

Hecate: Abuse Reporting in Secure Messengers with Sealed Sender

Rawane Issa
Boston University

Nicolas Alhaddad
Boston University

Mayank Varia
Boston University

Abstract

End-to-end encryption provides strong privacy protections to billions of people, but it also complicates efforts to moderate content that can seriously harm people. To address this concern, Tyagi et al. [CRYPTO 2019] introduced the concept of *asymmetric message franking* (AMF) so that people can report abusive content to a moderator, while otherwise retaining end-to-end privacy by default and compatibility with anonymous communication systems like Signal’s sealed sender.

In this work, we provide a new construction for asymmetric message franking called Hecate that is faster, more secure, and introduces additional functionality compared to Tyagi et al. First, our construction uses fewer invocations of standardized crypto primitives and operates in the plain model. Second, on top of AMF’s accountability and deniability requirements, we also add forward and backward secrecy. Third, we combine AMF with source tracing, another approach to content moderation that has previously been considered only in the setting of non-anonymous networks. Source tracing allows for messages to be forwarded, and a report only identifies the original source who created a message. To provide anonymity for senders and forwarders, we introduce a model of *AMF with preprocessing* whereby every client authenticates with the moderator out-of-band to receive a token that they later consume when sending a message anonymously.

1 Introduction

End-to-end encrypted messaging systems like Facebook Messenger, Signal, Telegram, Viber, and WhatsApp are used by billions of people [80] due to their powerful combination of cryptographic protections and ease of use. The security guarantees provided by encrypted messengers are both varied and valuable [77]: confidentiality and integrity from authenticated key exchange [16, 20, 53], deniability from the use of symmetric authenticated encryption [15, 29, 39], and forward and backward (aka post-compromise) security via key evolution [25, 40]. However, these very security guarantees

complicate efforts by secure messaging platforms to investigate reports of abuse or disinformation campaigns, which can have serious consequences for individuals and collective society [10, 31, 72, 74, 79].

To address these concerns, the security research community has developed three methods to augment end-to-end messengers with privacy-respecting technologies to assist with content moderation: message franking, source tracing, and automated identification. First, *message franking* [29, 32, 39, 54, 55, 75] allows recipients to manually report abusive messages with assurance that unreported messages retain all guarantees of secure messengers, and reported messages are both accountable (the moderator correctly identifies the message’s sender) and deniable (the moderator cannot prove this fact to anybody else). Second, *source tracing* [59, 76] allows the moderator to pinpoint the original source of a viral message rather than the person who forwarded the message to the eventual reporter. Finally, *automated identification* [11, 52] proactively matches messages against a moderator-provided list of messages using a private (approximate) set membership test, with possible interventions like rate-limiting or warning labels in case of a match [73]. We refer readers to [65] for more details about content moderation in encrypted settings.

This work contributes a new construction called Hecate that simplifies, strengthens, and unifies the first two content moderation techniques: asymmetric message franking and source tracing. We do not consider automated identification, focusing instead on abuse reporting schemes that empower the people who receive messages to choose the action they wish to take [49, 60]. To provide context for our work, we describe the nascent space of message franking and source tracing in more detail before explaining our improvements.

Prior work. There exists a long line of research into the security of end-to-end encrypted messaging systems (EEMS) at both the protocol design and software implementation layers (e.g., [4, 9, 13, 19, 23, 48]). Our work relies on these analyses in order to treat the underlying messaging protocol in a black-box manner and abstract away its details, so we can focus on

the additions provided by content moderation protocols.

Message franking constructions involve four parties: a sender and receiver of a message, plus the platform providing the secure messaging service and a moderator who acts on abuse reports (see Figure 1). *Symmetric* message franking protocols are limited to the setting in which the platform and moderator are the same entity and have sufficient network-level visibility to pinpoint the sender of each message. At a high level, these constructions operate as follows: when a sender submits a ciphertext corresponding to the message m , the platform signs an attestation binding the sender’s identity to a commitment $\text{com}(m)$ provided by the sender in the clear. The receiver also sees this commitment (e.g., as part of a robust encryption scheme [1, 33, 34]) and can check whether it is correct, dropping the packet if it is malformed. Subsequently, the receiver can report the message as abusive by opening all [29, 32, 39] or part [54] of the commitment; then, the platform can determine whether the message is abusive and take appropriate action.

