Ct of the underlying AEAD scheme. Formally, for every $(i, a, m, e) \in \text{chall}$, FS-Expl-In-Trans-Cts executes AEAD-Expl-Ct(K_i, a, m, e), where $K_i \leftarrow \mathcal{D}[i]$ for all $i < i_{P_2}$ and for all $i > i_{P_2}$, K_i is computed directly via $i - i_{P_2}$ KDF invocations on current KDF key $w_{i_{P_2}}$ in v_{P_2} . Algorithm FS-Expl-Vul-Cts proceeds similarly, instead computing K_i directly via i KDF invocations on initial KDF key w_0 .

Theorem 1. *Assume AE is a $(t', \varepsilon_{\text{aead}})$ -OT-CCA secure AEAD scheme and H is modelled as a random oracle. Then, the above FS-AEAD scheme FS-AEAD is (t, q, ε) -secure for $t \approx t'$ and*

$$\varepsilon \leq q \cdot \varepsilon_{\text{aead}} .$$

⁵For ease of description, the FS-AEAD state of the parties is not made explicit as a variable v .

Forward-Secure AEAD

<pre> Init-P₁(k) w ← k i_{P₁} ← 0 Init-P₂(k) w ← k i_{P₂} ← 0 D[·] ← ⊥ try-skipped(i) K ← D[i] D[i] ← ⊥ return K </pre>	<pre> FS-Send(a, m) i_{P₁} ++ (w, K) ← H(w) h ← (i_{P₁}, a) e ← Enc(K, h, m) return (i_{P₁}, e) skip(u) while i_{P₂} < u - 1 i_{P₂} ++ (w, K) ← H(w) D[i_{P₂}] ← K </pre>	<pre> FS-Rcv(a, c) (i, e) ← c K ← try-skipped(i) if K = ⊥ and i ≤ i_{P₂} error else if K = ⊥ and i > i_{P₂} skip(i) (w, K) ← H(w) i_{P₂} ← i h ← (i, a) m ← Dec(K, h, e) if m = ⊥ error return (i, m) </pre>
---	---	--

Figure 5: FS-AEAD scheme based on AEAD and a KDF H .

Proof. The proof is a straight-forward hybrid argument: Let H_0 denote the actual FS-AEAD security game.

- In hybrid experiment H_1 , the **win** condition inside the **deliver-P₂** oracle is removed. The perfect correctness of the proposed FS-AEAD scheme is easily seen by inspection and the correctness of the underlying AEAD scheme. Hence H_0 and H_1 are indistinguishable.
- Now, the security of H_1 follows from that of the underlying AEAD. It is obvious that injections for already received indices fail, since the FS-AEAD clearly outputs **error**. Moreover, injections for uncompromised indices ($i \notin \text{comp}$ and $i < \max\{\text{comp}\}$) fail by the security of the underlying AEAD (since the attacker does not have the corresponding key). Also, we program the random oracle so that all corrupted (random) AEAD keys (and subsequent KDF keys) are simulated correctly to the attacker (for example if the attacker queries **chall-P₁**, followed by **corr-P₁** then **corr-P₂**). Explainability follows from that of the AEAD.

□

B.3 Continuous Key Agreement

As in the work of Alwen et al. [ACD19], we separate out the primitive that generates the secrets for public-ratchet updates in **Signal** and **Signal⁺**, and call it *continuous key agreement (CKA)*. In this section, we present our definitions of CKA and instantiations from CDH in the presence of DDH oracle (often referred to as the *gap DH assumption* [OP01]).⁶

⁶We do not provide CKA schemes secure according to our definitions based on LWE or generic KEMs, as in ACD. However, we note that our stronger scheme CKA⁺ is intuitively at least as strong as their construction from generic KEMs, but more efficient.

B.3.1 Defining CKA

At a high level, CKA is a synchronous two-party protocol between P_1 and P_2 . Odd rounds i consist of P_1 sending and P_2 receiving a message T_i , whereas in even rounds, P_2 is the sender and P_1 the receiver. Each round i also produces a key I_i , which is output by the sender upon sending T_i and by the receiver upon receiving T_i .

Definition 6. A continuous-key-agreement (CKA) scheme is a quadruple of algorithms $\text{CKA} = (\text{CKA-Init-}P_1, \text{CKA-Init-}P_2, \text{CKA-S}, \text{CKA-R})$, where

- $\text{CKA-Init-}P_1$ (and similarly $\text{CKA-Init-}P_2$) takes a key k and produces an initial state $\gamma^{P_1} \leftarrow \text{CKA-Init-}P_1(k)$ (and γ^{P_2}),
- CKA-S takes a state γ , and produces a new state, message, and key $(\gamma', T, I) \stackrel{\$}{\leftarrow} \text{CKA-S}(\gamma)$, and
- CKA-R takes a state γ and message T and produces new state and a key $(\gamma', I) \leftarrow \text{CKA-R}(\gamma, T)$.

Denote by \mathcal{K} the space of initialization keys k and by \mathcal{I} the space of CKA keys I .

Correctness. A CKA scheme is correct if in the security game in Figure 6 (explained below), P_1 and P_2 always, i.e., with probability 1, output the same key in every round.

Security. The security property we will require a CKA scheme to satisfy is that conditioned on the transcript T_1, T_2, \dots , the keys I_1, I_2, \dots are unrecoverable. An attacker against a CKA scheme is required to be passive, i.e., may not modify the messages T_i . However, it is given the power to possibly (1) control the random coins used by the sender and (2) leak the current state of either party. Given the capabilities of the adversary, it is easy to see that some keys I_i would be recoverable. The security guarantee offered by the CKA scheme would then be that even given the transcript T_1, T_2, \dots , assuming certain “fine-grained” conditions around when the adversary controls the randomness used by parties and when the adversary learns the state of parties, most keys I_1, I_2, \dots are unrecoverable. It will also be the case that parties thus recover from such bad randomness and leakage issued by the adversary.

The formal security game for CKA is provided in Figure 6. It begins with a call to the **init** oracle, which initializes the states of both parties, and defines epoch counters t_{P_1} and t_{P_2} . Procedure **init** takes a value t^* , which determines in which round the challenge oracle may be called; the task of the adversary will be to recover the key I_{t^*} for that round.

Upon completion of the initialization procedure, the attacker gets to interact arbitrarily with the remaining oracles, as long as *the calls are in a “ping-pong” order*, i.e., a call to a send oracle for P_1 is followed by a receive call for P_2 , then by a send oracle for P_2 , etc. The attacker only gets to use the challenge oracle for epoch t^* . The attacker additionally has the capability of testing the consistency of T_t and I_t (i.e., whether the corresponding receiver in epoch t would produce key I_t on input message T_t).

The security game of Alwen et al. [ACD19] is parametrized by Δ_{CKA} , which stands for the number of epochs that need to pass after t^* until the states do not contain secret information pertaining to the challenge. Once a party reaches epoch $t^* + \Delta_{\text{CKA}}$, its state may be revealed to the attacker (via the corresponding corruption oracle). We avoid this and define two levels of

Security Games for CKA

init (t^*) $k \xleftarrow{\$} \mathcal{K}$ $\gamma_0^{P_1} \leftarrow \text{CKA-Init-P}_1(k)$ $\gamma_0^{P_2} \leftarrow \text{CKA-Init-P}_2(k)$ $t_{P_1}, t_{P_2} \leftarrow 0$ $\text{Recv-State}[*] \leftarrow \perp$	send-P₁ $t_{P_1} ++$ $(\gamma_{t_{P_1}}^{P_1}, T, I) \xleftarrow{\$} \text{CKA-S}(\gamma_{t_{P_1}}^{P_1})$ $\text{Recv-State}[t_{P_1} + 1] \leftarrow \gamma_{t_{P_1}}^{P_1}$ return (T, I)	chall-P₁ $t_{P_1} ++$ req $t_{P_1} = t^*$ $(\gamma_{t_{P_1}}^{P_1}, T, I) \xleftarrow{\$} \text{CKA-S}(\gamma_{t_{P_1}}^{P_1})$ return T
corr-P₁ req allow-corr_{P_1} return $\gamma_{t_{P_1}}^{P_1}$	send-P₁' (r) $t_{P_1} ++$ req $\text{allow-bad-rand}_{P_1}$ $(\gamma_{t_{P_1}}^{P_1}, T, I) \leftarrow \text{CKA-S}(\gamma_{t_{P_1}-1}^{P_1}; r)$ $\text{Recv-State}[t_{P_1} + 1] \leftarrow \gamma_{t_{P_1}}^{P_1}$ return (T, I)	test (t, T, I) req $\text{Recv-State}[t] \neq \perp$ if $(*, I) \leftarrow \text{CKA-R}(\text{Recv-State}[t], T)$ return 1 else return 0
	receive-P₁ $t_{P_1} ++$ $(\gamma_{t_{P_1}}^{P_1}, *) \leftarrow \text{CKA-R}(\gamma_{t_{P_1}-1}^{P_1}, T)$	

$$\text{allow-corr}_{P_1}, \text{allow-bad-rand}_{P_1} : \iff \begin{cases} t_{P_1} \neq t^* & t^* \text{ is odd} \\ t_{P_1} \neq t^* - 1 & t^* \text{ is even} \end{cases}$$

$$\text{allow-corr}_{P_2}, \text{allow-bad-rand}_{P_2} : \iff \begin{cases} t_{P_2} \neq t^* - 1 & t^* \text{ is odd} \\ t_{P_2} \neq t^* & t^* \text{ is even} \end{cases}$$

$$\text{allow-corr}_{P_1} : \iff t_{P_1} \neq t^* - 1 \vee t^* \text{ is odd}$$

$$\text{allow-bad-rand}_{P_1} : \iff t_{P_1} \neq t^* \vee t_{P_1} \neq t^* - 1$$

$$\text{allow-corr}_{P_2} : \iff t_{P_2} \neq t^* - 1 \vee t^* \text{ is even}$$

$$\text{allow-bad-rand}_{P_2} : \iff t_{P_2} \neq t^* \vee t_{P_2} \neq t^* - 1$$

Figure 6: Oracles corresponding to party P_1 of the CKA security game for a scheme $\text{CKA} = (\text{CKA-Init-P}_1, \text{CKA-Init-P}_2, \text{CKA-S}, \text{CKA-R})$; the oracles for P_2 are defined analogously. Conditions for the weaker security game, i.e., ε -security, are presented above those for the stronger game, i.e., $(\varepsilon, +)$ -security.

security, the former weaker than the latter. At the bottom of Figure 6, the conditions allow-corr_p and allow-bad-rand_p for the weaker version are presented above those of the stronger version. We define two levels of security in order to capture a stronger, more fine-grained security guarantee for CKA which will be useful in providing stronger security guarantees for **Signal** and **Signal⁺** as a whole when one considers the composition of all its building blocks, CKA being one of them. In the former, bad randomness is not allowed in the epochs t^* and $t^* - 1$, and corruptions are not allowed in the epoch t^* after invoking CKA-S (for the sender of epoch t^*) and before invoking CKA-R (for

the receiver of epoch t^*). In the latter, bad randomness is not allowed in the epochs t^* and $t^* - 1$, and corruption of the receiver of epoch t^* is not allowed before invoking CKA-R (for epoch t^*). There is no other difference between the two notions.

The game ends (not made explicit) once both states are revealed after the challenge phase. The attacker wins the game if it eventually outputs the correct secret key I_{t^*} corresponding to the challenge message T_{t^*} . The advantage of an attacker \mathcal{A} against a CKA scheme CKA is denoted by $\mathcal{A}^{\text{CKA}}(\mathcal{A})$ and $\mathcal{A}^{\text{CKA}^+}(\mathcal{A})$ for the weaker and stronger security guarantees, respectively. The attacker is parameterized by its running time t .

Definition 7. A CKA scheme CKA is (t, ε) -secure if for all t -attackers \mathcal{A} ,

$$\mathcal{A}^{\text{CKA}}(\mathcal{A}) \leq \varepsilon.$$

Definition 8. A CKA scheme CKA is $(t, \varepsilon, +)$ -secure if for all t -attackers \mathcal{A} ,

$$\mathcal{A}^{\text{CKA}^+}(\mathcal{A}) \leq \varepsilon.$$

Definition 9. A CKA scheme CKA is simply called ε -secure (resp. $(\varepsilon, +)$ -secure) if for every $t \in \text{poly}(\kappa)$, $\mathcal{A}^{\text{CKA}}(\mathcal{A}) \leq \varepsilon$ (resp. $\mathcal{A}^{\text{CKA}^+}(\mathcal{A}) \leq \varepsilon$) and $\varepsilon \in \text{negl}(\kappa)$, where κ is the security parameter.

Remark 1. Many natural CKA schemes satisfy an additional property that given a CKA message T and key I for a given round, it is possible to deterministically compute the corresponding state of the receiving party after her execution of CKA-R in that round. We model this explicitly by requiring a deterministic algorithm CKA-Der-R that takes a message T and key I and produces the correct state $\gamma' \leftarrow \text{CKA-Der-R}(T, I)$. All CKA schemes in this work are required to be natural (including the public-ratchet update mechanism of Signal).

B.3.2 Differences from [ACD19]

Fine-grained security guarantees. Recall the “CKA from DDH” scheme from [ACD19] (which is the public ratchet used in Signal), $\text{CKA} = (\text{CKA-Init-P}_1, \text{CKA-Init-P}_2, \text{CKA-S}, \text{CKA-R})$, that is instantiated in a cyclic group $G = \langle g \rangle$ as follows:

- The initial shared state $k = (h, x_0)$ consists of a (random) group element $h = g^{x_0}$ and its discrete logarithm x_0 . The initialization for P_1 outputs $h \leftarrow \text{CKA-Init-P}_1(k)$ and that for P_2 outputs $x_0 \leftarrow \text{CKA-Init-P}_2(k)$.
- The send algorithm CKA-S takes as input the current state $\gamma = h$ and proceeds as follows: It
 1. chooses a random exponent x ;
 2. computes the corresponding key $I \leftarrow h^x$;
 3. sets the CKA message to $T \leftarrow g^x$;
 4. sets the new state to $\gamma \leftarrow x$; and
 5. returns (γ, T, I) .
- The receive algorithm CKA-R takes as input the current state $\gamma = x$ as well as a message $T = h$ and proceeds as follows: It

1. computes the key $I = h^x$;
2. sets the new state to $\gamma \leftarrow h$; and
3. returns (γ, I) .

Now, let x_0 be the random exponent that is part of the initial shared state, and for $i > 0$, let x_i be the random exponent picked by CKA-S (which was run by P_1 for odd i , and P_2 for even i) in round i . Then, we have the following:

- The key for round i is $I_i = g^{x_{i-1}x_i}$.
- The message for round i is $T_i = g^{x_i}$.
- If i is odd, and P_1 has yet to invoke CKA-S, $\gamma^{P_1} = g^{x_{i-1}}$ and $\gamma^{P_2} = x_{i-1}$.
- If i is odd, and P_1 has invoked CKA-S, but P_2 has yet to invoke CKA-R, $\gamma^{P_1} = x_i$ and $\gamma^{P_2} = x_{i-1}$.
- If i is odd, and P_1 has invoked CKA-S, and P_2 has invoked CKA-R, $\gamma^{P_1} = x_i$ and $\gamma^{P_2} = g^{x_i}$.
- If i is even, and P_2 has yet to invoke CKA-S, $\gamma^{P_1} = x_{i-1}$ and $\gamma^{P_2} = g^{x_{i-1}}$.
- If i is even, and P_2 has invoked CKA-S, but P_1 has yet to invoke CKA-R, $\gamma^{P_1} = x_{i-1}$ and $\gamma^{P_2} = x_i$.
- If i is even, and P_2 has invoked CKA-S, and P_1 has invoked CKA-R, $\gamma^{P_1} = g^{x_i}$ and $\gamma^{P_2} = x_i$.

Based on the above, we make the following observations:

- If i is odd and P_1 is corrupted after invoking CKA-S, the adversary learns $\gamma^{P_1} = x_i$ and since it also has access to g^{x_j} for all $j \geq 1$, the adversary learns I_i and I_{i+1} .
- If i is even and P_1 is corrupted after invoking CKA-R, and P_2 used good randomness in picking x_i while invoking CKA-S in round i , the adversary learns $\gamma^{P_1} = g^{x_i}$, but since it only (assuming no other corruptions) has access to g^{x_j} for all $j \geq 1$, the adversary does not learn I_i (if P_1 also used good randomness in picking x_{i-1} while invoking CKA-S in round $i-1$) or I_{i+1} (if P_1 also uses good randomness in picking x_{i+1} while invoking CKA-S in round $i+1$).

Thus, the CKA keys for two rounds are compromised only in the case where the party that has last sent a message is corrupted, and not if the party has last received a message. This allows us to consider our stronger version of the CKA security game than the one described in [ACD19].

Non-malleability. Consider the following scenario in **Signal** or **Signal⁺**: It is P_1 's turn to start a new sending epoch, but she has not yet. Then her state is leaked, and afterwards, she sends the first message m_1 of the epoch with good randomness. Then, if P_2 started her last epoch with good randomness, and there are no other leakages, m_1 is required to remain private by $\mathcal{F}_{\text{Signal}}$ and $\mathcal{F}_{\text{Signal}^+}$, respectively. However, all authenticity for m_1 is lost—the adversary leaked on P_1 beforehand and thus could have generated the message herself. Therefore, we replace the indistinguishability definition of [ACD19] with our recoverability definition and require non-malleability of CKA messages via the **test** oracle—the adversary should not be able to maul them in order to learn about the actual session messages sent in **Signal** or **Signal⁺**. See the full security proof of Theorem 5 for **Signal** and **Signal⁺**, as well as Appendix D.2.4, for more details.

B.3.3 Instantiating CKA

First, we show that the above scheme is (t, ε) -secure. Note that it is *natural*, i.e., it supports a CKA-Der-R algorithm, namely, $\text{CKA-Der-R}(T, I) = T$. For the theorem below, let a group G be (t, ε) -secure if every attacker with running time at most t has advantage at most ε at solving CDH challenges in the presence of a DDH oracle.⁷

Theorem 2. *Assume group G is (t, ε) -CDH-secure in the presence of a DDH oracle. Then, the above CKA scheme CKA is (t, ε) -secure for $t \approx t'$. Furthermore, this CKA scheme also has dense transcripts.*

Proof. Assume w.l.o.g. that t^* is *odd*, i.e., P_1 sends the challenge; the case where t^* is even is handled analogously. Let (g^a, g^b) be a CDH challenge. The reduction simulates the CKA protocol in the straight-forward way but embeds the challenge into the CKA as follows:

- in epoch $t^* - 1$, it uses $T_{t^*-1} = g^a$ and $I_{t^*-1} = g^{xa}$, where x is the exponent used to simulate $T_{t^*-2} = g^x$.
- in epoch t^* , it uses $T_{t^*} = g^b$ and $I_{t^*} = g^{ab}$ which is the key the adversary is to recover; and
- in epoch $t^* + 1$, for exponent x' (possibly generated using adversarial randomness), it uses $T_{t^*+1} = g^{x'}$ and $I_{t^*+1} = g^{bx'}$.