The work of Tyagi et al. [75], which is the starting point for this paper and which we will henceforth refer to as TGLMR, introduces the notion of *asymmetric message franking* (AMF) that removes the limitations from above. Specifically, AMF can operate even when using Signal’s sealed sender [68] or an anonymous communication system (e.g., [3, 26, 28, 78]) that hides the identity of the sender or receiver from the platform. Furthermore, the system is secure whether the moderator and platform are operated by the same or different entities.

Inspired by designated-verifier signatures [46, 64], the TGLMR construction requires the sender to make a Diffie-Hellman tuple $\langle g, g^{\text{sk}_{\text{src}}}, g^{\text{k}_{\text{mod}}}, g^{\text{sk}_{\text{src}} \cdot \text{k}_{\text{mod}}} \rangle$ involving the moderator’s secret key and her own, as well as a non-interactive zero knowledge proof that the tuple is well-formed. TGLMR achieves accountability and deniability for the sender, but doesn’t provide forward and backward security due to the use of long-lived secret keys. Moreover, it is complex and expensive to implement (see §6), and requires a non-falsifiable knowledge of exponent assumption in the random oracle model. Finally, TGLMR does not easily generalize to more complex conversation graphs that allow for forwarding.

Another line of research investigates the ability for the moderator to trace the original source of messages that might have been forwarded several times within an EEMS. Tyagi, Miers, and Ristenpart [76] began this line of study with their Traceback scheme, which reveals to the moderator the entire path from the original source to the reporter, but it requires server-side storage proportional to the number of messages eligible to be traced. Two recent works provide *source tracing*, identifying only the original source of a reported message. First, Peale, Eskandarian, and Boneh [59] contribute a source tracing construction that inherits most security properties from the underlying EEMS (see Table 1). Using more expensive crypto operations, the stronger variant of their construction is the only one to date to achieve tree unlinkability — namely,

Construction	Features					Security Guarantees							
	Abuse reporting	Message forwarding	Source tracing	Trace info	Threshold report	Confidentiality	Anonymity	Tree unlinkability	Deniability	Forward security	Backward security	Unforgeability	Accountability
<i>Signal</i>	○	●	×	×	×	●	●	●	●	●	●	●	×
Tyagi et al. [75]	●	○	×	src	○	●	●	×	○	○	●	●	●
Traceback [76]	●	●	●	path	○	●	○	○	○	○	○	○	○
FACTS [55]	●	●	●	src	●	●	○	●	●	○	●	●	●
Peale et al. [59]	●	●	●	src	○	○	○	●	●	○	●	●	●
Hecate (<i>this work</i>)	●	●	●	src	○	●	○	○	●	●	●	●	●

●: fully provided, ○: provided but not proven, ◐: partially provided, ○: not provided, ×: not applicable

Table 1: A comparison of features and security properties provided by the Signal EEMS protocol as well as several abuse reporting constructions. Security properties are described in §2.2 and §5. For the anonymity column, ◐ refers to providing anonymity at the level of Signal’s sealed sender [68].

that a receiver who gets the same message twice cannot tell if they originate from the same or different sources. Second, the FACTS scheme by Liu et al. [55] provides source tracing along with a threshold reporting scheme so that the moderator only learns when sufficiently many complaints have been lodged against an abusive source client. However, none of these traceback or source tracing schemes [55, 59, 76] considers backward security as part of their security model. Also, none provides full anonymity of senders and receivers from the EEMS platform or moderator; FACTS is compatible with a network that provides one-sided anonymity, but it requires senders to identify themselves and request tokens from the moderator on the fly whenever they wish to send a message.

This leaves the following open question:

Can we design a protocol that simultaneously provides asymmetric message franking (AMF) and source tracing, achieves forward and backward security, maintains anonymity of senders and receivers to the extent provided by the underlying EEMS network, and only makes black-box use of standardized cryptography in the plain model?

In this work, we answer the question in the affirmative.