It is easy to verify that this correctly simulates the CKA experiment. Also note that the test oracle can be perfectly simulated with the help of a DDH oracle: if $\text{test}(t^*, T, I)$ is queried, the reduction simply queries the DDH oracle on (g^a, T, I) ; if $\text{test}(t^* + 1, T, I)$ is queried, the reduction simply queries the DDH oracle on (g^b, T, I) ; all other $\text{test}()$ queries can be directly simulated. □

B.3.4 Instantiating CKA⁺

A CKA scheme $\text{CKA}^+ = (\text{CKA-Init-P}_1, \text{CKA-Init-P}_2, \text{CKA-S}, \text{CKA-R})$ can be obtained assuming random oracles or circular-secure ElGamal in a cyclic group $G = \langle g \rangle$ (with exponent space \mathcal{X}) using a function $H : \mathcal{I} \rightarrow \mathcal{X}$ as follows:

- The initial shared state $k = (h, x_0)$ consists of a (random) group element $h = g^{x_0}$ and its discrete logarithm x_0 . The initialization for P_1 outputs $h \leftarrow \text{CKA-Init-P}_1(k)$ and that for P_2 outputs $x_0 \leftarrow \text{CKA-Init-P}_2(k)$.
- The send algorithm CKA-S takes as input the current state $\gamma = h$ and proceeds as follows: It
 1. chooses a random exponent x ;
 2. computes the corresponding key $I \leftarrow h^x$;
 3. sets the CKA message to $T \leftarrow g^x$;
 4. sets the new state to $\gamma \leftarrow x \cdot H(I)$; and

⁷This is often referred to as the *gap DH assumption* [OP01]: given g^a and g^b , where a and b are uniformly random and independent exponents, and oracle $\text{ddh}(\cdot, \cdot, \cdot)$, which on input $(X = g^x, Y = g^y, Z) \in G^3$, returns 1 if $Z = g^{xy}$ and 0 otherwise; the challenge is to compute g^{ab} .

5. returns (γ, T, I) .
- The receive algorithm CKA-R takes as input the current state $\gamma = x$ as well as a message $T = h$ and proceeds as follows: It
 1. computes the key $I = h^x$;
 2. sets the new state to $\gamma \leftarrow h^{H(I)}$; and
 3. returns (γ, I) .

Note that the above scheme is *natural*, i.e., it supports a CKA-Der-R algorithm, namely, $\text{CKA-Der-R}(T, I) = T^{H(I)}$. For the theorem below, let a group G be (t, ε) -secure if every attacker with running time at most t has advantage at most ε at solving CDH challenges in the presence of a DDH oracle.

Theorem 3. *Assume group G is (t, ε) -CDH-secure in the presence of a DDH oracle. Additionally, assume the existence of a random oracle H . Then, the above CKA scheme CKA is $(t, \varepsilon, +)$ -secure for $t \approx t'$. Furthermore, this CKA scheme also has dense transcripts.*

Proof. Assume w.l.o.g. that t^* is *odd*, i.e., P_1 sends the challenge; the case where t^* is even is handled analogously. Let (g^a, g^b) be a CDH challenge. The reduction simulates the CKA protocol in the straight-forward way but embeds the challenge into the CKA as follows:

- in epoch $t^* - 1$, it uses $T_{t^*-1} = g^a$ and $I_{t^*-1} = g^{xaH(I_{t^*-2})}$, where x is the exponent used to simulate $T_{t^*-2} = g^x$.
- in epoch t^* , it uses $T_{t^*} = g^b$ and $I_{t^*} = g^{abH(I_{t^*-1})}$ which is the key the adversary is to recover, as well as sets $\gamma_{t^*}^{\mathcal{P}_1} \leftarrow y$, for random y in \mathcal{X} ; and
- in epoch $t^* + 1$, for exponent x' (possibly generated using adversarial randomness), it uses $T_{t^*+1} = g^{x'}$ and $I_{t^*+1} = g^{yx'}$.

It is easy to verify that this correctly simulates the CKA experiment since H is a random oracle. In particular, randomly sampled y properly simulates $b \cdot H(I_{t^*})$: If the adversary does not query the random oracle on I_{t^*} then y is properly distributed. Moreover, when she makes a random oracle query for any I , the reduction can use the DDH oracle on $(g^{aH(I_{t^*-1})}, g^b, I)$ so that if indeed $I = I_{t^*}$, the reduction will know, and will then forward to its challenger $g^{ab} = I_{t^*}^{H(I_{t^*-1})^{-1}}$ before answering the CKA^+ attacker's query.

Similarly, the test oracle can be perfectly simulated with the help of the DDH oracle: if $\text{test}(t^*, T, I)$ is queried, the reduction simply queries the DDH oracle on $(g^{aH(I_{t^*-1})}, T, I)$; all other $\text{test}()$ queries can be directly simulated. □

Remark 2. *The above theorem can also be proved without assuming that H is a random oracle, but by assuming rather that El-Gamal is circularly secure in G . Specifically, we assume that for uniformly random and independent exponents a and b , the distributions $(g, g^a, g^b, b \cdot H(g^{ab}))$ and (g, g^a, g^b, y) are indistinguishable, where y is uniformly random in the exponent space \mathcal{X} . Based on this assumption (as opposed to H being a random oracle), we see that the embedding above still correctly simulates the CKA experiment, and essentially the same proof works for the security of the CKA scheme.*

B.3.5 Even stronger security

One can observe that CKA^+ is in fact even more secure than the $(t, \varepsilon, +)$ -security that we just proved for it. Although formalizing its full security is quite complex and does not have too much of an impact on our stronger Signal^+ protocol (see Section E.5 for an informal description of the impact that it does have, based on the below) we here informally describe a scenario in which CKA^+ does indeed achieve stronger security.

Consider the scenario in which element x_0 of the initial shared state $k = (g^{x_0}, x_0)$ is secure, P_2 never has good randomness during the session (of course for x_0 to be secure, if P_2 generated it, then she must have had good randomness at that point, but in the real-world, this may have happened a long time ago), and P_1 has good randomness for sampling her first exponent x_1 , but then never again. Moreover, assume that both P_1 and P_2 are never corrupted. Then *every* shared secret I_t will still be secure: $I_1 = g^{x_0 x_1}$ is secure by CDH (even with a DDH oracle) since the adversary only has g^{x_1} (and maybe g^{x_0}), but not x_0 nor x_1 . Shared secret $I_2 = g^{x_1 \text{H}(I_1) \cdot x_2}$ is also secure by CDH with a DDH oracle if we model H as a random oracle. First, observe that in order to compute I_2 , the adversary needs to query the random oracle on I_1 (for otherwise, I_2 is information-theoretically hidden). Now, since a reduction given g^{x_0} and g^{x_1} can give the adversary the same, then when the adversary queries the random oracle on I_1 , the reduction can query the DDH oracle on (g^{x_0}, g^{x_1}, I_1) and submit $I_1 = g^{x_0 x_1}$ to its challenger. Further observe that in order to compute any I_t , the adversary needs to query the random oracle on I_1 (for otherwise, all of I_2, \dots, I_t , where t is of course polynomial, are information-theoretically hidden). Thus, we can use the same reduction as above for every t .

Furthermore, even if P_1 is corrupted after she receives P_2 's epoch $t^* - 1$ message, for some t^* , then we can still guarantee security of I_{t^*-1} from CDH with a DDH oracle, if we model H as a random oracle (of course, all preceding epochs are secure too, from the above). When leaking $\gamma_{t^*-1}^{P_1}$, the reduction can simply sample random r and give g^r , in place of $g^{x_{t^*-1} \text{H}(I_{t^*-1})}$. If the adversary never queries the random oracle on $\text{H}(I_{t^*-1})$, then of course this is a perfect simulator. Now, as above, in order to compute I_{t^*-1} , the adversary still needs to query the random oracle on I_1 , so our same reduction will still go through. Moreover, if in this situation P_1 uses good randomness for her epoch t^* message T_{t^*} , then I_{t^*} will be secure too. We just showed that for the adversary to query the random oracle on I_{t^*-1} , meaning she has already computed I_{t^*-1} , the adversary needs to still query I_1 , which by our original reduction, we know it cannot. Given this information, a CDH reduction given g^r and g^{x_i} can simply give $\gamma_{t^*-1}^{P_1} = g^r$ instead of $g^{x_{t^*-1} \text{H}(I_{t^*-1})}$ for the leak on P_1 before epoch t^* , then send $T_{t^*} = g^{x_{t^*}}$. Then, when the adversary submits I_{t^*} , the reduction can simply forward it to its challenger.

Using similar arguments as above, if both parties use bad randomness for all subsequent epochs, and there are no more corruptions, then also all subsequent shared secrets I_t (and therefore all of them) will remain secure. Intuitively, this is because I_{t^*-1} is secure, so $\gamma_{t^*-1}^{P_2}$ is as well, and P_1 used good randomness in epoch t^* , thus of course γ_{t^*} is secure too, so we can use the above argument thereafter.

C Composition

In this section, we show how to compose the building blocks of Section B to construct Secure Messaging protocols Signal and Signal^+ that UC-realize our two ideal functionalities of Figure 2.

This section again often closely follows the work of [ACD19].

Initial Key Exchange Ideal Functionality. Before constructing **Signal** and **Signal**⁺ we must introduce an ideal functionality for an initial key exchange to be used upon initialization of a session of one of our protocols. While the actual **Signal** protocol uses the X3DH key exchange [MP16b], the focus of our work is to analyze the security and functionality of the double ratchet algorithm, and not X3DH. Therefore, we choose to present the following simple ideal functionality $\mathcal{F}_{\text{KE}}^{\text{CKA}}$ for key exchange that may be stronger than what X3DH offers, but nonetheless suffices for our purposes. The functionality is parameterized by a CKA protocol CKA.

On input (sid, **SETUP**) **from** P where $P \in \{P_1, P_2\}$: Send (sid, **SETUP**, P) to \mathcal{S} . When \mathcal{S} returns (sid, **SETUP**) then (i) sample $(\sigma_{\text{root}}, k) \xleftarrow{\$} \{0, 1\}^{2\lambda}$, (ii) execute $\gamma^{P_1} \leftarrow \text{CKA-Init-}P_1(k), \gamma^{P_2} \leftarrow \text{CKA-Init-}P_2(k)$, (iii) set $k^{P_1} \leftarrow (\sigma_{\text{root}}, \gamma^{P_1}), k^{P_2} \leftarrow (\sigma_{\text{root}}, \gamma^{P_2})$, and (iv) send (sid, **EXCHANGE**, k^{P_1}) to P and (sid, **EXCHANGE**, k^{P_2}) to \bar{P} .

C.1 Constructions

We now provide constructions of **Signal** and our modified **Signal**⁺ in the $\mathcal{F}_{\text{KE}}^{\text{CKA}}$ -hybrid model. As in the analysis of Alwen *et al.*, our presentation of **Signal** differs from the actual **Signal** protocol in a few minor aspects (see [ACD19, §5.2]), but the two are logically equivalent (and almost precisely the same otherwise). Indeed, the claim that the true **Signal** protocol UC-realizes $\mathcal{F}_{\text{Signal}}$ in the $\mathcal{F}_{\text{KE}}^{\text{CKA}}$ -hybrid model follows from Theorem D.2.1. Importantly, we consider the version of the true **Signal** protocol in which to increase security, parties defer new CKA secret key generation for their next sending epoch until they actually start that epoch, as opposed to when they receive the first message of the prior epoch [MP16a, §6.5]. Recall: we show in Section D.2.1 that while [ACD19] do the same, in their model it does not actually provide **Signal** any extra security.

As explained in the introduction, the main idea of the two schemes is that the two parties P_1 and P_2 keep track of the same root KDF Chain key σ_{root} and refresh it using the secrets of a CKA protocol that is run “in parallel”. The corresponding outputs of the root KDF Chain are used to generate initialization keys for FS-AEAD instances. The only difference between **Signal** and **Signal**⁺ is that the former uses a (t, ε) -secure CKA protocol while the latter uses a $(t, \varepsilon, +)$ -secure CKA protocol.

State. In **Signal** and **Signal**⁺, party P_1 (resp. P_2) keeps an internal state s_{P_1} (resp. s_{P_2}), which is initialized by **Init-}P_1** (resp. **Init-}P_2**) and used as well as updated by **Send** and **Rcv**. The state s_{P_1} of **SM** consists of the following values:

- The current key σ_{root} of the root KDF chain,
- States $v[0], v[1], v[2], \dots$ of the various FS-AEAD instances,
- The state γ of the corresponding CKA scheme,
- The current CKA message T_{cur} ,
- The number of messages sent in the last completed sending epoch of P_1 ℓ_{prv} , and
- An epoch counter t_{P_1} .

Init-P₁ (k^{P_1}) <pre> (σ_{root}, γ) ← k^{P₁} v[·] ← ⊥ T_{cur} ← ⊥ ℓ_{prv} ← 0 t_{P₁} ← 0 </pre>	Send-P₁ (m) <pre> if t_{P₁} <i>is even</i> t_{P₁} ++ (γ, T_{cur}, I) $\xrightarrow{\\$}$ CKA-S(γ) (σ_{root}, k) ← H(σ_{root}, I) v[t_{P₁}] ← FS-Init-S(k) h ← (t_{P₁}, T_{cur}, ℓ_{prv}) (v[t_{P₁}], e) ← FS-Send(v[t_{P₁}], h, m) return (h, e) </pre>	Rcv-P₁ (e) <pre> (h, e) ← c (t, T, ℓ) ← h req t <i>even</i> and t ≤ t_{P₁} + 1 if t = t_{P₁} + 1 ℓ_{prv} ← FS-Stop(v[t_{P₁}]) t_{P₁} ++ FS-Max(v[t - 2], ℓ) (γ, I) ← CKA-R(γ, T) (σ_{root}, k) ← H(σ_{root}, I) v[t] ← FS-Init-R(k) (v[t], i, m) ← FS-Rcv(v[t], h, e) if m = ⊥ error return (t, i, m) </pre>
---	--	---

Figure 7: Secure-messaging schemes **Signal** and **Signal⁺** in the $\mathcal{F}_{\text{KE}}^{\text{CKA}}$ -hybrid model based on (i) an FS-AEAD scheme FS-AEAD, (ii) (t, ε) - and $(t, \varepsilon, +)$ -secure CKA schemes, respectively, and (iii) a KDF Chain using HKDF H. We assume w.l.o.g. that P_1 initializes the session. The initialization keys k^{P_1} and k^{P_2} are provided by a session of $\mathcal{F}_{\text{KE}}^{\text{CKA}}$, which is also initialized by P_1 . The figure only shows the algorithms for P_1 ; P_2 's algorithms are analogous, with “even” replaced by “odd”.

We may refer to some such components of the state of P_1 throughout (e.g., in the proof of Theorem 5 in Section E) using “dot-notation”. For example, to refer to the epoch counter of party P_1 , we will write “ $s_{P_1}.t_{P_1}$ ”. Recall (cf. Section B.2) that once the maximum number of messages has been received for an epoch according to FS-Max, a session “erases” itself from the memory, and similarly when FS-Stop is called on a particular FS-AEAD session, it is erased. For simplicity, removing the corresponding $v[t]$ from memory is not made explicit in either case. The state s_{P_2} is defined analogously.

The algorithms. The algorithms of schemes **Signal** and **Signal⁺** are depicted in Figure 7 and described in more detail below. For ease of description, the algorithms **Send** and **Rcv** are presented as **Send-P₁** and **Rcv-P₁**, which handle the case where the algorithm is invoked by P_1 ; the case where the algorithm is invoked by P_2 works analogously. Moreover, to improve readability, the state s_{P_1} is not made explicit in the description: it consists of the variables set by the initialization algorithm. We also assume w.l.o.g. that P_1 initializes the session.

- **Initialization:** In the initialization procedure **Init-P₁**, P_1 initializes a session of $\mathcal{F}_{\text{KE}}^{\text{CKA}}$ to obtain initialization key $k^{P_1} = (\sigma_{\text{root}}, \gamma^{P_1})$. It consists of the initial root KDF Chain key σ_{root} and the initial CKA state of P_1 , γ^{P_1} . Furthermore, **Init-P₁** also sets the initial epoch $t_{P_1} \leftarrow 0$, $\ell_{\text{prv}} \leftarrow 0$, and T_{cur} to a default value.

As pointed out above, in **Signal** and **Signal⁺**, P_1 and P_2 run a CKA protocol in parallel to sending their messages; **Signal** uses a (t, ε) -secure CKA protocol while **Signal⁺** uses a $(t, \varepsilon, +)$ -secure CKA

protocol. To that end, P_1 's first message includes the first message T_1 output by CKA-S. All subsequent messages sent by P_1 include T_1 until some message received from P_2 includes T_2 . At that point P_1 would run CKA-S again and include T_3 with all her messages, and so on (cf. Section B.3).

Upon either sending or receiving T_i for odd or even i , respectively, the CKA protocol also produces a random value I_i , which P_1 absorbs into the root KDF Chain. The resulting output k is used as key for a new FS-AEAD epoch.

- **Sending messages:** Procedure $\text{Send-}P_1$ allows P_1 to send a message to P_2 . As a first step, $\text{Send-}P_1$ determines whether it is P_1 's turn to send the next CKA message, which is the case if t_{P_1} is even. Whenever it is P_1 's turn, $\text{Send-}P_1$ runs CKA-S to produce her next CKA message T and key I , which is absorbed into the root KDF Chain. The resulting value k is used as the key for a new FS-AEAD epoch, in which P_1 acts as sender.

Irrespective of whether it was necessary to generate a new CKA message and generate a new FS-AEAD epoch, $\text{Send-}P_1$ creates a header $h = (t_{P_1}, T_{\text{cur}}, \ell_{\text{prv}})$ (see below for how ℓ_{prv} is computed in Rcv), and uses the current epoch $v[t_{P_1}]$ to compute a ciphertext for (h, m) (where h is treated as associated data).

- **Receiving messages:** When a ciphertext $c = (h, e, \ell)$ with header $h = (t, T, \ell)$ is processed by $\text{Rcv-}P_1$, there are two possibilities:
 - $t \leq t_{P_1}$ (and t even): In this case, ciphertext c pertains to an existing FS-AEAD epoch, in which case FS-Rcv is simply called on $v[t]$ to process e . If the maximum number of messages has been received for session $v[t]$, the session is removed from memory.
 - $t = t_{P_1} + 1$ and t_{P_1} odd: Here, the receiver algorithm advances t_{P_1} by incrementing it and processes T with CKA-R. This produces a key I , which is absorbed into the root KDF chain to obtain a key k with which to initialize a new epoch $v[t_{P_1}]$ as receiver. Then, e is processed by FS-Rcv on $v[t_{P_1}]$. Note that Rcv also uses FS-Max to store ℓ as the maximum number of messages in the previous receive epoch. It also terminates its most recent sending epoch by calling FS-Stop and stores the number of messages in the old epoch in ℓ_{prv} , which will be sent along inside the header for every message of the next sending epoch.

Irrespective of whether a new CKA message was received and a new epoch created, if e is rejected by FS-Rcv , the algorithm raises an exception (**error**), which causes the entire state s_{P_1} to be rolled back to what it was before $\text{Rcv-}P_1$ was called.

C.2 Vulnerability of Signal with Respect to $\mathcal{F}_{\text{Signal}^+}$

Recall that the only difference between Signal and Signal^+ is that the former uses a (t, ε) -secure CKA protocol while the latter uses a $(t, \varepsilon, +)$ -secure CKA protocol. As we will show in the next section, Signal UC-realizes functionality $\mathcal{F}_{\text{Signal}}$, whereas Signal^+ UC-realizes stronger functionality $\mathcal{F}_{\text{Signal}^+}$.

Here, we will explicitly show the vulnerability of Signal with CKA protocol CKA which prevents it from realizing the stronger $\mathcal{F}_{\text{Signal}^+}$ functionality. Briefly recall from Section B.3.2 how CKA protocol CKA works: When P_1 sends a new CKA message, she generates random exponent x then sends message $T_1 \leftarrow g^x$ and computes shared secret $I_1 \leftarrow T_0^x$, where T_0 is the message P_2 sent

in the previous epoch. Then, when P_2 sends the next message, she generates random exponent y and sends message $T_2 \leftarrow g^y$ so that the shared secret for that round is $I_2 \leftarrow T_1^y = g^{xy}$. Thus, P_1 needs to save exponent x after sending T_1 so that when she receives T_2 , she can compute $I_2 = T_2^x$. Therefore, if she is corrupted after she sends T_1 , but before she receives T_2 , then the adversary obtains x , rendering I_1 insecure since the adversary itself can compute $I_1 \leftarrow T_0^x$.

The insecurity of I_1 in the above scenario is at the root of the vulnerability in **Signal**. If P_1 starts a new sending epoch, then she computes CKA message T_1 and secret I_1 as above, then FS-AEAD initialization key k for the epoch as $(\cdot, k) \leftarrow H(\sigma_{\text{root}}, I_1)$, using the current root KDF chain key σ_{root} that she has in her state. She may then proceed to use the FS-AEAD initialized with key k to encrypt and send several messages m_1, \dots, m_ℓ . As highlighted above, before P_1 receives a message for P_2 's next sending epoch, she needs to keep exponent x that she used to generate CKA message T_1 in her state, so that she can compute the CKA shared secret I_2 that P_2 uses to compute the FS-AEAD initialization key for her next epoch. Therefore, consider the following two state leakages of P_1 : (i) before she starts her new epoch and (ii) before she receives a new message for P_2 's next epoch. The first leakage will indeed give the adversary σ_{root} and the second leakage will give the adversary x so that it can first compute I_1 as above, then FS-AEAD initialization key for that epoch as $(\cdot, k) \leftarrow H(\sigma_{\text{root}}, I_1)$, and finally decrypt all of the messages P_1 has sent in her epoch.

Notice that $\mathcal{F}_{\text{Signal}}$ will indeed provide the ideal adversary with these *vulnerable* messages upon the second leakage (the dashed part of Figure 2), while $\mathcal{F}_{\text{Signal}^+}$ will not. Thus, while **Signal** can UC-realize $\mathcal{F}_{\text{Signal}}$, it cannot UC-realize $\mathcal{F}_{\text{Signal}^+}$. **Signal**⁺'s use of CKA⁺ allows it to UC-realize stronger functionality $\mathcal{F}_{\text{Signal}^+}$, because even with both leakages on P_1 described above, the shared secret I_1 remains secure, as long as good randomness is used in the generation of T_0 and T_1 (c.f. Section B.3.4).

We emphasize that while our improved **Signal**⁺ protocol does UC-realize $\mathcal{F}_{\text{Signal}^+}$ (without any additional communication and only small additional computation to that of **Signal**), **Signal** itself should already, as understood by the Double Ratchet whitepaper [MP16a]: As a result of PCS, the messages in the new epoch which P_1 starts should be secure, despite *only* the leakage before she starts the epoch. Moreover, as a result of FS, the messages which P_1 already sent in her new epoch should remain secure despite *only* the second leakage of P_1 . In reality, we however show that PCS and FS of **Signal** with respect to these two leakages no longer hold when *both* of them occur.

D Limitations of ACD's Secure Messaging Security Notion

In this section we demonstrate the following *four* distinct (sometimes contrived) modifications to the **Signal** secure messaging protocol (c.f. Figure 7) which, when instantiated with **correct** and **secure** underlying FS-AEAD and CKA schemes, are vulnerable to natural attacks – formally, we show that they are insecure with respect to our (weaker) ideal functionality $\mathcal{F}_{\text{Signal}}$. Despite this, we show that the modified protocols remain ACD-secure as long as the building blocks are ACD-secure (see Section D.1 for the definitions of ACD-security.)

1. **Postponed FS-AEAD Key Deletion.**
2. **Postponed CKA Key Deletion.**
3. **Eager CKA Randomness Sampling.**
4. **Malleable CKA.**

All transformations discussed in this section satisfy the correctness while affect the security – intuitively breaking either forward security or post-compromise security.

D.1 Definitions from ACD [ACD19]

To formally show that the transformed protocols continue to satisfy ACD’s definitions we need to borrow definitions from ACD. While ACD’s game-based secure message definition is completely different than our simulation based definition (Figure 2), our FS-AEAD and CKA definitions are quite similar to theirs. The differences are discussed in the respective sections (c.f. Section B.2 and Section B.3). To distinguish from our definitions, we refer to the a scheme which is secure according to ACD’s definition as *ACD-secure*. Next we present ACD’s definitions, often taken verbatim from ACD.

D.1.1 Secure Messaging (ACD)

ACD’s game-based definition of secure messaging consists of (potentially asymmetric) initialization algorithms `Init-A` and `Init-B`, a generic sending algorithm `Send`, and a generic receiving algorithm `Rcv`. Each party maintains a state to use across invocations of the sending and receiving algorithms. Importantly, the receiving algorithm additionally outputs an epoch number and message index which is used to determine the order in which the sending party transmitted their messages. We present these algorithms formally in Definition 10.

Definition 10. A secure-messaging (SM) scheme consists of four probabilistic algorithms $SM = (\text{Init-A}, \text{Init-B}, \text{Send}, \text{Rcv})$, where

- `Init-A` (and similarly `Init-B`) takes a key k and outputs a state $s_A \leftarrow \text{Init-A}(k)$,
- `Send` takes a state s and a message m and produces a new state and a ciphertext $(s', c) \xleftarrow{\$} \text{Send}(s, m)$, and
- `Rcv` takes a state s and a ciphertext c and produces a new state, an epoch number, an index, and a message $(s', t, i, m) \leftarrow \text{Rcv}(s, c)$.

Game-Based Secure Messaging. We reproduce the game-based secure messaging security notion from ACD in Figure 8. The security game consists of an initialization procedure `init`, two sending oracles `transmit-P1` (normal transmission) and `chall-P1` (challenge transmission), two receive oracles `deliver-P1` (honest delivery) and `inject-P1` (for forged ciphertexts), and a corruption oracle `corr-P1`. These oracles (with the exception of `init`) are defined with respect to party P_1 . The oracles for P_2 are defined analogously. We remark that, ACD’s definition considers Alice (**A**) and Bob (**B**) instead of parties P_1 and P_2 .

Due to the complexity of capturing secure messaging with game-based definitions, a number of game management functions are provided. These consist of a epoch management function `ep-mgmt`, randomness sampling function `sam-if-nec`, and record keeping functions `record` and `delete` which can be called by the attacker. Additionally, the `safe-chP` and `safe-inj` control the adversary’s ability to make challenges and inject respectively.

The game is parametrized by an integer Δ_{SM} which relates to how fast parties recover from state compromise. The advantage of \mathcal{A} against an SM scheme SM is denoted by $\mathcal{A}_{\text{sm}, \Delta_{SM}}^{\text{SM}}(\mathcal{A})$. The

Security Game for Secure Messaging

<p>init</p> <ul style="list-style-type: none"> $k \xleftarrow{\\$} \mathcal{K}$ $s_{P_1} \leftarrow \text{Init-}P_1(k)$ $s_{P_2} \leftarrow \text{Init-}P_2(k)$ $(t_{P_1}, t_{P_2}) \leftarrow (0, 0)$ $i_{P_1}, i_{P_2} \leftarrow 0$ $t_L \leftarrow -\infty$ $\text{trans, chall, comp} \leftarrow \emptyset$ $b \xleftarrow{\\$} \{0, 1\}$ <p>corr-P_1</p> <ul style="list-style-type: none"> req $P_2 \notin \text{chall}$ $\text{comp} \xleftarrow{\pm} \text{trans}(P_2)$ $t_L \leftarrow \max(t_{P_1}, t_{P_2})$ return s_{P_1} 	<p>transmit-$P_1(m, r)$</p> <ul style="list-style-type: none"> $(r, \text{flag}) \leftarrow \text{sam-if-nec}(r)$ ep-mgmt(P_1, flag) $i_{P_1} ++$ $(s_{P_1}, c) \leftarrow \text{Send}(s_A, m; r)$ record(P_1, norm, m, c) return c <p>chall-$P_1(m_0, m_1, r)$</p> <ul style="list-style-type: none"> $(r, \text{flag}) \leftarrow \text{sam-if-nec}(r)$ ep-mgmt(P_1, flag) req safe-chp_{P_1} and $m_0 = m_1$ $i_{P_1} ++$ $(s_{P_1}, c) \leftarrow \text{Send}(m_b; r)$ record($P_1, \text{chall}, m_b, c$) return c 	<p>deliver-$P_2(c)$</p> <ul style="list-style-type: none"> req $(B, t, i, m, c) \in \text{trans}$ for some t, i, m $(s_{P_1}, t', i', m') \leftarrow \text{Rcv}(s_{P_1}, c)$ if $(t', i', m') \neq (t, i, m)$ win if $(t, i, m) \in \text{chall}$ $m' \leftarrow \perp$ $t_{P_1} \leftarrow \max(t_{P_1}, t)$ delete(t, i) return (t', i', m') <p>inject-$P_2(c)$</p> <ul style="list-style-type: none"> req $(P_2, c) \notin \text{trans}$ and safe-inj $(s_{P_1}, t', i', m') \leftarrow \text{Rcv}(s_{P_1}, c)$ if $m' \neq \perp$ and $(P_2, t', i') \notin \text{comp}$ win $t_A \leftarrow \max(t_{P_1}, t')$ delete(t', i') return (t', i', m')
<p>ep-mgmt (P, flag)</p> <ul style="list-style-type: none"> if $P = P_1$ and t_P even or $P = P_2$ and t_P odd if $\text{flag} = \text{bad}$ and $\neg \text{safe-chp}$ $t_L \leftarrow t_P + 1$ $t_P ++$ $i_P \leftarrow 0$ 	<p>record (P, flag, m, c)</p> <ul style="list-style-type: none"> $\text{rec} \leftarrow (P, t_P, i_P, m, c)$ $\text{trans} \xleftarrow{\pm} \text{rec}$ if $\neg \text{safe-chp}$ $\text{comp} \xleftarrow{\pm} \text{rec}$ if $\text{flag} = \text{chall}$ $\text{chall} \xleftarrow{\pm} \text{rec}$ 	<p>delete (t, i)</p> <ul style="list-style-type: none"> $\text{rec} \leftarrow (P, t, i, m, c)$ for some P, m, c $\text{trans, chall, comp} \xleftarrow{-} \text{rec}$ <p>$\text{safe-chp} := \iff t_P \geq t_L + \Delta_{SM}$</p> <p>$\text{safe-inj} := \iff \min(t_A, t_B) \geq t_L + \Delta_{SM}$</p>
<p>sam-if-nec (r)</p> <ul style="list-style-type: none"> $\text{flag} \leftarrow \text{bad}$ if $r = \perp$ $r \xleftarrow{\\$} \mathcal{R}$ $\text{flag} \leftarrow \text{good}$ return (r, flag) 		

Figure 8: Game-based notion of secure messaging (SM) security from [ACD19]. Oracles correspond to party P_1 of the SM security game for a scheme $SM = (\text{Init-}P_1, \text{Init-}P_2, \text{Send}, \text{Rcv})$. The oracles for P_2 are defined analogously. We note that the syntax above is slightly changed to parties P_1/P_2 from A/B for consistency with our own definitions.

attacker is parameterized by its running time t , the total number of queries q it makes, and the maximum number of epochs q_{ep} it runs for. We define SM security in Definition 11.

Definition 11. A secure-messaging scheme SM is $(t, q, q_{ep}, \Delta_{SM}, \varepsilon)$ -secure if for all (t, q, q_{ep}) -at-

ACD's Security Game for CKA

<pre> init (t^*) $k \xleftarrow{\\$} \mathcal{K}$ $\gamma_0^{P_1} \leftarrow \text{CKA-Init-P}_1(k)$ $\gamma_0^{P_2} \leftarrow \text{CKA-Init-P}_2(k)$ $t_{P_1}, t_{P_2} \leftarrow 0$ $b \xleftarrow{\\$} \{0, 1\}$ corr-P₁ req allow-corr or finished_{P₁} return $\gamma_{t_{P_1}}^{P_1}$ </pre>	<pre> send-P₁ $t_{P_1} ++$ $(\gamma_{t_{P_1}}^{P_1}, T, I) \xleftarrow{\\$} \text{CKA-S}(\gamma_{t_{P_1}}^{P_1})$ return (T, I) send-P₁'(r) $t_{P_1} ++$ req allow-corr $(\gamma_{t_{P_1}}^{P_1}, T, I) \leftarrow \text{CKA-S}(\gamma_{t_{P_1}-1}^{P_1}; r)$ return (T, I) receive-P₁ $t_{P_1} ++$ $(\gamma_{t_{P_1}}^{P_1}, *) \leftarrow \text{CKA-R}(\gamma_{t_{P_1}-1}^{P_1}, T)$ </pre>	<pre> chall-P₁ $t_{P_1} ++$ req $t_{P_1} = t^*$ $(\gamma_{t_{P_1}}^{P_1}, T, I) \xleftarrow{\\$} \text{CKA-S}(\gamma_{t_{P_1}-1}^{P_1})$ if $b = 0$ return $(T_{t_{P_1}}, I_{t_{P_1}})$ else $I \xleftarrow{\\$} \mathcal{I}$ return $(T_{t_{P_1}}, I)$ </pre>
---	---	--

allow-corr_{P₁}, $:\iff \max(t_{P_1}, t_{P_2}) \leq t^* - 2$
finished_{P₁}, $:\iff t_P \geq t^* + \Delta_{\text{CKA}}$

Figure 9: Oracles corresponding to party P₁ of the CKA security game for a scheme CKA = (CKA-Init-P₁, CKA-Init-P₂, CKA-S, CKA-R); the oracles for P₂ are defined analogously; conditions for the weaker security game are presented above those for the stronger.

tackers \mathcal{A} ,

$$\mathcal{A}_{\text{sm}, \Delta_{\text{SM}}}^{\text{SM}}(\mathcal{A}) \leq \varepsilon .$$

Definition 12. A secure-messaging scheme SM is called simply secure with Δ_{SM} if it is $(\text{poly}(\kappa), \text{poly}(\kappa), \text{poly}(\kappa), \Delta_{\text{SM}}, \text{negl}(\kappa))$ -secure.

For more details regarding the secure messaging game-based definition, we refer the reader to Section 3 of the full version of ACD [ACD18].

D.1.2 ACD's CKA Definition

The syntax and the correctness of ACD's CKA is exactly the same as ours (c.f. Definition 6). The security notion, however, diverges as we attempt to capture a more fine-grained notion. Nevertheless, we stick to the game-based notion analogous to ACD. We present ACD's game-based definition in Figure 9, which is adjusted to our notations of parties P₁/P₂ in contrast with ACD's **A/B**. The differences with our CKA definition are discussed in Section B.3. Among them, two main differences are: (i) their definition is a indistinguishability-based definition and hence has a challenge bit in the game, whereas our definition (c.f. Definition 9) is a recoverability-based definition; (ii) their definition is parameterized by a Δ_{CKA} for the recovery period after a corruption whereas we capture recovery in a more fine-grained manner without such parameter. The parameter Δ_{CKA} stands for the number of epochs that need to pass after the challenge epoch t^* until the states do

not contain secret information pertaining to the challenge. Once a party reaches epoch $t^* + \Delta_{\text{CKA}}$, its state may be revealed to the attacker (via the corresponding corruption oracle). The game ends (implicitly captured above) once both states are revealed after the challenge phase. The attacker wins the game if it eventually outputs a bit $b' = b$. The advantage of an attacker A against the above game is denoted by $\mathcal{A}_{\text{ACD}, \Delta_{\text{CKA}}}^{\text{CKA}}(A)$. We define the ACD-security as follows:

Definition 13. A CKA scheme CKA is $(t, \Delta_{\text{CKA}}, \varepsilon)$ -ACD-secure if for all attackers A :

$$\mathcal{A}_{\text{ACD}, \Delta_{\text{CKA}}}^{\text{CKA}}(A) \leq \varepsilon$$

and whenever $t \in \text{poly}(\kappa)$ and $\varepsilon \in \text{negl}(\kappa)$ then we say that the scheme is simply ACD-secure with Δ_{CKA} .

For more details regarding the ACD's CKA definition we refer the reader to Section 4.1 of the full version of ACD [ACD18].

D.1.3 ACD's FS-AEAD Definition

The syntax and the correctness of ACD's FS-AEAD is almost the same as ours (c.f. Definition 3) except that our FS-Stop algorithm is slightly different from theirs. However, similar to ACD we do not require explicit definitions of FS-Stop in this section. The difference comes up later while describing transformation T_1 – we defer this discussion until Section D.2.2. Our security notion diverges here as well although we stick to the game-based notion analogous to ACD. We present ACD's game-based definition of FS-AEAD in Figure 10, which is adjusted to our notations of parties P_1/P_2 in contrast with ACD's \mathbf{A}/\mathbf{B} . We discuss the differences between our definition and ACD's definition in Section B.2. In short, the main difference is that the adversary has much stronger corruption abilities. The advantage of an attacker A against an FS-AEAD scheme FS-AEAD is denoted by the expression $\mathcal{A}_{\text{ACD}}^{\text{FS-AEAD}}(A)$. The attacker is parameterized by its running time t and the total number of queries q it makes.

Definition 14. An FS-AEAD scheme FS-AEAD is (t, q, ε) -ACD-secure if for all (t, q) -attackers A ,

$$\mathcal{A}_{\text{ACD}}^{\text{FS-AEAD}}(A) \leq \varepsilon$$

and whenever $t, q \in \text{poly}(\kappa)$ and $\varepsilon \in \text{negl}(\kappa)$ we simply say that FS-AEAD is ACD-secure.