Our contributions. In this work, we provide a new definition and construction for asymmetric message franking (AMF) that is more general, more secure, and faster than previous work. To achieve this goal, we revisit the decision by TGLMR [75] to “restrict attention to non-interactive schemes for which franking, verification, and judging requires sending just a single message.” On its face, this restriction seems natural because end-to-end encrypted messengers are designed to work asynchronously in situations with limited network connectivity, so one-round (online) protocols are desirable. However, this restriction also appears to direct the solution

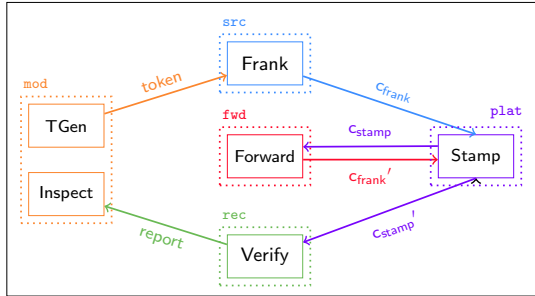


Figure 1: Diagram of Hecate’s data flow for a message m , from the source (top) to the forwarder (middle) and then to a reporting receiver (bottom). The commands match our definition of AMF with preprocessing (Def. 1), and the variables token , c_{frank} , c_{stamp} , and report are defined in Fig. 2.

space toward expensive crypto tools like designated-verifier signatures and zero knowledge proofs.

Our core insight is to introduce an *AMF with preprocessing* model as shown in Fig. 1. As before, the online work of message franking and transmission requires only one round of communication from the source to platform to receiver. Beforehand, we allow the source and moderator to engage in a single data-independent preprocessing interaction to produce *tokens* that can be consumed during the online phase. Preprocessing can be batched to produce many tokens at once, it can be performed during off-peak hours when the source’s device is connected to power and wifi, and it should be performed in advance rather than on the fly in order to avoid network-level traffic linking attacks [57]. As with MPC [7] or PIR [8], we show that adding a preprocessing round to AMF allows for more efficient protocols, and in particular allows us to answer the open question from above.

Concretely, we contribute an AMF scheme called Hecate. Our construction leverages the fact that, with preprocessing, the communication path of reported messages begins and ends with the content moderator. Ergo, we can use techniques from (faster) symmetric message franking whereby the moderator can prepare a token (e.g., a symmetric encryption of the source’s identity) that is only intelligible to its future self. The token is passed through the sender \rightarrow platform \rightarrow receiver communication flow of an EEMS; that said, end-to-end encryption prevents the platform from viewing the token.

Hecate also supports *source tracing*, in which receivers can forward messages along with their corresponding tokens. Any recipient can choose to report an abusive message; this only requires sending one communication to the moderator.

A big challenge in our construction is to combine message forwarding with our AMF *backward* (or post-compromise) security requirement, which states that an attacker who previously (but no longer) controlled a source’s device cannot blame the source for new messages. To our knowledge, this work is the first one to consider and formalize backward se-

curity within AMF. As we will discuss in more detail in §2, the challenge in combining AMF with backward security stems from the fact that indirect recipients of a message (after forwarding) cannot rely on the backward security of the underlying encrypted messaging protocol to tell whether the token is produced by the now-uncorrupted source or the attacker.

In summary, we make four contributions in this work.

- We rigorously define AMF with preprocessing (§3). We generalize the definition from TGLMR, formalize forward and backward security, and add source tracing.
- We provide a construction called Hecate (§4). It requires only a few black-box calls to standard crypto primitives.
- We formalize and prove (§5) that Hecate achieves all of the security guarantees shown in Table 1.
- We implement Hecate (§6) and integrate it into a Signal client. We show that Hecate’s performance compares favorably to prior work and is imperceptible in practice.

Before continuing, we wish to stress that any decision to use content moderation within end-to-end encrypted messengers requires weighing all of its potential benefits and risks, including the limitations of Hecate and prior works (see §7), and the risk of abuse by (or coercion of) the moderator. This is a complex policy question whose discussion should involve computer scientists, but not only computer scientists. We take no stance on the policy question in this work; instead, we observe that these policy discussions are already ongoing [2, 17, 62] and that a sub-optimal understanding of the technological possibilities may push a service provider or nation-state policymakers toward a worse policy decision. We undertake this research in order to demonstrate the feasibility of alternatives to blunt privacy-inhibiting legislation.

2 Overview

In this section, we describe our objectives for an asymmetric message franking (AMF) system. We begin by describing the setting and threat model for our work, and then we provide a high-level description of the security requirements and a brief description of how our Hecate protocol will achieve them.