For more details regarding the ACD's FS-AEAD definition we refer the reader to Section 4.2 of the full version of ACD [ACD18].

D.1.4 ACD's PRF-PRNG Definition

Instead of modelling the function which defines the root KDF chain as a random oracle (as we do), ACD propose a primitive called PRF-PRNG, which resembles both a pseudo-random function (PRF) and a pseudorandom number generator with input (PRNG). They also provide a (deterministic) construction in the standard model for PRF-PRNGs. In Section B.1.1, we explain the insufficiency of a standard model construction of PRF-PRNGs for realizing our $\mathcal{F}_{\text{Signal}}$ ideal functionality. Namely, the adversary's unrestricted ability to leak parties' states requires programming a random oracle to explain these states, and further to provide non-malleability properties (as in CCA-security for hashed ElGamal [ABR01, CS03, KM04]).

It is trivial to see that a PRF-PRNG can be modelled as a random oracle in the ACD composition theorem which we include in the next section. For more details regarding ACD's PRF-PRNG definition, we refer the reader to Section 4.3 of the full version of ACD [ACD18].

ACD's Security Game for FS-AEAD

<pre> init $k \xleftarrow{\\$} \mathcal{K}$ $v_{P_1} \leftarrow \text{FS-Init-S}(k)$ $v_{P_2} \leftarrow \text{FS-Init-R}(k)$ $i_{P_1} \leftarrow 0$ $\text{corr}_{P_1} \leftarrow \text{false}$ $\text{trans, chall, comp} \leftarrow \emptyset$ $b \xleftarrow{\\$} \{0, 1\}$ corr-P₁ $\text{corr}_{P_1} \leftarrow \text{true}$ return v_{P_1} corr-P₂ req $\text{chall} = \emptyset$ end (v_{P_1}, v_{P_2}) </pre>	<pre> transmit-P₁ (a, m) $i_{P_1} ++$ $(v_{P_1}, e) \leftarrow$ $\text{FS-Send}(v_{P_1}, a, m)$ record(good, a, m, e) return e chall-P₁ (a, m_0, m_1) req $\neg \text{corr}_{P_1}$ and $m_0 = m_1$ $i_{P_1} ++$ $(v_{P_1}, e) \leftarrow$ $\text{FS-Send}(v_{P_1}, a, m_b)$ record$(\text{chall}, a, m_b, e)$ return e </pre>	<pre> deliver-P₂ (a, e) req $(i, a, m, e) \in \text{trans}$ for some i, m $(v_{P_2}, i', m') \leftarrow$ $\text{FS-Rcv}(v_{P_2}, a, e)$ if $(i', m') \neq (i, m)$ win if $(i, m) \in \text{chall}$ $m' \leftarrow \perp$ delete(i) return (i', m') inject-P₂ (a, e) req $(a, e) \notin \text{trans}$ $(v_{P_2}, i', m') \leftarrow$ $\text{FS-Rcv}(v_{P_2}, a, e)$ if $m' \neq \perp$ and $i' \notin \text{comp}$ win delete(i') return (i', m') </pre>	
<hr/> <pre> record (flag, a, m, e) $\text{rec} \leftarrow (i_{P_1}, a, m, e)$ $\text{trans} \stackrel{\pm}{\leftarrow} \text{rec}$ if corr_{P_1} $\text{comp} \stackrel{\pm}{\leftarrow} \text{rec}$ if $\text{flag} = \text{chall}$ $\text{chall} \stackrel{\pm}{\leftarrow} \text{rec}$ </pre>			<pre> delete (i) $\text{rec} \leftarrow (i, a, m, e)$ for m, a, e s.t. $(i, a, m, e) \in \text{trans}$ $\text{trans, chall, comp} \stackrel{-}{\leftarrow} \text{rec}$ </pre>

Figure 10: Oracles corresponding to party P_1 of the FS-AEAD security game for a scheme $\text{FS-AEAD} = (\text{FS-Init-S}, \text{FS-Init-R}, \text{FS-Send}, \text{FS-Rcv})$; the oracles for P_2 are defined analogously.

D.1.5 ACD's Composition

Now, deriving from the concrete version of ACD's composition theorem we state the asymptotic version of ACD's composition theorem (adjusted for our notations/syntax and modelling of the PRF-PRNG as a random oracle).

Theorem 4. *Assume that:*

- CKA is a ACD-secure CKA scheme with Δ_{SM}
- FS-AEAD is a ACD-secure FS-AEAD scheme, and
- H is modelled as a random oracle.

Then the Signal protocol (c.f. Figure 7) is ACD-secure with $\Delta_{\text{SM}} = 2 + \Delta_{\text{CKA}}$.

The transformed protocol $T_1(\text{Signal})$

<pre> Init-P₁(k^{P_1}) ($\sigma_{\text{root}}, \gamma$) $\leftarrow k^{P_1}$ $v[\cdot] \leftarrow \perp$ $T_{\text{cur}} \leftarrow \perp$ $\ell_{\text{prv}} \leftarrow 0$ $t_{P_1} \leftarrow 0$ </pre>	<pre> Send-P₁(m) if t_{P_1} is even $\ell_{\text{prv}} \leftarrow \text{FS-Stop}(v[t_{P_1} - 1])$ $t_{P_1} ++$ ($\gamma, T_{\text{cur}}, I$) $\stackrel{\\$}{\leftarrow}$ CKA-S(γ) (σ_{root}, k) $\leftarrow H(\sigma_{\text{root}}, I)$ $v[t_{P_1}] \leftarrow \text{FS-Init-S}(k)$ $h \leftarrow (t_{P_1}, T_{\text{cur}}, \ell_{\text{prv}})$ ($v[t_{P_1}], e$) $\leftarrow \text{FS-Send}(v[t_{P_1}], h, m)$ return (h, e) </pre>	<pre> Rcv-P₁(e) (h, e) $\leftarrow c$ (t, T, ℓ) $\leftarrow h$ req t even and $t \leq t_{P_1} + 1$ if $t = t_{P_1} + 1$ $\ell_{\text{prv}} \leftarrow \text{FS-Stop}(v[t_{P_1}])$ $t_{P_1} ++$ FS-Max($v[t - 2], \ell$) (γ, I) $\leftarrow \text{CKA-R}(\gamma, T)$ (σ_{root}, k) $\leftarrow H(\sigma_{\text{root}}, I)$ $v[t] \leftarrow \text{FS-Init-R}(k)$ ($v[t], i, m$) $\leftarrow \text{FS-Rcv}(v[t], h, e)$ if $m = \perp$ error return (t, i, m) </pre>
---	---	---

Figure 11: Signal-based secure-messaging scheme as described by ACD in Page 30 of [ACD18] where the PRF-PRNG algorithms are replaced with hash functions. The differences with our version of Signal are highlighted in blue: basically the only difference is that FS-Stop is called inside the sending algorithm in ACD instead of inside the receiving algorithm in our protocol. The figure only shows the algorithms for P_1 ; P_2 's algorithms are analogous, with “even” replaced by “odd”.

We will be using this to show security of the transformed protocols. For more details we refer to Section 5.3 of the full version of ACD [ACD18].

D.2 The Transformations to Signal and Their (In)Security

Now we are ready to present the four transformations, each of which we show to be insecure according to our weaker ideal functionality $\mathcal{F}_{\text{Signal}}$, but secure with respect to ACD's definition (Def. 11).

D.2.1 T_1 : Postponed FS-AEAD Key Deletion

In this section we present a transformation that slightly modifies the usage of FS-AEAD scheme in the Signal protocol. This transformation alters the timing of the deletion of FS-AEAD states. In particular, instead of deleting the (old) FS-AEAD secret state (a.k.a. sending chain) when switching from a sending epoch to a receiving epoch, this protocol does it when the next sending epoch is started. The transformed protocol $T_1(\text{Signal})$ is described in Figure 11. Remarkably, ACD's presentation of the “Signal-based Secure-messaging protocol” is actually the same as this transformed protocol; their composition theorem already proved (c.f. Theorem 4) that this protocol is secure according to their definitions. However, it turns out that it is insecure according to our ideal functionality $\mathcal{F}_{\text{Signal}}$. We show this by presenting an “injection attack” which is formalized below.

Lemma 1. *Suppose that the protocol $T_1(\text{Signal})$ is instantiated with a correct and secure CKA*

scheme (with any Δ_{CKA}) and a correct and secure FS-AEAD scheme. Then there exists a PPT adversary against $T_1(\text{Signal})$ protocol in the real world for which there is no PPT simulator that realizes the functionality $\mathcal{F}_{\text{Signal}}$ in the ideal world.

Proof. We propose an explicit attack strategy for an environment and a PPT adversary in the real world. We exploit the fact that in protocol $T_1(\text{Signal})$ (alternatively ACD’s **Signal**-based secure-messaging protocol), parties do not update the sending chain until the next call to **Send**. This means that if a party is compromised within the receiving epoch t (which is in-between two sending epochs $t - 1$ and $t + 1$), it leaks secrets pertaining to the immediately past sending epoch $t - 1$. We formalize this strategy as follows:

1. Once the setup is completed send a message m_1 from P_1 in epoch-1. Get the ciphertext c_1 .
2. Deliver the ciphertext $c_1 = (h, e)$ to P_2 where $h = (t, T, \ell)$.
3. Send a message m_2 from P_2 . Get the ciphertext c_2 in epoch-2
4. Deliver the ciphertext c_2 to P_1 .
5. Compromise P_1 . Get the secret state which includes $v[1]$.
6. Choose an arbitrary m' and then use $v[1]$ to compute $(v[1]', e') \leftarrow \text{FS-Send}(v[1], h, m')$.
7. Send m' on behalf of P_1 using the injection ciphertext $c' = (h, e')$.
8. Finally, deliver c' to P_2 .

Let us now explain **how the attack works and why it cannot be simulated**. Observe that P_2 accepts the *injected* message m' by executing $\text{Rcv-}P_1(c')$ due to the correctness of the underlying schemes. The simulator cannot make this delivery successfully because this message m' was not in the transmission list $P_1.T$ and hence will be skipped by the ideal functionality in the first step of the delivery. This concludes the proof. □

Note that, we **do not need to prove** that $T_1(\text{Signal})$ is secure, when instantiated with secure FS-AEAD and CKA schemes, as in Theorem 4, because this was already proven by ACD.

D.2.2 T_2 : Postponed CKA Key Deletion

We present our second transformation T_2 here, which slightly modifies the CKA scheme in the **Signal** protocol. In particular, the modification holds on to the shared key derived by the CKA-R algorithm for the entire receiving epoch and is deleted only when the next time CKA-S is run. Given any CKA scheme for party P ($\text{CKA-Init-}P, \text{CKA-S}, \text{CKA-R}$) we define the modified scheme ($\text{CKA-Init-}P', \text{CKA-S}', \text{CKA-R}'$) as:

- $\text{CKA-Init-}P'$ is the same as $\text{CKA-Init-}P$
- $\text{CKA-S}'$ takes a state γ and then:
 1. parses γ as $(I_{\text{prv}}, \gamma^*)$ and set $I_{\text{prv}} \leftarrow \perp$;

2. runs $(\gamma', T, I) \stackrel{\S}{\leftarrow} \text{CKA-S}(\gamma^*)$;

and then it outputs (γ', T, I)

• CKA-R' takes inputs (γ, T) and then:

1. run $(\gamma^*, I) \leftarrow \text{CKA-R}(\gamma, T)$;

2. $I_{\text{prv}} \leftarrow I$

3. $\gamma' \leftarrow (I_{\text{prv}}, \gamma^*)$

and outputs (γ', I)

The transformed protocol $T_2(\text{Signal})$ can be obtained just by replacing the CKA algorithms with the modified ones and hence a detailed description is omitted. We now prove that this modified protocol is insecure in our framework.

Lemma 2. *Suppose that the T_2 -modified CKA scheme described above is instantiated with a secure CKA (Def. 9) scheme. Moreover, consider that the protocol $T_2(\text{Signal})$ is instantiated with a secure FS-AEAD scheme (Def. 5). Then there exists a PPT adversary against the $T_2(\text{Signal})$ protocol for which there is no PPT simulator that realizes the functionality $\mathcal{F}_{\text{Signal}}$ in the ideal world.*

Proof. The attacker's strategy is to compromise the sender right before a message is sent and then compromise the receiver in the next epoch. Using the information obtained from both compromises the adversary can recover all messages sent in this epoch – this information can not be simulated in the ideal world to realize the ideal functionality $\mathcal{F}_{\text{Signal}}$. The attacker's strategy works as follows.

1. Send a message m_1 from party P_1 to P_2 (with good randomness) and deliver it to P_2 in epoch-1.
2. Compromise P_2 and obtain σ_{root} .
3. In epoch 2, send a message m_2 from P_2 to P_1 (with good randomness) and deliver it to P_1 . Get the CKA message T and ciphertext $c = (h, e)$.
4. Corrupt P_1 and obtain its state that includes $\gamma = (I_{\text{prv}}, \gamma^*)$.
5. Then recover m_2 as follows:

- $\dots, k \leftarrow \text{H}(\sigma_{\text{root}}, I_{\text{prv}})$;
- $v[2] \leftarrow \text{FS-Init-S}(k)$;
- $\dots, m_2 \leftarrow \text{FS-Rcv}(v[2], h, e)$

We argue **why this attack is possible**. First, observe that the new step added to the protocol (highlighted in blue) stores the previous secret CKA key, I into a variable I_{prv} before this is overwritten (and thus deleted) by $(\gamma, I) \leftarrow \text{CKA-R}(\gamma, T)$ in the Rcv- P_1 (and Rcv- P_2) algorithm. Furthermore, (i) the root-key σ_{root} that is obtained by the first compromise of P_2 is the same as the root key σ_{root} used by the receiver's algorithm Rcv- P_1 and (ii) I_{prv} obtained by the compromise of P_1 is the same as the I used to encrypt m by P_2 due to the correctness of the underlying primitives. Therefore the plaintext can be derived correctly by the attacker.

Finally let us argue **how this can not be simulated**. To see this, observe that ciphertext c is delivered and there are no further corruptions besides the one on P_2 before she sends c , and the one on P_1 after she receives c . The first leak of course does not reveal m_2 to the simulator (since it has not yet been sent). More formally, from the description of $\mathcal{F}_{\text{Signal}}$ we observe that when the epoch changes, that is $\text{New}(P_2, \text{TURN}, P_2.T)$ returns 1, then the flag $P_2.CUR_SLEK$ is reset to 0 in Step 3 of Figure 2 – this is because P_1 and P_2 use good randomness while starting epochs 1 and 2 so that even with $P_1.CUR_SLEK = 1$, both $P_1.PLEK = P_2.PLEK = 0$. Also, the second leak does not reveal m_2 : c is delivered so m_2 is no longer in $P_2.T$ and m_2 is never in $P_1.V$ (it is in $P_2.V$ which is not leaked to the simulator). Thus the simulator only knows the length of m_2 and cannot properly simulate its leakage in the real world. \square

Lemma 3. *Suppose that the T_2 modified CKA scheme CKA' is instantiated with an ACD-secure CKA scheme with $\Delta_{\text{CKA}} \geq 1$. Also assume that $T_2(\text{Signal})$ is instantiated with a secure FS-AEAD scheme. Then $T_2(\text{Signal})$ is a ACD-secure secure-messaging scheme with $\Delta_{\text{SM}} = \Delta_{\text{CKA}} + 2$.*

Proof. We show that if there exists a PPT adversary A with non-negligible advantage against the modified CKA scheme, then we can construct a PPT adversary (or reduction) B that breaks the underlying CKA scheme with non-negligible advantage. Both security games are played with respect to the ACD security game (c.f. Definition 13). Taking a closer look we observe that all queries from A can be simulated by B using the challenger, except a corruption query following a receive query. The state at this time additionally contains I_{prv} . Nevertheless, for any epoch $t \neq t^*$, I_{prv} is just obtained by making a send-query which returns $I = I_{\text{prv}}$. But if the query sequence is $\text{chall-}P_1 \rightarrow \text{receive-}P_2 \rightarrow \text{corr-}P_2$, then I_{prv} is not always obtained (because when $b = 1$ a random value is obtained instead). However, due to the restriction $\Delta_{\text{CKA}} \geq 1$ such corruption query would not return anything because both allow-corr and finished_{P_2} return 0 since $t_{P_2} = t^*$. This concludes the proof. \square

D.2.3 T_3 : Eager CKA Randomness Sampling

This transformation slightly modifies how the CKA scheme is used in the Signal protocol without making any change to the CKA itself. In particular, the modification samples the randomness that is to be used in the CKA-S algorithm of the next sending epoch while it is still in the prior receiving epoch. We note that this transformation reflects the primary description of the Double Ratchet algorithm in its white paper (in the main body) [MP16a]. However, Perrin and Marlinspike later state that better security can be achieved if the randomness is indeed sampled within the actual sending epoch (as in our protocol). We present the modified protocol $T_3(\text{Signal})$ in Figure 12 and prove that this is insecure with respect to our ideal functionality $\mathcal{F}_{\text{Signal}}$.

We now prove the following:

Lemma 4. *Suppose that the protocol $T_3(\text{Signal})$ is instantiated with a secure CKA (Def. 9) scheme and a secure FS-AEAD scheme (Def. 5). Then there exists a PPT adversary against the $T_3(\text{Signal})$ protocol for which there is no PPT simulator that realizes the functionality $\mathcal{F}_{\text{Signal}}$ in the ideal world.*

Proof. We propose an explicit attack strategy for an environment and an adversary in the real world. The adversary never instructs the parties to use bad randomness.

1. Send a message m_1 from P_1 . Get the ciphertext c_1 .

The transformed protocol $T_3(\text{Signal})$

<pre> Init-P₁ (k^{P_1}) ($\sigma_{\text{root}}, \gamma$) $\leftarrow k^{P_1}$ $v[\cdot] \leftarrow \perp$ $T_{\text{cur}} \leftarrow \perp$ $\ell_{\text{priv}} \leftarrow 0$ $t_{P_1} \leftarrow 0$ $r \xleftarrow{\\$} \mathcal{R}$ </pre>	<pre> Send-P₁ (m) if t_{P_1} is even $t_{P_1} ++$ ($\gamma, T_{\text{cur}}, I$) $\leftarrow \text{CKA-S}(\gamma; r)$ (σ_{root}, k) $\leftarrow \text{H}(\sigma_{\text{root}}, I)$ $v[t_{P_1}] \leftarrow \text{FS-Init-S}(k)$ $h \leftarrow (t_{P_1}, T_{\text{cur}}, \ell_{\text{priv}})$ ($v[t_{P_1}], e$) $\leftarrow \text{FS-Send}(v[t_{P_1}], h, m)$ return (h, e) </pre>	<pre> Rcv-P₁ (c) (h, e) $\leftarrow c$ (t, T, ℓ) $\leftarrow h$ req t even and $t \leq t_{P_1} + 1$ if $t = t_{P_1} + 1$ $\ell_{\text{priv}} \leftarrow \text{FS-Stop}(v[t_{P_1}])$ $t_{P_1} ++$ $\text{FS-Max}(v[t-1], \ell)$ (γ, I) $\leftarrow \text{CKA-R}(\gamma, T)$ $r \xleftarrow{\\$} \mathcal{R}$ (σ_{root}, k) $\leftarrow \text{H}(\sigma_{\text{root}}, I)$ $v[t] \leftarrow \text{FS-Init-R}(k)$ ($v[t], i, m$) $\leftarrow \text{FS-Rcv}(v[t], h, e)$ if $m = \perp$ error return (t, i, m) </pre>
---	--	--

Figure 12: The difference from Signal are highlighted in blue.