2.1 Setting and Threat Model

In this work, we consider an EEMS that might contain network-level anonymity protections such as Signal’s sealed sender [68] or Tor [28]. We focus on two party point-to-point communication; that said, our techniques translate directly to Signal’s group messaging protocol as described in §7.

Within the context of any single message transmission, we refer to the participating clients using the following terminology: the *source* who initially produced the message within the messaging platform, the *receivers* who receive the message and can optionally decide to report it (in which case we call them a *reporter*) and the *forwarders* who are receivers that

decide to send the message along to others. All clients only possess the computational power of a phone.

Due to forwarding, the communication graph of each message has the structure of a tree rooted at the source. Several messages can be sent concurrently, and a client can have different roles in the communication trees of different messages.

In addition to the messenger clients, our model contains two (possibly separate, and more computationally powerful) entities that everyone can communicate with: (1) the *platform* that provides the messaging service, and (2) the *content moderator* that receives reports and helps victims of abuse. We consider the platform and moderator as possibly separate so that our model can capture settings where a platform outsources moderation tasks to other, more qualified organizations (e.g., Facebook’s oversight board [58]). We emphasize that our model and construction do not rely upon separation of these roles in any way; they remain valid even if the platform and moderator are operated by the same entity.

In general, the parties in the system view all other parties as potentially malicious and colluding together. Every party wants confidentiality and integrity to the strongest extent possible, even if some or all of their counterparty, the platform, and the moderator are colluding against them. In particular, we wish to retain all of the security goals that end-to-end encrypted messengers provide, as detailed in §2.2 and §5.

The parties’ relationships toward the moderator are subtle and merit further discussion. The moderator and platform view each other as semi-honest; looking ahead to our Hecate construction, the moderator trusts the accuracy of any timestamp applied by the platform but it need not trust the platform for any other purpose. Clients have a choice: if they view the moderator and platform as malicious and colluding, then they must be assured of limits to the moderator’s power; or, if they view the moderator as semi-honest then they must be assured that the moderator can perform its role.

In this work, a malicious attacker has the power to compromise one or more parties, in which case it can observe these parties’ local state (e.g., cryptographic keys) and run arbitrary software for the duration of their control of a victim’s machine. A semi-honest party, by contrast, is assumed to perform all actions honestly, and the only objective against such a party is data minimization. We presume that the software implementing the encrypted messenger faithfully reproduces the intended specification so that the adversary cannot control the behavior of honest parties. Put another way, supply chain attacks and formal verification are out of scope of this work.

The objective of holding senders accountable for reported messages creates a tension with the security goals of end-to-end encrypted messengers. In particular, clients no longer receive confidentiality, deniability, anonymity, or other privacy guarantees against the moderator for reported messages. Moreover, an AMF scheme imposes a limit on forward security, because messages sent in the past now can be revealed to the moderator in the present. Our objective is to ensure

security up to these fundamental limits. We emphasize that even if the moderator is malicious and colluding with some clients, *all of the security guarantees for end-to-end encrypted messaging continue to hold for all unreported messages communicated between non-colluding clients*. Moreover, even for reported messages, security holds against all other parties who are not colluding with the moderator.

Another tension exists between content moderation and network anonymity. For example, *sealed sender* is a feature introduced by the Signal protocol to hide the identity of the sender from the platform. It offers sender confidentiality and minimizes the amount of metadata stored by the platform. But if the sender can deny ever sending a message, then can we hold anyone responsible for sending an abusive message? TGLMR [75] resolved this dilemma using zero-knowledge signatures. In this paper, we contribute a construction based solely on black-box use of standard crypto primitives.

2.2 Security goals

In an asymmetric message franking scheme, we aim to provide all of the security and privacy goals of encrypted messengers [23, 77]. Some EEMS goals (cf. §3.1) are already consistent with content moderation, in which case AMF constructions can use these properties and must ensure that they don’t weaken them. To give a concrete example for our Hecate protocol: we use the EEMS as a black box, and we will take advantage of the receiver’s ability to authenticate the sender’s identity. On the other hand, some security goals are not fully compatible with content moderation, in which case we aim to make the smallest modification possible.

Below, we describe each security goal from Table 1 and highlight the extent to which it is impacted by content moderation. These security goals apply to all clients who construct properly formatted messages that adhere to the encrypted messaging protocol, whether or not their messages are subsequently reported. That is: even though malicious parties in a crypto protocol receive no security guarantees, the mere act of sending a reported message does not render a client malicious.

- *Confidentiality*. Anyone not involved in the creation, forwarding, or reporting of a message m must not learn anything about m except an upper bound on its size.
- *Anonymity*. The AMF scheme should not allow the platform, receiver, or moderator to learn anything about the source and forwarding path of a message beyond what they would learn from the underlying EEMS or a report.
- *Deniability*. Every user should be able to deny a claim about the contents of an unreported message made by any adversary (even other recipients of this message). Also, reported messages are deniable to anyone other than the moderator, and if the moderator’s keys are breached then they become deniable to everyone.

- *Forward security.* Adversaries that compromise users' state in the present should not be able to deduce anything about messages exchanged in the past. This goal does not apply to messages that happen to remain on the phone in the present, which can still be read and reported.
- *Backward security.* Once a client recovers from a compromise event, then the compromised state becomes 'useless' after a short delay. That is, the adversary cannot subsequently originate a new message that (if reported) would cause the moderator to blame the victim client.
- *Unforgeability.* The adversary cannot send a message that appears to be sent by another party. An honest receiver will reject any malformed or tampered messages.
- *Accountability.* If a message passes a receiver's verification check and is subsequently reported, the moderator will trace it back to its original source. That is: nobody can falsely accuse someone who wasn't the source of a message, and the true source cannot evade detection and yet also have the message verified by the receiver.

2.3 Protocol Overview

In this section, we give a high level overview of our Hecate protocol in two stages (with and without message forwarding) and explain informally how it satisfies our security goals.

Hecate without forwarding. At a high level, our Hecate construction can be thought of as an interactive variant of designated-verifier signatures. Given a message m , the source constructs a 2-out-of-2 secret sharing, say $H(m) = x_1 \oplus x_2$. In Hecate, the moderator binds x_1 to the source's identity (which on its own reveals nothing about m), and then the source binds x_1 to x_2 without using any long-lived keys.

As shown in Fig. 2: since one of the two shares can be sampled even before m is known, during preprocessing the moderator selects x_1 as an encryption of the source's identity (which appears random to everyone else), samples an ephemeral digital signature keypair (sk_e, pk_e) , and signs both x_1 and pk_e . The tuple of x_1 , pk_e , and their signature constitute the preprocessing token $token$. During the online phase, the source uses the ephemeral key sk_e to sign x_2 ; we refer to the pair of x_2 and its signature as another token $token$.

The source provides both tokens to the receiver within the payload of an ordinary Signal packet, as shown in Fig. 2; ignore the other elements of the franked ciphertext C_{frank} for now. Any receiver can check on its own whether the signatures are valid and the underlying values x_1 and x_2 combine to form the real message m that the receiver also gets from the underlying Signal communication; if verification fails, then the message is malformed, so it is dropped without displaying on the receiver's device. If a verified message is later reported, the two tokens together will convince the moderator that the source was the originator of message m .

Achieving our security goals. Many of our security guar-

antees follow directly from the corresponding property of the underlying EEMS, so we focus on the most challenging goals here. Hecate provides accountability for the same reason as symmetric messaging franking schemes: the moderator created an authenticated encryption of the source's identity for its future self. Forward security holds because ephemeral signing keys sk_e from the past were deleted before a compromise event in the present. Deniability can be shown in two parts: if the moderator's keys are breached then anyone can produce signatures for any choices of x_1 and x_2 , and otherwise the source's identity is hidden within the encrypted token so anyone could have 'forged' signatures of an (x_1, x_2) pair using her own tokens rather than those of the real source.

Backward (or post-compromise [25]) security is more challenging to address, and it is worth pausing for a moment to discuss what this guarantee means in the context of content moderation. If an adversary corrupts the source's phone, it *can* produce messages whose reports blame the source; this is inevitable. Our goal is to ensure that once the source recovers control of her phone, then (perhaps after a short delay δ) any new message produced by the adversary cannot implicate the honest source. To provide this guarantee within Hecate, the moderator includes a timestamp within its attestation to x_1 , and receivers drop any message where this timestamp is too old. This ensures that an adversary cannot use stale pre-processing tokens long after the compromise event.