2. Deliver the ciphertext c_1 to P_2 .
3. Compromise P_2 . Get the secret state which includes $r, \gamma^{P_2}, \sigma_{\text{root}}$.
4. Send a message m_2 from P_2 . Get the ciphertext $c_2 = (h_2, e_2)$.
5. Use γ and r and then compute:
 - (a) $(\dots, I) \leftarrow \text{CKA-S}(\gamma^{P_2}; r)$;
 - (b) $(\dots, k) \leftarrow \text{H}(\sigma_{\text{root}}, I)$;
 - (c) $v[2] \leftarrow \text{FS-Init-R}(k)$;
 - (d) $(\dots, m_2) \leftarrow \text{FS-Rcv}(v[2], h_2, e_2)$;
6. Output m_2 .

It is easy to see that the attack works due to the correctness of the underlying schemes. Furthermore, note that the simulator can not obtain m_2 in the ideal world, because compromising/leaking on P_2 before it sends m_2 will not reveal m_2 (as in the proof of Lemma 2). \square

Lemma 5. *Suppose that $T_3(\text{Signal})$ is instantiated with an ACD-secure CKA scheme with Δ_{CKA} and an ACD-secure FS-AEAD scheme. Then $T_3(\text{Signal})$ is a ACD-secure secure-messaging scheme with $\Delta_{\text{SM}} = \Delta_{\text{CKA}} + 2$*

Proof. We construct a reduction \mathbf{B} against a challenger, which is running the secure-messaging game with the original **Signal** protocol; \mathbf{B} uses an adversary \mathbf{A} who is trying to gain a non-negligible advantage in a secure-messaging game running protocol $T_3(\text{Signal})$. Note that, the only difference in the two schemes is that in the modified scheme the randomness r is sampled one epoch earlier.

Therefore, when a **corr-P** query is received from A followed by a **deliver-P** query, then B must simulate the randomness r , which it uses in the next epoch (when P becomes a sender). B simulates this by sampling a random r and giving that to A. However, in order to do so, it uses the transmit oracle with randomness r : **transmit-P**(m, r). This does not work if A makes a **chall-P** call right after. However, since $\Delta_{\text{SM}} \geq 2$, this can not happen. Hence, the simulation can be done. The other oracles are simulated straightforwardly by using the challenger as there is no other change in the protocol. This concludes the proof. \square

D.2.4 T_4 : Malleable CKA

For our proof that **Signal** UC-realizes $\mathcal{F}_{\text{Signal}}$, we crucially assume that the underlying CKA protocol provides non-malleability guarantees. More specifically, recall that in Definition 7 for CKA in this paper, we include the **test**(t, T, I) oracle, which outputs 1 if the corresponding receiver in epoch t would on input CKA message T output I ; and 0 otherwise. In order to prove that the **Signal** CKA is secure with respect to this definition, we used the assumption that CDH is hard even given a DDH oracle.

In this section, we elaborate on the need for such non-malleability, by first showing that the **Signal** building blocks can be used to build PKE, even if some of the secret information of **Signal** parties is not hidden. Then, we apply T_4 to **Signal**, which replaces the CKA secure with respect to Definition 7, with one that permits the existence of an adversary that given a ciphertext encrypting m from the above PKE scheme, can successfully maul it into a new ciphertext encrypting $m + 1$. Based on this, we show that $T_4(\text{Signal})$ does not UC-realize $\mathcal{F}_{\text{Signal}}$.

We further observe that if our derived PKE scheme is instantiated using the building blocks of the real **Signal** protocol, then this PKE is (a variant of) hashed ElGamal. Hashed ElGamal is known to be CPA-secure based on the DDH assumption, but however, is only known to be CCA-secure under the stronger assumption that CDH is hard given a DDH oracle. Thus such malleability attacks could *in theory* exist against **Signal**. This is why we need to assume CDH-security with a DDH oracle for our proof that **Signal** UC-realizes $\mathcal{F}_{\text{Signal}}$, and not just DDH-security.

On the other hand, **Signal** does not need such non-malleability guarantees to achieve ACD-security, since ACD prove it secure based on a CKA definition without such a **test**() oracle (for the **Signal** CKA, this is proved using only DDH). Thus, generally, if an instantiation of **Signal** with a CKA that prevents against malleability is transformed via T_4 into an instantiation of **Signal** with a CKA that is prone to malleability, then **Signal** is still ACD-secure but is not secure according to $\mathcal{F}_{\text{Signal}}$. More specifically, we highlight a gap in the provable security of **Signal**, based on whether *only* DDH is secure, or if also CDH is secure with a DDH oracle.

PKE from Signal building blocks with (partially) exposed secrets. Here, we observe that PKE can be constructed from the **Signal** building blocks, even if the root KDF chain key σ_{root} is not hidden. Assume that there are some publicly known values σ and h .

- The generation algorithm **Gen**() first samples random CKA initialization key k , then sets $(sk, pk) \leftarrow (\text{CKA-Init-P}_2(k), \text{CKA-Init-P}_1(k))$.
- The encryption algorithm **Enc**(pk, m) first computes $(\cdot, T, I) \leftarrow \text{CKA-S}(pk)$, then $k \leftarrow \text{H}(\sigma, I)$, $v \leftarrow \text{FS-Init-S}(k)$, $(\cdot, e) \leftarrow \text{FS-Send}(v, (T, h), m)$, and finally $c \leftarrow (T, e)$.

- The decryption algorithm $\text{Dec}(sk, c)$ parses $(c_1, c_2) \leftarrow c$ and computes $(\cdot, I) \leftarrow \text{CKA-R}(sk, c_1)$, then $k \leftarrow \text{H}(\sigma, I)$, $v' \leftarrow \text{FS-Init-R}(k)$, and finally $(\cdot, \cdot, m) = \text{FS-Rcv}(v', (c_1, h), c_2)$.

One can observe that if the above PKE scheme is instantiated using the CKA of `Signal`, then the PKE scheme is essentially hashed ElGamal: The secret key is exponent a , the public key is g^a , exponent b (and g^b) is sampled during encryption, then AEAD key k is derived by both the encryptor and decryptor by hashing shared DDH secret g^{ab} (and public information), and finally m is encrypted using the AEAD with key k .

Intuitively, CPA-security of this PKE scheme comes from that of the underlying CKA, since I should be hidden given just the CKA state of the CKA sender before transmitting, and the resulting CKA message that is in the ciphertext (c.f. Section B.3). However, it is unclear if this PKE scheme provides CCA-security, particularly if the underlying CKA does not achieve the security definition with a `test()` oracle—how does the reduction answer decryption queries? In fact, there could, for example, exist some PPT adversary \mathcal{A} that on input ciphertext from our above PKE scheme, $c \leftarrow \text{Enc}(pk, m)$, and public information σ_{root}, h , can with non-negligible probability maul this ciphertext to create c' such that $m + 1 \leftarrow \text{Dec}(sk, c')$.

Malleability attack on `Signal`. Now assume that $T_4(\text{Signal})$ uses a CKA scheme that is secure with respect to the definition of ACD in Definition 13, but for which such an attacker \mathcal{A} above does exist. Then, although $T_4(\text{Signal})$ is ACD-secure with this CKA scheme, it does not realize $\mathcal{F}_{\text{Signal}}$.

Lemma 6. *Assume that malleability adversary \mathcal{A} above exists. Then there exists a PPT adversary against `Signal` for which there is no PPT simulator that realizes $\mathcal{F}_{\text{Signal}}$ in the ideal world.*

Proof. We propose an explicit attack strategy for an environment and an adversary in the real world.

1. After Setup (initialized by P_1), corrupt P_1 . This leaks σ_{root} to the attacker.
2. Send a message m_1 from P_1 . Get the ciphertext $c_1 = (h_1, e_1)$, $h_1 = (t, \ell, T)$.⁸
3. Compute $c'_1 \leftarrow \mathcal{A}((T, e_1), \sigma_{\text{root}}, (t, \ell))$.
4. Deliver c'_1 to P_2 .
5. When P_2 outputs m'_1 , the attacker outputs $m'_1 - 1$.

Now, since \mathcal{A} is successful with non-negligible probability, then with non-negligible probability, P_2 decrypts c'_1 to $m_1 + 1$ and then our attacker outputs m_1 .

However, in the ideal world, the simulator does not get m_1 from the functionality and therefore only with negligible probability can create ciphertext c_1 such that P_2 will decrypt it to $m_1 + 1$ and output that. More formally, from the description of $\mathcal{F}_{\text{Signal}}$, we observe that when P_1 starts the first epoch, although flag `P2.CUR.SLEK` is set to 1 after the first leakage, if P_1 has good randomness then both `P1.PLEK` = `P2.PLEK` = 0, so `P1.CUR.SLEK` is set to 0 and thus the simulator is only given the length of m . \square

⁸We reorder for simplicity; the attack can easily be performed on the original order as presented in ACD.

Intuitively, this attack cannot be executed according to the ACD definition, since such an injection of mauled ciphertext c'_1 after corruption of P_1 is not allowed. More formally, after the corruption of P_1 in the ACD definition, t_L is set to 0. Since for **Signal**, they prove security with respect to $\Delta_{SM} = 3$, t_{P_2} will still be 0 and thus **safe-inj** will evaluate to **false**. Thus, no matter if such a malleability attack can be performed, the same level of security can be proved for **Signal** according to the ACD definition.

E UC Security of **Signal** and **Signal**⁺

In this section, we show that **Signal** and **Signal**⁺ UC-realize $\mathcal{F}_{\text{Signal}}$ and $\mathcal{F}_{\text{Signal}^+}$, respectively, in the $\mathcal{F}_{\text{KE}}^{\text{CKA}}$ -hybrid model.

\mathcal{S}_{t^*}

Notation: The simulator algorithm interacts with the functionality $\mathcal{F}_{\text{Signal}}$ and the hybrid algorithm H_{t^*} . The algorithm initializes lists of *in-transit* ciphertexts $P_{1.T}$, [and *vulnerable ciphertexts* $P_{1.V}$] sent by P_1 to P_2 to ϕ . Analogously, lists $P_{2.T}$ [and $P_{2.V}$] are also initialized to ϕ . The algorithm also initializes leakage flags of both P_1 and P_2 for their corresponding (i) public ratchet secrets: $P_1.\text{PLEK}, P_2.\text{PLEK}$, (ii) current sending epoch symmetric secrets: $P_1.\text{CUR_SLEK}, P_2.\text{CUR_SLEK}$, and (iii) previous sending epoch symmetric secrets: $P_1.\text{PREV_SLEK}, P_2.\text{PREV_SLEK}$, all to 0. Further, it initializes bad-randomness flags $P_1.\text{BAD}, P_2.\text{BAD}$ to 0. Finally, it initializes the turn flag **TURN** as \perp . \mathcal{S}_{t^*} also keeps track of $\mathcal{S}_{t^*}.k_t$; the FS-AEAD initialization key for each epoch t once it is generated by the corresponding sender of that epoch.

- **On input** (sid, **SETUP**, P) **from** H_{t^*} where $P \in \{P_1, P_2\}$: Sample initialization keys $k^P, k^{\bar{P}}$ according to $\mathcal{F}_{\text{KE}}^{\text{CKA}}$ and once \mathcal{A} approves the interaction (i) run $s_P \leftarrow \text{Init-}P_1(k^P), s_{\bar{P}} \leftarrow \text{Init-}P_2(k^{\bar{P}})$, (ii) set **TURN** $\leftarrow P$, and (iii) return $(s_P, s_{\bar{P}})$.
- **On input** (sid, mid, **IN_TRANSIT**, P, $|m|, m'$) **from** H_{t^*} where $P \in \{P_1, P_2\}$:
 1. Set $r \leftarrow \perp$.
 2. If **New**(P, **TURN**, P.T)^a then:
 - (a) Set (i) $\beta \leftarrow \bar{P}.\text{CUR_SLEK}$, (ii) $P.\text{PLEK} \leftarrow P.\text{BAD}$, and (iii) $P.\text{CUR_SLEK} \leftarrow \bar{P}.\text{CUR_SLEK} \wedge (P.\text{PLEK} \vee \bar{P}.\text{PLEK})$.
 - (b) If $P.\text{BAD} = 0$ then sample random $r \xleftarrow{\$} \mathcal{R}$; otherwise ask \mathcal{A} for random r' and set $r \leftarrow r'$.
 3. If $m' = \perp$ then set $m \leftarrow 0^{|m'|}$; otherwise, set $m \leftarrow m'$.
 4. Run $(s_P, c) \leftarrow \text{Send}(s_P, m; r)$, add (sid, mid, $s_P.t, s_P.v[s_P.t].i$, **IN_TRANSIT**, c , **TURN**) to P.T, [and if $\beta \vee (P.V \neq \emptyset)$ then add (s

Finally return (s_P, c) .
- **On input** (sid, **DELIVER**, P, c) **from** H_{t^*} :
 1. Set $t_{\text{prev}} \leftarrow s_{\bar{P}}.t$ and run $(s_{\bar{P}}, t, i, m) \leftarrow \text{Rcv}(s_{\bar{P}}, c)$. If $m = \perp$ then return \perp .
 2. If $t = t_{\text{prev}} + 1$ then set (i) **TURN** $= \bar{P}$, (ii) P.T $\leftarrow \text{Flip}(P, P.T)$,^b (iii) $P.\text{PREV_SLEK} \leftarrow 0$, (iv) $\bar{P}.\text{PREV_SLEK} \leftarrow \bar{P}.\text{CUR_SLEK}$, (v) $\bar{P}.\text{CUR_SLEK} \leftarrow 0$, (vi) $\bar{P}.\text{PLEK} \leftarrow 0$, [and (vii) $\bar{P}.V \leftarrow \emptyset$].
 3. Find (sid, mid, t, i , **IN_TRANSIT**, c, γ) in P.T, remove it from P.T, and return $(s_P, \text{sid}, \text{mid}, \text{DELIVER}, P, m, 0, \perp, \perp)$. If no such entry is found:
 - (a) If $t = t_{\text{prev}} + 1$ then:
 - i. If $\nexists (\text{sid}, \cdot, t, \cdot, \text{IN_TRANSIT}, \cdot, P) \in P.T$ or for $(\text{sid}, \cdot, t, \cdot, \text{IN_TRANSIT}, c', P) \in P.T, (h', e') \leftarrow c', (t', T', \ell') \leftarrow h', T' \neq T$, return $(s_{\bar{P}}, \text{sid}, (t, i), \text{INJECT}, P, m, 1, \perp, \perp)$.
 - ii. Otherwise, return $(s_{\bar{P}}, \text{sid}, (t, i), \text{INJECT}, P, m, 0, P, \ell')$, where ℓ' is as above.
 - (b) Otherwise return $(s_{\bar{P}}, \text{sid}, (t, i), \text{INJECT}, P, m, 0, \perp, \perp)$.

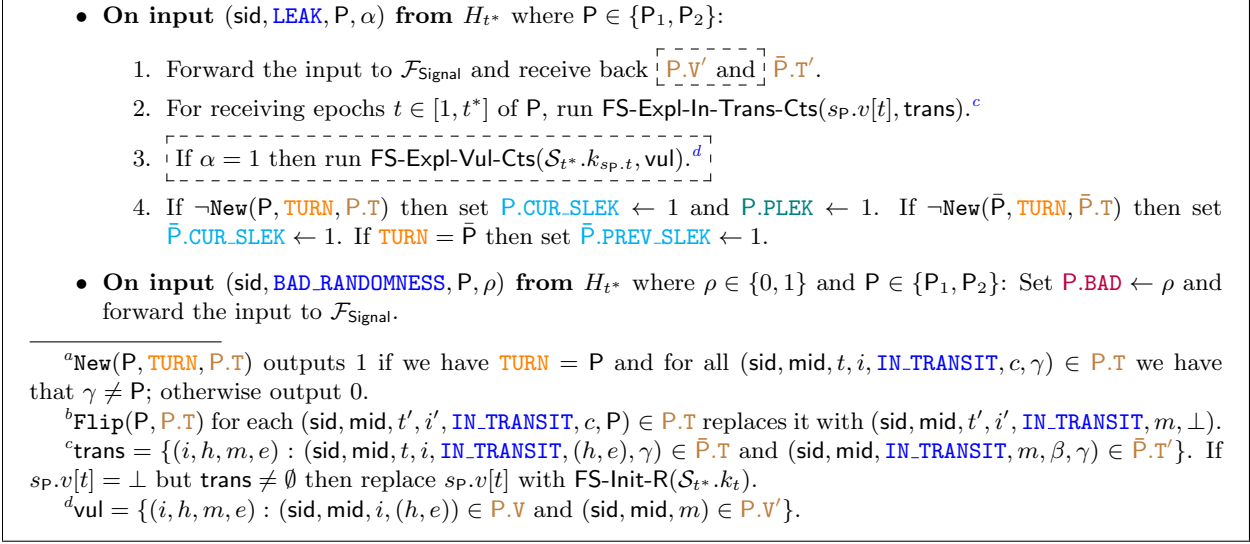


Figure 13: Simulator algorithm \mathcal{S}_{t^*} .

Theorem 5. *Assume that*

- CKA is a ϵ_{CKA} -secure natural CKA scheme,
- CKA^+ is a $(\epsilon_{\text{CKA}}, +)$ -secure natural CKA scheme,
- FS-AEAD is a $\epsilon_{\text{FS-AEAD}}$ -secure FS-AEAD scheme, and
- H is modelled as a random oracle.

Then protocols **Signal** and **Signal**⁺ UC-realize functionalities $\mathcal{F}_{\text{Signal}}$ and $\mathcal{F}_{\text{Signal}^+}$ in the $\mathcal{F}_{\text{KE}}^{\text{CKA}}$ -hybrid and $\mathcal{F}_{\text{KE}}^{\text{CKA}^+}$ -hybrid models, respectively, with security loss $\epsilon = T \cdot (\epsilon_{\text{CKA}} + \epsilon_{\text{FS-AEAD}})$, where T is the number of epochs the attacker initiates.

In the following, we may simply refer to the ideal functionality as $\mathcal{F}_{\text{Signal}}$, but such statements will hold for $\mathcal{F}_{\text{Signal}^+}$ too.

E.1 Hybrid Algorithms H_{t^*} and the Simulator \mathcal{S}

To prove the theorem, we proceed in a series of hybrids H_0, H_1, \dots, H_T , where T is the number of epochs which the adversary initiates. Hybrid H_0 is the real world protocol **Signal** or **Signal**⁺ which interacts with the real world adversary \mathcal{A} . Hybrid H_{t^*} runs stateful simulator algorithm \mathcal{S}_{t^*} for (i) setup, (ii) generation of messages from epochs 1 through t^* using input from $\mathcal{F}_{\text{Signal}}$, (iii) delivery of messages from epochs 1 through t^* , and (iv) for leakages of parties, the explanation of in-transit messages for epochs 1 through t^* and vulnerable messages if the party is not in an epoch after t^* . For generation of messages from epochs after t^* , hybrid H_{t^*} uses the corresponding information written to the input tapes of P_1 and P_2 in the real world and executes as in the real world. Also for delivery of messages for epochs later than t^* and leakages of states at any point, hybrid H_{t^*} behaves as in the real world. Formally, hybrid H_{t^*} (for $t^* > 0$) is defined as follows:

- **On input** $(\text{sid}, \text{SETUP}, P)$ **from** $\mathcal{F}_{\text{Signal}}$ where $P \in \{P_1, P_2\}$: Forward the input to \mathcal{S}_{t^*} and once $(s_P, s_{\bar{P}})$ is received back, send $(\text{sid}, \text{SETUP})$ back to $\mathcal{F}_{\text{Signal}}$.
- **On input** $(\text{sid}, \text{mid}, \text{IN_TRANSIT}, P, |m|, m')$ **from** $\mathcal{F}_{\text{Signal}}$ where $P \in \{P_1, P_2\}$: Forward the input to \mathcal{S}_{t^*} , receive back (s_P, c) , and send c to \mathcal{A} .
- **On input** $(\text{sid}, \text{mid}, \text{SEND}, m)$ **from** \mathcal{Z} to $P \in \{P_1, P_2\}$: Run $(s_P, c) \leftarrow \text{Send}(s_P, m)$ (asking \mathcal{A} for randomness if necessary) and send c to \mathcal{A} .
- **On input** $(\text{sid}, \text{DELIVER}, P, c)$ **from** \mathcal{A} :
 1. Parse $(h, e) \leftarrow c, (t, T, \ell) \leftarrow h$.
 2. If $t \leq t^*$ then forward the input to \mathcal{S}_{t^*} and:
 - (a) If \mathcal{S}_{t^*} returns \perp then skip.
 - (b) Otherwise, parse the return as $(s'_{\bar{P}}, \text{sid}, \text{mid}, \text{INSTRUCTION}, P, m, \delta, \gamma, \ell')$.
 - (c) If $s_{\bar{P}}.t \leq t^*$, then set $s_{\bar{P}} \leftarrow s'_{\bar{P}}$.
 - (d) If $s_{\bar{P}}.t > t^*$ then set $s_{\bar{P}}.v[t] \leftarrow s'_{\bar{P}}.v[t]$.
 - (e) If **INSTRUCTION** = **DELIVER** then forward $(\text{sid}, \text{mid}, \text{DELIVER}, P, m)$ to $\mathcal{F}_{\text{Signal}}$. Otherwise, forward $(\text{sid}, \text{mid}, \text{INJECT}, P, m, \delta, \gamma)$ to $\mathcal{F}_{\text{Signal}}$ and if $\delta = 0$ then when $\mathcal{F}_{\text{Signal}}$ passes activation back to H_{t^*} , send $(\text{sid}, \text{mid}, \text{DELIVER}, P, m)$ to $\mathcal{F}_{\text{Signal}}$. Also, if $\ell' \neq \perp$ then while $\ell' < \ell$: $\ell' ++$ and send $(\text{sid}, (t - 2, \ell'), \text{INJECT}, P, \perp, 0, \perp)$ to $\mathcal{F}_{\text{Signal}}$. If $\delta = 1$ then H_{t^*} will use $s_{\bar{P}}$ to communicate directly with \mathcal{A} on behalf of \bar{P} : It will generate all messages from \bar{P} directly using the normal **Send** algorithm and receive all messages for \bar{P} directly using the normal **Rcv** algorithm (starting with this one).
 3. Otherwise, run $(s_{\bar{P}}, t, i, m) \leftarrow \text{Rcv}(s_{\bar{P}}, c)$ and write $(\text{sid}, (t, i), m)$ to the output tape of \bar{P} .
- **On input** $(\text{sid}, \text{LEAK}, P)$ **from** \mathcal{A} where $P \in \{P_1, P_2\}$:
 1. If $s_P.t \leq t^*$ then set $\alpha \leftarrow 1$; otherwise, set $\alpha \leftarrow 0$.
 2. Forward $(\text{sid}, \text{LEAK}, P, \alpha)$ to \mathcal{S}_{t^*} and send s_P to \mathcal{A} .
- **On input** $(\text{sid}, \text{BAD_RANDOMNESS}, P, \rho)$ **from** \mathcal{A} where $\rho \in \{0, 1\}$ and $P \in \{P_1, P_2\}$: If $s_P.t \leq t^*$ then forward the input to \mathcal{S}_{t^*} .

Observe that H_T is the simulator \mathcal{S} for the ideal world (it simply uses \mathcal{S}_T for everything).

Lemma 7. *If we make the same assumptions as in Theorem 5, then for $t^* \in [T]$, hybrids H_{t^*-1} and H_{t^*} are indistinguishable with security loss $\varepsilon = \varepsilon_{\text{FS-AEAD}} + \varepsilon_{\text{CKA}}$.*

To show H_{t^*-1} and H_{t^*} are indistinguishable, we first make the simplifying assumptions that (i) P_1 initializes the session and (ii) P_1 sends in epoch t^* . The other cases are handled analogously. Now, for those adversaries \mathcal{A} that leak on parties P_1 and P_2 and/or give them bad randomness such that after P_1 sends the first message in epoch t^* , it will be that $P_1.\text{CUR_SLEK} = 1$: it is clear that from the correctness of CKA and FS-AEAD that $H_{t^*-1} \equiv H_{t^*}$, so we are done. We omit an explicit reduction to the security game of FS-AEAD (which captures its correctness guarantees) for brevity. Otherwise, we divide the types of adversaries such that the above does not hold into the three types of the following subsections.

E.2 Type 1 Adversaries

Type 1 adversaries \mathcal{A} are those adversaries such that:

1. Before P_2 has accepted any epoch t^* message, \mathcal{A} successfully forges an epoch t^* message with CKA message T in the header such that either
 - (a) P_1 has not sent any epoch t^* ciphertext yet, or
 - (b) For the CKA message T_{t^*} that is contained in the ciphertexts P_1 has sent in epoch t^* , $T_{t^*} \neq T$;

and \mathcal{A} has not queried the random oracle on $(\sigma_{\text{root}}^{t^*-1}, I)$, where I is the corresponding CKA secret that P_2 would output upon receiving message with T ; or

2. Before P_1 has accepted any epoch $t^* + 1$ message, \mathcal{A} successfully forges an epoch $t^* + 1$ message with CKA message T in the header such that either
 - (a) P_2 has not sent any epoch $t^* + 1$ ciphertext yet, or
 - (b) For the CKA message T_{t^*+1} that is contained in the ciphertexts P_2 has sent in epoch $t^* + 1$, $T_{t^*+1} \neq T$;

and \mathcal{A} has not queried the random oracle on $(\sigma_{\text{root}}^{t^*}, I)$, where I is the corresponding CKA secret that P_1 would output upon receiving message with T .

Lemma 8. *If we make the same assumptions as in Theorem 5, then Type 1 adversaries only succeed with probability $\varepsilon_{\text{FS-AEAD}}$.*

Proof. We provide a reduction to the security of the underlying FS-AEAD scheme FS-AEAD to show that successful Type 1 adversaries only exist with negligible probability. Specifically, assuming that there is some successful Type 1 adversary \mathcal{A} , we construct reduction algorithm $\mathcal{B}_{\text{FS-AEAD},1}$ that has non-negligible probability against FS-AEAD in the FS-AEAD security game, thus reaching a contradiction. $\mathcal{B}_{\text{FS-AEAD},1}$ simulates hybrid H_{t^*-1} or H_{t^*} for \mathcal{A} , using the FS-AEAD security game oracle **init** to initialize the FS-AEAD state of P_2 [or P_1] for injection in epoch t^* [or $t^* + 1$], and oracle **inject- P_2** for the corresponding injections.

More formally, $\mathcal{B}_{\text{FS-AEAD},1}$ first samples random $b \xleftarrow{\$} \{0, 1\}$ and **inj-guess** $\xleftarrow{\$} \{P_1, P_2\}$ corresponding to the party to whom \mathcal{A} will successfully inject a message. It then proceeds as in H_{t^*-1} with the following exceptions:

- **On input** (sid, mid, **SEND**, m) **from** \mathcal{Z} **to** $P \in \{P_1, P_2\}$: If (i) $s_P.t = t^*$ before $(s_{P_1}.v[t^*], e) \leftarrow \text{FS-Send}(s_P.v[t], h, m)$ would normally be executed within **Send**, (ii) $P_1.\text{CUR-SLEK} = 0$, and (iii) $b = 1$; then replace m with $m \leftarrow 0^{|m|}$. Otherwise, proceed as in H_{t^*-1} .
- **On input** (**DELIVER**, P, c) **from** \mathcal{A} :
 1. Parse $(h, e) \leftarrow c$, $(t, T, \ell) \leftarrow h$.
 2. If $t = t^* \wedge T \neq T_{t^*}$ [**inj-guess** = P_2] then query **inject- P_2** (h, e). If there are no more queries of this form then send random $b' \xleftarrow{\$} \{0, 1\}$ to the challenger.

3. $\boxed{\text{If } t = t^* + 1 \wedge T \neq T_{t^*+1} \wedge \text{inj-guess} = \text{P}_1}$ then query **init** followed by **inject-P₂**(h, e).
4. Otherwise proceed as in H_{t^*-1+b} .

Now, before epoch t^* , H_{t^*-1} and H_{t^*} do not diverge and thus we simulate them perfectly. In epoch t^* , before any additional leakage, if $b = 0$ we encrypt the real message m and thus simulate H_{t^*-1} perfectly; if $b = 1$ we encrypt the all 0 message and thus simulate H_{t^*} perfectly. Since \mathcal{A} does not query the random oracle on $(\sigma_{\text{root}}^{t^*-1}, I)$, for I corresponding to T , the corresponding FS-AEAD initialization key k used in both H_{t^*-1} and H_{t^*} is uniformly random and unknown to \mathcal{A} and thus the key sampled by the FS-AEAD challenger is distributed correctly. So, when \mathcal{A} eventually submits a successful forgery, $\mathcal{B}_{\text{FS-AEAD},1}$ will pass it to the challenger and the challenger will declare that $\mathcal{B}_{\text{FS-AEAD},1}$ has won the game, a contradiction.

$\boxed{\text{If } \mathcal{A} \text{ waits until epoch } t^* + 1 \text{ to submit a successful forgery}}$ then to deliver any well-formed epoch t^* message (either honest or an injection for which \mathcal{A} has queried the random oracle on the corresponding (σ, I)), $\mathcal{B}_{\text{FS-AEAD},1}$ acts as in H_{t^*-1} or H_{t^*} according to the sampled bit b . Finally, it will act as above for epoch $t^* + 1$ attempted forgeries and pass the successful forgery to the challenger, a contradiction. \square

E.3 Type 2 Adversaries

Type 2 adversaries \mathcal{A} are those that are not Type 1 and query the random oracle on $(\sigma_{\text{root}}^{t^*-1}, I_{t^*})$ without beforehand:

1. Any leak on P_2 after P_1 sends the first message in epoch t^* , but before P_2 receives any message that causes it to advance to epoch t^* ; or
2. $\boxed{\text{Any leak on } \text{P}_1 \text{ after } \text{P}_1 \text{ sends the first message in epoch } t^*, \text{ but before } \text{P}_1 \text{ receives any message that cau}}$

Lemma 9. *If we make the same assumptions as in Theorem 5, then Type 2 adversaries only succeed with probability ε_{CKA} .*

Proof. We provide a reduction to the security of the underlying CKA scheme CKA to show that successful Type 2 adversaries only exist with negligible probability. Specifically, assuming that there is some successful Type 2 adversary \mathcal{A} , we construct reduction algorithm \mathcal{B}_{CKA} that has non-negligible probability against CKA in the corresponding CKA security game, thus reaching a contradiction. \mathcal{B}_{CKA} simulates hybrid H_{t^*-1} or H_{t^*} for \mathcal{A} , using the CKA security game oracle **init** to initialize the CKA states of P_1 and P_2 , oracles **corr-P₁**, **corr-P₂** to handle leakages of CKA states, oracle **receive-P₁**, **receive-P₂** to handle CKA message reception, and oracles **send-P₁**, **send-P₂** and **send-P₁'**, **send-P₂'** to handle CKA secret and message generation, except for in epoch t^* , in which **chall-P₁** is used instead.

More formally, \mathcal{B}_{CKA} first samples random $b \xleftarrow{\$} \{0, 1\}$ and sets $\text{RODict}[\cdot] \leftarrow \perp$, then proceeds as in H_{t^*-1} , with the following exceptions:

- **On input** (sid, **SETUP**, P) **from** $\mathcal{F}_{\text{Signal}}$: Replace the executions of CKA-Init-P₁ and CKA-Init-P₂ within the executions of Init-P₁ and Init-P₂, respectively, in \mathcal{S}_{t^*} with (i) an oracle query to **init**(t^*), (ii) set $\gamma_{\text{P}} \leftarrow \text{corr-P}_1$, and (iii) if $t^* \neq 1$ then set $\gamma_{\bar{\text{P}}} \leftarrow \text{corr-P}_2$.
- **On input** (sid, mid, **IN_TRANSIT**, P, $|m|, m'$) **from** $\mathcal{F}_{\text{Signal}}$: Replace the execution of CKA-S($\gamma^{\text{P}}; r$) within the Send execution of \mathcal{S}_{t^*} with the following:

1. If $\mathbf{P.BAD} = 0$ then an oracle call to $(T_t, I_t) \leftarrow \mathbf{send-P}$.
 2. Else if $\mathbf{P.BAD} = 1 \wedge s_{\mathbf{P}}.t = t^* - 2$ (before sending) then abort.
 3. Else if $\mathbf{P.BAD} = 1$ then an oracle call to $(T_t, I_t) \leftarrow \mathbf{send-P}'(r')$, where r' is the randomness acquired from \mathcal{A} in \mathcal{S}_{t^*} .
 4. In both non-abort cases: (i) if $s_{\mathbf{P}}.t < t^* - 1$ after sending, set $\gamma_{\mathbf{P}} \leftarrow \mathbf{corr-P}$ and (ii) use the output T_t, I_t of $\mathbf{send-P}$ or $\mathbf{send-P}'$ as in H_{t^*-1} .
- **On input** $(\text{sid}, \text{mid}, \mathbf{SEND}, m)$ **from** \mathcal{Z} **to** $\mathbf{P} \in \{\mathbf{P}_1, \mathbf{P}_2\}$:
 1. If $s_{\mathbf{P}}.t = t^* - 1$ (before sending) then if $\mathbf{P.BAD} = 1$, abort; otherwise, execute \mathbf{Send} normally with the following exceptions:
 - (a) Instead of executing $(s_{\mathbf{P}}.\gamma, T_{t^*}, I) \leftarrow \mathbf{CKA-S}(s_{\mathbf{P}}.\gamma)$, query oracle $T_{t^*} \leftarrow \mathbf{chall-P}$ and set $s_{\mathbf{P}}.\gamma \leftarrow \mathbf{corr-P}$.
Note: if we are analyzing \mathbf{Signal} , then $\mathbf{corr-P}$ will simply exit with no return.
 - (b) Instead of executing $(s_{\mathbf{P}}.\sigma_{\text{root}}, k) \leftarrow \mathbf{H}(s_{\mathbf{P}}.\sigma_{\text{root}}, I)$, simply sample σ and k uniformly at random and set $s_{\mathbf{P}}.\sigma_{\text{root}} \leftarrow \sigma$.
 - (c) If $b = 1 \wedge \mathbf{P}_1.\mathbf{CUR_SLEK} = 0$ then set $m = 0^{|m|}$.
 2. If $s_{\mathbf{P}}.t = t^*$ (before sending) and $b = 1 \wedge \mathbf{P}_1.\mathbf{CUR_SLEK} = 0$ then set $m = 0^{|m|}$.
 3. Proceed as in H_{t^*-1} .
 - **On input** $(\mathbf{DELIVER}, \mathbf{P}, c)$ **from** \mathcal{A} : Parse $(h, e) \leftarrow c, (t, T, \ell) \leftarrow h$ then:
 1. If $t < t^*$ then execute $(\gamma'_{\bar{\mathbf{P}}}, I) \leftarrow \mathbf{CKA-R}(\gamma_{\bar{\mathbf{P}}}, T)$ as in H_{t^*-1} (inside \mathcal{S}_{t^*}) normally. If the state is not rolled back within \mathcal{S}_{t^*} then additionally make an oracle call to $\mathbf{receive-\bar{P}}$ and set $\gamma_{\bar{\mathbf{P}}} \leftarrow \gamma'_{\bar{\mathbf{P}}}$ (otherwise, $\gamma_{\bar{\mathbf{P}}}$ stays the same as before).
 2. If $t = t^* \wedge T = T_{t^*}$ then run $\mathbf{Rcv}(s_{\bar{\mathbf{P}}}, c)$ except instead of executing $(s_{\bar{\mathbf{P}}}. \gamma, I) \leftarrow \mathbf{CKA-R}(s_{\bar{\mathbf{P}}}. \gamma, T)$ and $(s_{\bar{\mathbf{P}}}. \sigma_{\text{root}}, k) \leftarrow \mathbf{H}(s_{\bar{\mathbf{P}}}. \sigma_{\text{root}}, I_{t^*})$:
 - (a) Query $\mathbf{receive-P}_2$ (if $\mathcal{B}_{\mathbf{CKA}}$ has not yet for T_{t^*}),
 - (b) Use k sampled by $\mathcal{B}_{\mathbf{CKA}}$ above for FS-AEAD initialization, and
 - (c) If $s_{\bar{\mathbf{P}}}$ is not rolled back, then set $s_{\bar{\mathbf{P}}}. \gamma \leftarrow \mathbf{corr-P}_2$ and $s_{\bar{\mathbf{P}}}. \sigma_{\text{root}} \leftarrow \sigma$, where σ is that which was sampled by $\mathcal{B}_{\mathbf{CKA}}$ above.
 3. If $t = t^* \wedge T \neq T_{t^*}$ then run $\mathbf{Rcv}(s_{\bar{\mathbf{P}}}, c)$ except instead of executing $(s_{\bar{\mathbf{P}}}. \gamma, I) \leftarrow \mathbf{CKA-R}(s_{\bar{\mathbf{P}}}. \gamma, T)$ and $(s_{\bar{\mathbf{P}}}. \sigma_{\text{root}}, k) \leftarrow \mathbf{H}(s_{\bar{\mathbf{P}}}. \sigma_{\text{root}}, I_{t^*})$; for every $(s_{\bar{\mathbf{P}}}. \sigma_{\text{root}}, I)$ such that $\mathbf{RODict}[(s_{\bar{\mathbf{P}}}. \sigma_{\text{root}}, I)] \neq \perp$:
 - (a) Query $\mathbf{test}(t^*, T, I)$.
 - (b) If the challenger returns 1 then (i) parse $(\sigma, k) \leftarrow \mathbf{RODict}[(s_{\bar{\mathbf{P}}}. \sigma_{\text{root}}, I)]$ (ii) set $s_{\bar{\mathbf{P}}}. \sigma_{\text{root}} \leftarrow \sigma$, (iii) set $s_{\bar{\mathbf{P}}}. \gamma \leftarrow \mathbf{CKA-Der-R}(T, I)$ and (iv) use k for the rest of $\mathbf{Rcv}(s_{\bar{\mathbf{P}}}, c)$.