Hecate with message forwarding. Next, we allow forwarding of messages and consider source tracing, in which the moderator should identify only the original source of a reported message. For the most part, our construction is already amenable to source tracing: a forwarder can simply include the original source's tokens within a forwarded message rather than generating new tokens that would implicate herself. However, our timestamp-based solution to backward security now fails because the age of x_1 is insufficient to determine whether the original source had control of her cryptographic keys at the moment that the *original message* was sent (as opposed to the time of the forwarding).

We solve this problem by appending a timestamp time as the data traverses through the platform. To verify backward security, it suffices to verify whether the pre-processed token (which contains the identity of the source to blame) was produced close in time to the message transmission. Timestamps for forwarded messages are disregarded; future recipients only care about the timestamp from the original source.

It only remains to bind the source timestamp to the message, so that it cannot be tampered later. Note that we cannot reveal x_2 to the platform, or else the platform and moderator together could recover the content of messages. Blind signatures are a possible solution to allow the platform to timestamp-and-sign obliviously, but constructions that only require one message received and sent by the platform require trusted setup [35], non-standard crypto assumptions [37,38], or a concretely slow

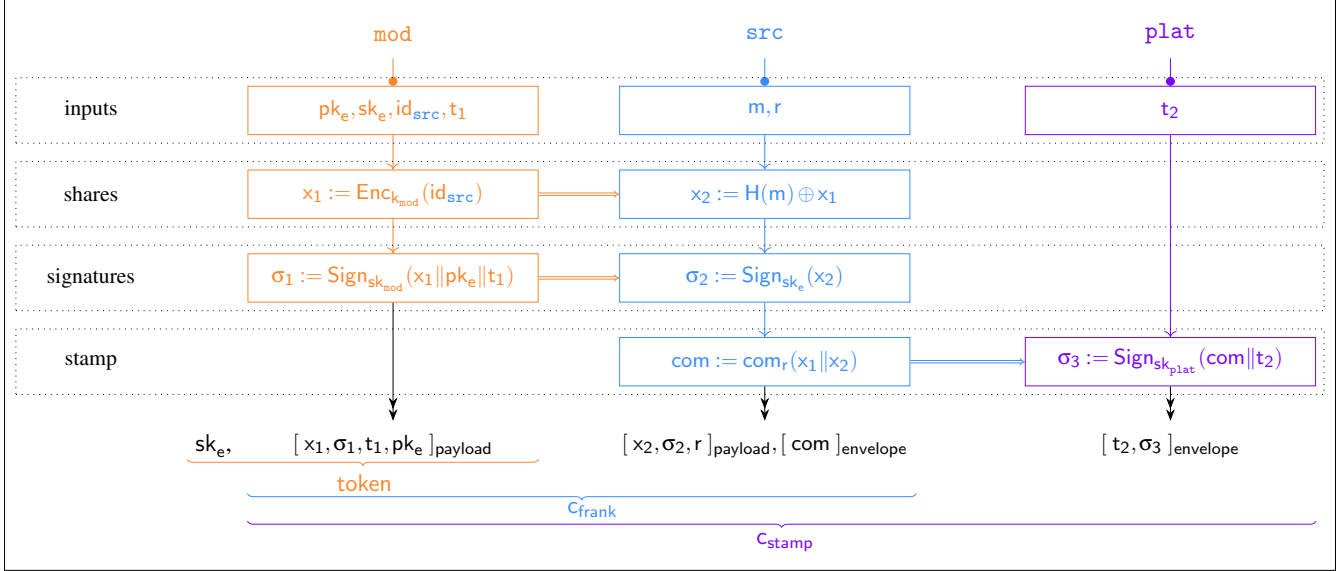


Figure 2: The construction of the different parts of a franked cipher. The outputs of the diagram correspond to each party’s contribution to the eventual stamped cipher C_{stamp} . The source constructs the franked cipher C_{frank} using a token provided by the moderator during preprocessing. We denote by payload and envelope two different parts of the ciphertext as defined in the Signal sealed sender protocol [68]; the platform and receiver can read the envelope whereas only the receiver can read the payload.

runtime with non-black-box reductions [36, 50]. Instead, we take advantage of the fact that the platform’s actions need only be verified by recipients who already know x_1 and x_2 , so it suffices for the platform to produce a signature σ of the current time together with a commitment to the two shares. The corresponding decommitment randomness can be sent to the receiver within the encrypted messenger payload, so that the recipient can verify that it is well-formed.