If no \mathbf{test} query returns 1 then skip.
 4. If $t = t^* + 1 \wedge T = T_{t^*+1}$ then run $\mathbf{Rcv}(s_{\bar{\mathbf{P}}}, c)$ except instead of executing $(s_{\bar{\mathbf{P}}}. \gamma, I) \leftarrow \mathbf{CKA-R}(s_{\bar{\mathbf{P}}}. \gamma, T)$ and $(s_{\bar{\mathbf{P}}}. \sigma_{\text{root}}, k) \leftarrow \mathbf{H}(s_{\bar{\mathbf{P}}}. \sigma_{\text{root}}, I_{t^*+1})$:
 - (a) Query $\mathbf{receive-P}_1$ (if $\mathcal{B}_{\mathbf{CKA}}$ has not yet for T_{t^*+1}),

- (b) Use k that \mathcal{B}_{CKA} naturally computes when sending the first message of epoch $t^* + 1$ for P_2 , and
- (c) If $s_{\bar{\text{P}}}$ is not rolled back, then set $s_{\bar{\text{P}}}. \gamma \leftarrow \mathbf{corr}\text{-}\mathbf{P}_2$ and $s_{\bar{\text{P}}}. \sigma_{\text{root}} \leftarrow s_{\text{P}}. \sigma_{\text{root}}$.
- 5. $\boxed{\text{If } t = t^* + 1 \wedge T \neq T_{t^*+1}}$ then run $\text{Rcv}(s_{\bar{\text{P}}}, c)$ except instead of executing $(s_{\bar{\text{P}}}. \gamma, I) \leftarrow \text{CKA}\text{-R}(s_{\bar{\text{P}}}. \gamma, T)$ and $(s_{\bar{\text{P}}}. \sigma_{\text{root}}, k) \leftarrow \text{H}(s_{\bar{\text{P}}}. \sigma_{\text{root}}, I_{t^*+1})$; for every $(s_{\bar{\text{P}}}. \sigma_{\text{root}}, I)$ such that $\text{RODict}[(s_{\bar{\text{P}}}. \sigma_{\text{root}}, I)] \neq \perp$:
 - (a) Query $\mathbf{test}(t^* + 1, T, I)$.
 - (b) If the challenger returns 1 then (i) parse $(\sigma, k) \leftarrow \text{RODict}[(s_{\bar{\text{P}}}. \sigma_{\text{root}}, I)]$ (ii) set $s_{\bar{\text{P}}}. \sigma_{\text{root}} \leftarrow \sigma$, (iii) set $s_{\bar{\text{P}}}. \gamma \leftarrow \text{CKA}\text{-Der}\text{-R}(T, I)$, and (iv) use k for the rest of $\text{Rcv}(s_{\bar{\text{P}}}, c)$.

If no \mathbf{test} query returns 1 then skip.

- 6. Otherwise, proceed as in H_{t^*-1+b}

- **On input (QUERY, (σ, I)) from \mathcal{A} :** If $s_{\text{P}_1}. t \geq t^*$ and $\mathbf{test}(t^*, T_{t^*}, I) = 1$ then send I to the challenger (and the game ends). Otherwise,
 - 1. If $\text{RODict}[(\sigma, I)] = \perp$ then sample random (σ', k) and set $\text{RODict}[(\sigma, I)] \leftarrow (\sigma', k)$.
 - 2. Return $\text{RODict}[(\sigma, I)]$.

First note that since we model H as a random oracle, while \mathcal{A} has not queried the random oracle on $(s_{\text{P}}. \sigma_{\text{root}}^{t^*-1}, I_{t^*})$, the output $(s_{\text{P}}. \sigma_{\text{root}}^{t^*}, k)$ is always uniformly random to \mathcal{A} . Moreover, if \mathcal{A} does make that query, \mathcal{B}_{CKA} forwards I_{t^*} to the challenger and the game ends. Therefore, \mathcal{B}_{CKA} properly simulates both H_{t^*-1} and H_{t^*} in the view of \mathcal{A} when it samples uniformly random $(s_{\text{P}}. \sigma_{\text{root}}^{t^*}, k)$ in epoch t^* . Before epoch t^* , we send messages just as in the two hybrids and are able to corrupt the sender right afterwards (except for after sending the first message of epoch $t^* - 1$) in the CKA game to obtain the CKA state of the sender by definition. Thus, when the other party receives a message for the first time in an epoch, we are able to use the real state to see how they would act in both hybrids. And if the receiver's state is not rolled back, this must mean that the CKA message is the same as the honestly generated one for that epoch by the sender, and so we can query the $\mathbf{receive}\text{-}\bar{\text{P}}$ oracle to advance their state in the CKA game (since, by definition of Type 2 adversaries, P_1 must make it to t^* without a takeover of either party).

Note that Type 2 adversaries that make \mathcal{B}_{CKA} abort in Step 2 of an $\mathbf{IN_TRANSIT}$ instruction or Step 1 of a \mathbf{SEND} instruction only exist with negligible probability. This is because if $\text{P}_2. \mathbf{BAD} = 1$ before P_2 sends the first message of epoch $t^* - 1$ or $\text{P}_1. \mathbf{BAD} = 1$ before P_1 sends the first message of epoch t^* then for $\text{P}_1. \mathbf{CUR_SLEK}$ to be 0 after P_1 sends, it must be that $\text{P}_2. \mathbf{CUR_SLEK} = 0$ beforehand. Therefore, it must also be the case that after P_2 sent the first message of epoch $t^* - 1$, $\text{P}_2. \mathbf{CUR_SLEK} = 0$, and no leakages on either party occurred in the interim. So, we know from the proof of indistinguishability of H_{t^*-2} and H_{t^*-1} that \mathcal{A} could not have queried the random oracle on $(\sigma_{\text{root}}^{t^*-2}, I_{t^*-1})$ (since if it did then $\text{P}_2. \mathbf{CUR_SLEK}$ would not be 0), and thus $\sigma_{\text{root}}^{t^*-1}$ is uniformly random and unknown to \mathcal{A} . Hence, the probability of a Type 2 adversary later querying the random oracle on $(\sigma_{\text{root}}^{t^*-1}, I_{t^*})$ is negligible. If \mathcal{B}_{CKA} does not abort then in epoch t^* : before any additional leakages, if $b = 0$, \mathcal{B}_{CKA} encrypts m and thus properly simulates H_{t^*-1} ; otherwise, it encrypts $0^{|m|}$ and thus properly simulates H_{t^*} . Also, all message deliveries are directly simulated as in H_{t^*-1+b} .

It is clear that \mathcal{B}_{CKA} handles bad randomness appropriately and further for corruptions, \mathcal{B}_{CKA} can always provide the correct state to \mathcal{A} since Type 2 adversary \mathcal{A} never attempts to leak on P_2

when she is in epoch $t^* - 1$ or P_1 when she is in epoch t^* . It is furthermore clear that once \mathcal{B}_{CKA} gets the random oracle query $(\sigma_{\text{root}}, I_{t^*})$ from \mathcal{A} , **test** will return 1 and \mathcal{B}_{CKA} will win the CKA game, a contradiction. \square

E.4 Type 3 Adversaries

Type 3 adversaries are all other adversaries that are not type 1 or 2.

Lemma 10. *If we make the same assumptions as in Theorem 5, then Type 3 adversaries only succeed with probability $\varepsilon_{\text{FS-AEAD}}$.*

Proof. For Type 3 adversaries, to show $H_{t^*-1} \approx H_{t^*}$, we provide a reduction to the security of the underlying FS-AEAD scheme FS-AEAD. Specifically, assuming that there is some Type 3 adversary \mathcal{A} that distinguishes between hybrids H_{t^*-1} and H_{t^*} with non-negligible probability, we construct reduction algorithm $\mathcal{B}_{\text{FS-AEAD},2}$ that has non-negligible probability against FS-AEAD in the FS-AEAD security game, thus reaching a contradiction. $\mathcal{B}_{\text{FS-AEAD},2}$ simulates hybrid H_{t^*-1} or H_{t^*} for \mathcal{A} , using the FS-AEAD security game oracle **init** to initialize the FS-AEAD states of P_1 and P_2 in epoch t^* ; oracles **corr-P₁**, **corr-P₂**, and **corr-init-key** to handle leakages of FS-AEAD states of P_1 , P_2 , and the initialization key; and oracles **chall-P₁** and **transmit-P₁** to handle FS-AEAD message transmission.

More formally, $\mathcal{B}_{\text{FS-AEAD},2}$ proceeds as in H_{t^*-1} , with the following exceptions:

- **On input** $(\text{sid}, \text{mid}, \text{SEND}, m)$ **from** \mathcal{Z} **to** $P \in \{P_1, P_2\}$:
 1. If $s_P.t = t^* - 1$ before **Send** would normally be executed then replace the execution of $(s_P.\sigma_{\text{root}}, k) \leftarrow H(s_P.\sigma_{\text{root}}, I_{t^*})$ and $s_{P_1}.v[t^*] \leftarrow \text{FS-Init-S}(k)$ with (i) oracle query **init**, (ii) sample uniformly random σ , and (iii) set $s_P.\sigma_{\text{root}} \leftarrow \sigma$.
 2. If $s_P.t = t^*$ before $(s_{P_1}.v[t^*], e) \leftarrow \text{FS-Send}(s_P.v[t], h, m)$ would normally be executed within **Send**, and $\mathcal{B}_{\text{FS-AEAD},2}$ has not yet queried corruption oracles of the FS-AEAD game, then replace its execution with $e \leftarrow \text{chall-P}_1(h, m, 0^{|m|})$ and use the output as in H_{t^*-1} .
 3. If $s_P.t = t^*$ before $(s_{P_1}.v[t^*], e) \leftarrow \text{FS-Send}(s_P.v[t], h, m)$ would normally be called within **Send**, and $\mathcal{B}_{\text{FS-AEAD},2}$ has queried corruption oracles of the FS-AEAD game, then replace its execution with $e \leftarrow \text{transmit-P}_1(h, m)$ and use the output as in H_{t^*-1} .
 4. Otherwise, proceed as in H_{t^*-1} .
- **On input** $(\text{DELIVER}, P, c)$ **from** \mathcal{A} :
 1. Parse $(h, e) \leftarrow c, (t, T, \ell) \leftarrow h$.
 2. If $t = t^*$ then:
 - (a) If $s_{P_2}.t = t^* - 1$ before **Rcv** would normally be called and $T \neq T_{t^*}$ then proceed as in H_{t^*-1} ; otherwise replace the execution of $(s_{P_2}.\sigma_{\text{root}}, k) \leftarrow H(s_{P_2}.\sigma_{\text{root}}, I_{t^*})$ and $s_{P_2}.v[t^*] \leftarrow \text{FS-Init-R}(k)$ with set $s_{P_2}.\sigma_{\text{root}} \leftarrow \sigma$.
 - (b) If $\exists(\text{sid}, \cdot, t^*, \cdot, \text{IN_TRANSIT}, c, \cdot) \in P_1.\mathbf{T}$ then query oracle $(i, m) \leftarrow \text{deliver-P}_2(h, e)$ and use the output as in H_{t^*-1} (including possibly rolling back the state of P_2).
 - (c) Otherwise, query oracle $(i, m) \leftarrow \text{inject-P}_2(h, e)$ and use the output as in H_{t^*-1} (including possibly rolling back the state of P_2).

- (d) If $s_{P_2}.v[t^*] \neq \perp \wedge m \neq \perp$ then additionally execute $s_{P_2}.v[t^*] \leftarrow \mathbf{corr-P}_2$.
3. Otherwise proceed as in H_{t^*-1} (which is the same as in H_{t^*}).

• **On input** $(\text{sid}, \mathbf{LEAK}, P)$ **from** \mathcal{A} :

1. If $(P = P_1 \wedge s_{P_1}.t = t^* \wedge P.V \neq \emptyset) \vee (P = P_2 \wedge s_{P_2}.t = t^* - 1)$ then query $k \leftarrow \mathbf{corr-init-key}$ and program the random oracle so that $(\sigma, k) \leftarrow H(\sigma_{\text{root}}^{t^*-1}, I_{t^*})$.
2. If $P = P_1 \wedge s_{P_1}.t = t^* \wedge P.V = \emptyset$ then query $s_{P_1}.v[t^*] \leftarrow \mathbf{corr-P}_1$.
3. If $P = P_2 \wedge s_{P_2}.t = t^*$ then query $s_{P_2}.v[t^*] \leftarrow \mathbf{corr-P}_2$.
4. Otherwise proceed as in H_{t^*-1} .

• **On input bit** b **from** \mathcal{A} : Forward b to the challenger.

We now show that when the challenge bit b of the FS-AEAD security game is 0, $\mathcal{B}_{\text{FS-AEAD},2}$ properly simulates H_{t^*-1} and when it is 1, $\mathcal{B}_{\text{FS-AEAD},2}$ properly simulates H_{t^*} . Since \mathcal{A} is a Type 3 adversary, she either does not query the random oracle on $(s_{P_2}.v^{t^*-1}, I_{t^*})$, or before she makes such a query she leaks on P_2 when $s_{P_2}.t = t^* - 1$ or leaks on P_1 when $s_{P_1}.t = t^*$. Before such a leakage (and thus such a random oracle query), in both H_{t^*-1} and H_{t^*} the corresponding output $(s_{P_2}.v^{t^*}, k)$ is uniformly random and unknown to \mathcal{A} . Thus implicitly using uniformly random k for FS-AEAD initialization via the challenger and using randomly sampled σ for the updated σ_{root} of epoch t^* properly simulates both hybrids. Furthermore, of course encrypting in epoch t^* before this point using the FS-AEAD $\mathbf{chall-P}_1()$ oracle of course properly simulates H_{t^*-1+b} , where b is the challenge bit of the FS-AEAD security game. If such a leakage does happen then we properly program the random oracle on $(\sigma_{\text{root}}^{t^*-1}, I_{t^*})$ if needed, and properly explain the FS-AEAD states and ciphertexts using FS-Expl-In-Trans-Cts and FS-Expl-Vul-Cts.

It is clear that delivery of epoch t^* messages by $\mathcal{B}_{\text{FS-AEAD},2}$ simulates H_{t^*-1} properly. If \mathcal{A} unsuccessfully forges an epoch t^* message (before P_2 accepts a message for the epoch) then the behavior of H_{t^*-1} and H_{t^*} are identical (both reject the message). Since \mathcal{A} is a Type 3 adversary, if she does successfully forge the first message that P_2 receives in epoch t^* with CKA message $T \neq T_{t^*}$, then it must be that \mathcal{A} queried the random oracle on $(\sigma_{\text{root}}^{t^*-1}, I)$, where I is the corresponding CKA secret associated with T . As in the proof of Lemma 9, in order for \mathcal{A} to do this with non-negligible probability, then it must be that of course $P_1.\mathbf{TAKEOVER_POSS} = 1$; for otherwise, after P_2 sent the first message of epoch $t^* - 1$, $P_2.\mathbf{CUR_SLEK} = 0$ and so $\sigma_{\text{root}}^{t^*-1}$ would be uniformly random and unknown to \mathcal{A} . Thus such forgeries are also successful in H_{t^*} . Furthermore, if P_2 accepts the first message of epoch t^* with proper CKA message T_{t^*} , then the reduction simulates delivery for H_{t^*} properly due to the underlying correctness and security of FS-AEAD: namely, messages are correctly decrypted with the correct index i , only one message for each index i successfully decrypts (even with injections), and injections/modifications of in-transit messages can only happen if the FS-AEAD secret state is leaked.

If the challenger declares **win** or \mathcal{A} guesses the challenge bit b and $\mathcal{B}_{\text{FS-AEAD},2}$ passes it along, then $\mathcal{B}_{\text{FS-AEAD},2}$ wins the FS-AEAD security game, a contradiction. \square

Proof of Lemma 7. Follows immediately from Lemmas 8, 9, and 10. \square

Proof of Theorem 5. Follows immediately from Lemma 7. \square

E.5 Even Stronger Security of Signal⁺ with CKA⁺

We observe that Signal⁺ instantiated with our specific CKA scheme CKA⁺ actually has even *stronger* security, but do not write $\mathcal{F}_{\text{Signal}^+}$ to reflect it, in order to tame the ideal functionality’s complexity. Roughly, in $\mathcal{F}_{\text{Signal}^+}$ (Figure 2), if P₂ starts a new epoch $t^* - 1$ with bad randomness, then the functionality sets $\text{P}_2.\text{PLEK} \leftarrow 1$. Next, if P₁ receives a message in P₂’s new epoch $t^* - 1$ and then the adversary leaks her state, the functionality of course sets $\text{P}_2.\text{CUR_SLEK} \leftarrow 1$. So, when P₁ starts new epoch t^* , the functionality will set $\text{P}_1.\text{CUR_SLEK} \leftarrow 1$, since both $\text{P}_2.\text{CUR_SLEK}$ and $\text{P}_2.\text{PLEK}$ are 1. Therefore, all messages of epoch t^* are deemed insecure by $\mathcal{F}_{\text{Signal}^+}$.

However, recall the stronger security of CKA⁺, which we informally highlighted in Section B.3.5. If we assume secure initialization, good randomness when P₁ sends the first message of epoch 1, no leakage except that of epoch $t^* - 1$ on P₁ above, and again good randomness when P₁ sends the first message of epoch t^* , then the CKA secret I_t for every epoch t remains secure. This is true even if P₂ always uses bad randomness, and besides from epochs 1 and t^* , P₁ also always uses bad randomness. Thus, the symmetric ratchet of every epoch except $t^* - 1$ remains secure (in-transit messages at the time of corruption are also insecure), and so notably the above messages of epoch t^* that $\mathcal{F}_{\text{Signal}^+}$ deems insecure, are indeed secure.

We remark that if P₂ is leaked before receiving a message for epoch $t^* - 1$, then all security is lost: the adversary by correctness learns I_{t^*-1} and can thus also compute the CKA state of P₂, $\gamma_{t^*-1}^{\text{P}_2} = x_{t^*-1} \cdot \text{H}(I_{t^*-1})$, where x_{t^*-1} is the CKA exponent that P₂ samples for epoch $t^* - 1$. So, along with $\sigma_{\text{root}}^{t^*-1}$, the adversary will also have all future CKA shared secrets, and so all future messages will be insecure.

However, in the (not completely far-fetched) scenario in which initialization is secure, P₂ never has good randomness, and P₁ is leaked *after* receiving an epoch $t^* - 1$ message; if P₁ at least uses good randomness for epochs 1 and t^* , then all messages remain secure. On the other hand, in Signal, all CKA secrets I_t are of course leaked in this situation, so once P₁ is corrupted, all security is lost.

We also note that for *further* stronger security, after P sends the first message of an epoch t , instead of setting $\gamma_t^{\text{P}} \leftarrow x_t \cdot \text{H}(I_t)$, where x_t is the sampled exponent of epoch t , one could set $\gamma_t^{\text{P}} \leftarrow x_t \cdot \text{H}(\sigma_{\text{root}}^{t-1}, I_t)$ (or just expand the root KDF computation $(\sigma_{\text{root}}^t, k) \leftarrow \text{H}(\sigma_{\text{root}}^{t-1}, I_t)$ that is already in Signal⁺). Intuitively, this should provide even stronger security, but we did not find any simple examples where security is stronger, other than the above.

F The Model in Detail

In this section we provide details of the model, a lot of which is taken from Bitanksy et al. [BCH12].

F.1 UC Security: A Brief Overview

We summarize the UC security framework of Canetti [Can01]. For brevity and simplicity, we describe a somewhat restricted variant á la [BCH12]; still, the summary is intended to provide sufficient detail for verifying the treatment in this work. The description below is taken almost verbatim from [BCH12]. Further elaboration and justification of definitional choices appears in [Can01].

F.1.1 The Basic Model of Computation

We first present the underlying model of computation, which provides the basic mechanics on top of which the notion of protocol security is defined.

Interactive Turing Machines (ITM). The basic computing element is an Interactive Turing Machine (ITM), which represents a program written for a distributed system. The UC framework uses a formalism of an ITM that augments the original formalism of [Gol04, GMR89] with some additional structure, for the purpose of capturing protocols in multi-party, multi-instance systems. Specifically, an ITM is a Turing machine (one may consider any standard definition such as [Sip97]) with the following additional constructs. It has three *special tapes* that represent three different types of information coming from external sources: (i) The *input tape* represents information coming from the “calling protocol” (for now consider a protocol to be just another ITM, we shall formalize it below); (ii) the *communication tape* represents information coming from other parties over untrusted communication links; (iii) the *subroutine output tape* represents information coming from “subroutine protocols” in a trusted way. In addition, an ITM has a special *identity tape* which cannot be written on by the ITM transition function, or program. The contents of the identity tape is interpreted as three values: The program of the ITM, represented in some canonical form; a session-identifier (sid), representing a specific protocol session; and a party identifier (pid), representing an identity of a party within that session. In the context of this work we will use the pid as an indicator of the physical device on which the ITI runs. That is, all the ITIs that represent processes that run on the same physical device (and can thus leak in a correlated way) have the same pid. Finally, to the standard ITM syntax we add the ability to perform an *external write* instruction. The semantics of this instruction are defined below.

Systems of ITMs. Running programs in a distributed system is captured as follows. An *ITM instance* (also called an ITI) $\mu \leftarrow (M, \text{id})$ is an ITM M (alternatively, a program) along with a string $\text{id} \leftarrow (\text{sid}, \text{pid})$, called the identity of μ . An ITI represents a running instance of the ITM M where the identity id is written in its identity tape. A *system of ITMs* is a pair (M, C) where M is an ITM and $C : \{0, 1\}^* \rightarrow \{0, 1\}^*$ is a control function that determines the effect of the external write commands.

An *execution of a system* (M, C) of ITMs, on input x , consists of a sequence of activations of ITIs. Initially, the system consists of a single ITI with ITM M , some fixed identity (say, $\text{id}_0 = (0, 0)$), and x written on the input tape. This ITI, $\mu_0 \leftarrow (M, \text{id}_0)$ is then activated.

In each activation of an ITI, the active ITI runs its program. The execution ends when the initial ITI μ_0 halts. The output of the execution is the output of the initial ITI.

It remains to specify the effect of the external-write operation of an arbitrary ITI μ^* , that is active at any given point. This operation specifies three things: (i) a target ITI μ (say); (ii) one of its tapes $\tau \in \{\text{input}, \text{communication}, \text{subroutine output}\}$; (iii) and the data δ to be written. When an external-write operation is carried out, the control function C is applied to the sequence of external write requests in the execution so far. Then:

1. If C returns 1 then:
 - (a) If an ITI with the same identity as the target ITI does not exist in the system then a new ITI with the given specification μ .

(b) The data δ is written to the specified tape τ of ITI μ (this is uniquely determined).

The active ITI μ^* becomes inactive and the target ITI μ becomes active.

2. If C returns 0 or the ITI μ^* halts, then the initial ITI μ_0 is activated.
3. If C returns another value, parse that as a description of an ITM M . The effect is the same as in Case-1 except, the program of μ is replaced by M instead of the original value specified in the command.

Subroutines. An ITI μ_{sub} is a *subroutine* of another ITI μ_{main} in an execution if μ_{main} wrote to the input tape of μ_{sub} or μ_{sub} wrote to the subroutine output tape of μ_{main} .

Protocols and protocol instances. A protocol π is formalized as a *single ITM*, that represents the programs to be run by all the intended participants. A protocol may specify different roles, in that case the single ITM describes the programs for all the roles and the role is given to the specific party as part of an input. For example, in a secure message transmission protocol, the sender and the receiver has different roles and therefore run different programs. An *instance* (or session) of a protocol π with session identifier sid , within a system of ITMs, is the set of ITIs that run the program π and whose session identifier is sid .

Polynomial Time ITMs. We consider ITMs that run in probabilistic polynomial time (PPT), where PPT is defined as follows: an ITM M is PPT if there exists a constant $c > 0$ such that, for any ITI μ with program M , at any point during its run, and for any contents of the random tape, the overall number of steps taken is at most n^c , where n is the overall number of bits written on the input tape of μ minus the overall number of bits written by μ to input tapes of other ITIs. An execution of a system of ITM is said to be run in PPT if the initial ITM is PPT and the control function is computable in probabilistic polynomial time by any TM.

F.1.2 Security of Protocols.

Recall that real world protocols that securely implement a given task are defined via comparison with an ideal process for carrying out the task. Formalizing this notion is done in several steps, as follows. First, we define the process of executing a protocol in the presence of an adversarial environment. We then define what it means for one protocol to “emulate” another protocol. Next, we define the ideal functionality for carrying out the task. A protocol is said to securely carry out the task if it emulates the idealized protocol for that task.

The model for protocol execution. The model for executing a protocol π is parametrized by a security parameter $\kappa \in \mathbb{N}$, and three ITMs: the protocol implementation π an adversary \mathcal{A} , which represents the adversarial activity against a *single instance* of π , and an environment \mathcal{Z} , which represents the rest of the system. Specifically, to *execute the protocol* π on input x , execute the system of ITMs $(\mathcal{Z}, C_{\mathcal{A},\pi})$, $C_{\mathcal{A},\pi}$ being the control function for the protocol π in presence of an adversary \mathcal{A} — it remains to describe this control function, namely the external write capabilities of each ITI.

In essence, the definition of $C_{\mathcal{A},\pi}$ captures a model where a *single instance* of π interacts with \mathcal{Z} and \mathcal{A} , in that \mathcal{Z} controls the inputs to parties and reads the outputs. All communication (via

the communication tapes) *must pass through* \mathcal{A} . In addition, the parties of π can create subroutine ITIs, can write to the input tapes of the subroutines, and receive outputs from the subroutine on the subroutine output tapes of the calling parties. More precisely:

- *External writes by the environment:* The environment can write only to the tapes of other ITIs. The program of the first ITI invoked by the environment is set (by the control function) to be the program of the adversary \mathcal{A} . The programs of all the other ITIs that the environment writes to are set to be the protocol π . In addition, the SIDs of all the ITIs invoked by the environment (except \mathcal{A}) must be the same, which implies that all of those ITIs with the same sid belong to the *same instance of* π .
- *External writes by the adversary:* The adversary can write only to the communication tapes of ITIs. In addition, it is *not allowed* to create new ITIs; namely, if the adversary performs an external write request with non-existing target ITI, the control function returns 0.
- *External writes by other ITIs:* An ITI μ other than the environment and the adversary can write only to: (i) the subroutine output tapes of ITIs that have previously written to the input tape of μ ; (ii) the input tapes of ITIs that μ has invoked; (iii) and to the input tapes of ITIs with the same session ID as μ . In addition, it can write to the communication tape of the adversary. (Writing to input tapes of ITIs is the same SID will become useful when defining ideal protocols.)

We also use the convention that creation of a new ITI must be done by writing to the input tape of that ITI; the data written in this activation must start with 1^κ , where κ is the security parameter.

Modeling party corruption specific to our setting. Since the modeling of party corruption will be central to our modeling of *leakage* and *bad randomness*, we describe it in more detail. Formally, party corruption is modeled as a special message sent by the adversary to the corrupted party (ITI). Different types of corruptions (e.g., passive, Byzantine) are modeled as parameters in the corruption message. The response of a party (ITI) to an incoming corruption message is formally treated as part of the protocol specification. This modeling has the advantage that general notions and theorems such as UC emulation, the UC theorem, and the universality of the dummy adversary apply regardless of the specific corruption model. However, some additional formalism is necessary in order to make sure that the formal corruption operation corresponds to the generally accepted intuitive notion of party corruption. Specifically, we assume that an ITM is *corruption compliant* if its program consists of a main program σ (which can be thought of as an “operating system” of sorts), and a subroutine π which represents the actual program run by the ITM. The main program relays all inputs, incoming messages, and subroutine outputs to π , with the exception of the corruption messages sent by the adversary. The behavior of σ upon receipt of the corruption message essentially determine the corruption model.

Let us specify the behavior of σ for two salient types of corruption. In the case of passive party corruption, σ behaves as follows. When an ITI μ receives the first corruption message, the σ part of the code of μ reports that μ has been corrupted to all the ITIs that have written on μ 's input tape. Upon receipt of all other corruption messages, σ returns to the adversary the entire current state of π . Note that π is never notified of the corruption message; this captures the intuitive concept that a party is generally *not aware of being passively corrupted*. Also, note that here the adversary has to explicitly ask for each new report of internal state; however this formalism is chosen for

convenience only and is of no real consequence. Our **leakage** (alternatively state-compromise) is modeled as a restricted form of passive corruption, in that σ returns the secret-state only once; in particular, unlike the standard passive corruption this is just a one-time affair in which the future states of μ are not returned. Nevertheless, leakage requests can be made multiple time, each of which is addressed with a one-time state-return (plus reporting to the other ITIs that have written to μ 's input tape).

In the case of Byzantine/malicious/active corruptions, it is assumed that the corruption message from the adversary includes in it a description of an ITM M . Here σ behaves the same as before, except that it immediately replaces the code M instead of π . Our **bad-randomness** (or adversarially chosen randomness) corruption can be thought of as a restricted form of malicious corruption, in that the adversary instead provides description of a specific ITM, which runs the code of π except that now the code does not read from its own randomness tape, but takes the adversarial randomness provided as part of the ITM. This is also notified to all ITIs that have written to μ 's input tape. Such bad randomness is used until the adversary sends another corruption request with an ITM that runs exactly the code of π reading from its own randomness tape – again this is reported to all ITI's that have written to μ 's input tape.

Let $\text{EXEC}_{\pi, \mathcal{A}, \mathcal{Z}}(\kappa, z)$ denote the output distribution of the environment \mathcal{Z} when interacting with parties running protocol π on security parameter κ and input z . Let $\text{EXEC}_{\pi, \mathcal{A}, \mathcal{Z}}$ denote the ensemble $\{\text{EXEC}_{\pi, \mathcal{A}, \mathcal{Z}}(\kappa, z)\}_{\kappa \in \mathbb{N}, z \in \{0,1\}^*}$.

Protocol emulation. Informally, we say that a protocol π *UC-emulates* protocol π' if for any adversary \mathcal{A} there exists an adversary \mathcal{A}' such that no environment \mathcal{Z} , on any input, can tell with non-negligible probability whether it is interacting with \mathcal{A} and parties running π or it is running π' . This means that, from the view of the environment, running protocol π is “just as good” as interacting with π' .⁹ This notion is formalized as follows. A distribution ensemble is called binary if it consists of distributions over $\{0,1\}$. We have:

Definition 15. *Two binary distribution ensembles $\{X(\kappa, a)\}_{\kappa \in \mathbb{N}, a \in \{0,1\}^*}$ and $\{Y(\kappa, a)\}_{\kappa \in \mathbb{N}, a \in \{0,1\}^*}$ are called indistinguishable (written $X \approx Y$) if for any $c, d \in \mathbb{N}$ there exists $\kappa_0 \in \mathbb{N}$ such that for all $\kappa > \kappa_0$ and for all $a \in \{0,1\}^{\kappa^d}$ we have:*

$$|\Pr[X(\kappa, a)] - \Pr[Y(\kappa, a)]| < \kappa^{-c}$$

Definition 16 (Protocol emulation). *Let π and π' be two protocols. We say that π UC emulates π' if for any adversary \mathcal{A} there exists an adversary \mathcal{A}' such that for any environment \mathcal{Z} that outputs a value in $\{0,1\}$ we have:*

$$\text{EXEC}_{\pi, \mathcal{A}, \mathcal{Z}} \approx \text{EXEC}_{\pi', \mathcal{A}', \mathcal{Z}}$$

This work makes use of the following simplified formulation of UC emulation. Let the *dummy adversary* \mathcal{D} be the adversary that merely reports to the environment all the messages sent by the parties, and follows the instructions of the environment regarding which messages to deliver to parties. Then, it is enough to prove security with respect to the dummy adversary. That is:

⁹To be precise, the definition of protocol emulation only quantifies over balanced environments. An environment is balanced if at any point in time the overall length of input to the adversary is at least some polynomial in the overall length the of inputs given to the rest of the ITIs in the system. As explained in [Can01], failing to make this restriction makes the definition unreasonably strong, and also causes technical problems with the composition theorem.

Definition 17 (Protocol emulation with the dummy adversary). *Let π and π' be two protocols. We say that π UC emulates π' with the dummy adversary if there exists an adversary \mathcal{A}' such that for any environment \mathcal{Z} that outputs a value in $\{0, 1\}$ we have:*

$$\text{EXEC}_{\pi, \mathcal{D}, \mathcal{Z}} \approx \text{EXEC}_{\pi', \mathcal{A}', \mathcal{Z}}$$

Claim 1 ([Can01]). *Protocol π UC-emulates protocol π' iff π UC-emulates π' with respect to dummy adversary.*

Ideal functionalities and ideal protocol. A key ingredient in the ideal process for a given task is the *ideal functionality* that captures the desired behavior, or in other words, the specification of that task. The ideal functionality is modeled as an ITM (representing a “trusted party”) that interacts with the parties and the (ideal) adversary. For convenience, the process for realizing an ideal functionality is represented as a special type of protocol, called an *ideal protocol* denoted $I_{\mathcal{F}}$. The adversary interacting with the ideal protocols is called the simulator and denoted \mathcal{S} . *Executing the ideal protocol $I_{\mathcal{F}}$* on input x formally means executing the system of ITMs $(\mathcal{Z}, C_{\mathcal{S}, I_{\mathcal{F}}}), C_{\mathcal{S}, I_{\mathcal{F}}}$ being the control function for the protocol $I_{\mathcal{F}}$ in presence of the adversary \mathcal{S} . Informally, the execution works of the ideal protocol $I_{\mathcal{F}}$ works as follows: all parties *simply hand their inputs* to an ITI with program \mathcal{F} plus a session ID that is equal to the local session ID, and party ID set to some fixed value, say \perp . Whenever a party in $I_{\mathcal{F}}$ receives a value from \mathcal{F} on its subroutine output tape, it immediately copies this value to the subroutine output tape of the ITI that invoked it. We call the parties of the ideal protocol *dummy parties*.

Definition 18 (Realizing functionalities). *Let π be a protocol, and let \mathcal{F} be an ideal functionality. We say that π UC-realizes \mathcal{F} if π UC-emulates $I_{\mathcal{F}}$, the ideal protocol for \mathcal{F} .*

Ideal functionalities and party corruption. An ideal functionality represents an ideal specification, rather than an actual program that runs on an actual, physical device. Thus, party corruption messages sent to an ideal functionality do not directly represent physical corruption. Instead, the behavior of an ideal functionality upon receipt of corruption messages from the adversary specifies the security requirements from the realizing protocols upon party corruption.

In general, ideal functionalities can modify their behavior in arbitrary ways as a function of the corruption requests received from the adversary so far. (For instance, an ideal functionality may allow the adversary to modify sensitive information as soon as more than some number of parties have been corrupted.) Still, we define some “standard” behavior of an ideal functionality in face of corruption. Specifically, we say that an ideal functionality \mathcal{F} is *standard corruption* if:

1. An instance of \mathcal{F} with SID sid keeps some “ideal local state” state_P for each dummy party (sid, P) that interacts with this instance of \mathcal{F} . (Here P is the PID of this dummy party.)
2. Upon receipt of the first “corrupt P” message from the adversary, \mathcal{F} first notifies the dummy party (sid, P) that it has been corrupted. Next, in each future “corrupt p” message, \mathcal{F} returns to the adversary the contents of the ideal state state_P .

It is stressed that a standard corruption functionality can specify additional instructions to be performed upon receipt of a corruption message; it can also alter its overall behavior as exemplified above.

Universal Composition. Let ρ^ϕ be a protocol that uses one or more instances of some protocol ϕ as a subroutine, and let π be a protocol that UC-emulates ϕ . The composed protocol ρ^π is constructed by modifying the program of ρ^ϕ so that calls to ϕ are replaced by calls to π . Similarly, subroutine outputs coming from π are treated as subroutine outputs coming from ϕ . The universal composition theorem says that protocol ρ^π behaves essentially the same as the original protocol ρ^ϕ . That is:

Theorem 6 (Universal Composition [Can01]). *Let π, ϕ, ρ be protocols, such that π UC-emulates ϕ . Then the protocol ρ^π UC-emulates ρ^ϕ . In particular, if ρ^ϕ UC-realizes an ideal functionality \mathcal{F} , then so does ρ^π .*

We note that the universal composition theorem hold only in the case where protocols π and ϕ are *modular*. Essentially, a protocol is modular if in no instance s of this protocol there is a subroutine ITI I of some ITI which is part of the instance (or a subroutine thereof), where I receives input from or sends outputs to and ITI that is not a descendant of a member of instance s . Alternatively, modular protocols are also called subroutine respecting [Can01].