3 Definitions

In this section, we present a nearly black-box model of the Signal protocol that we use in this work, and then we detail a new definition for an asymmetric message franking scheme that generalizes TGLMR [75]. We refer readers to Appendix A for formal definitions of the cryptographic building blocks we use in this work: commitments, authenticated encryption, digital signatures, and collision-resistant hash functions.

3.1 Modeling an EEMS

End-to-end encrypted messaging systems (EEMS) using the Signal protocol [70] are a complex, delicate combination of standard cryptographic primitives. Starting with Cohn-Gordon et al. [23], there is a long line of research that analyzes the security of EEMS constructions such as the two-party Signal protocol itself (e.g., [4, 9, 48]), modified versions that provide stronger guarantees (e.g., [45, 47, 61]), and extensions to support group messaging (e.g., [18, 22, 63]).

In this work, we wish to treat an EEMS as a black box and consider only an API-level description of its operation and underlying security guarantees. For this reason, we follow the UC models of an EEMS as shown by the recent works of Bienstock et al. [13] and Canetti et al. [19], rather than the game-based definitions in the literature (e.g., [4, 9]). The ideal functionality $\mathcal{F}_{\text{Signal}}$ in these works models the creation, evolution, and destruction of communication ‘sessions’ between different pairs of parties, and keeps track of the long-term and ephemeral state that parties must hold for each operation.

We follow this abstract model, with three changes. First, we allow a sender to attach public information outside of the encrypted message on the envelope that is visible to the platform, as shown in Fig. 2. As defined by Signal’s sealed sender construction [68], the *envelope* constitutes the outer shell of an encrypted message and contains any data that is visible to the platform such as the ciphertext and the recipient’s address. Second, in order to support an anonymous network, we allow for inputs involving the parties’ identities to be optional. Third, we add an explicit forgery method to highlight the fact that an EEMS achieves deniable authentication [30] (which is implicitly true in universally composable security models [21]): that is, the sender and receiver can forge a transcript showing that a message originated with the other party.

Hence, our abstract model of Signal involves three methods. All of these methods implicitly use the state of the party (or parties) that participate in each method.

- $\text{send}_{\text{eems}}(m^*; \text{id}_{\text{src}}, \text{id}_{\text{rec}}) \rightarrow c$: Run by the source client id_{src} with message m^* , this method sends a ciphertext

c to the platform. This message m^* might contain payload and envelope components, similarly to how Signal’s sealed sender operates. We sometimes omit the latter two inputs when they are clear from context.

- $\text{deliver}_{\text{eems}}(c; \text{id}_{\text{rec}}) \rightarrow m^*$: An interactive protocol in which the platform delivers a ciphertext c to the receiver id_{rec} . If this receiver was the intended target of a previous $\text{send}_{\text{eems}}$ that produced c , then they can decrypt using their local state to recover m^* . Here, it is unclear whether to include id_{rec} as an input: for an anonymous communication channel it is important that the platform not know id_{rec} , but for non-anonymous networks it may be required. We leave id_{rec} as an optional parameter, and throughout this work we focus on the stronger setting in which the network is anonymous so this input is omitted.
- $\text{forge}_{\text{eems}}(m^*; \text{id}_{\text{src}}, \text{id}_{\text{rec}}) \rightarrow c$: A forgery algorithm executed by a party id_{rec} and requiring its state $\text{state}_{\text{rec}}$. It forges a transcript that looks as though the message m^* were sent by its counterparty id_{src} in an EEMS communication, with a destination of id_{rec} . The parameters id_{src} and id_{rec} are optional for the same reasons as $\text{send}_{\text{eems}}$, and they will be omitted from this work.

3.2 Defining AMF with Preprocessing

Next, we present a rigorous definition for an asymmetric message franking system with preprocessing. This definition extends the one from TGLMR [75] in two ways. First, it includes an (optional) out-of-band communication between the moderator and sender, which results in a one-time *token* that is consumed when sending a message. Second, it is designed in a modular fashion so that it can be built on top of any EEMS that adheres to the model in §3.1.

Definition 1. An asymmetric message franking scheme with preprocessing AMF = (KGen, TGen, Frank, Forward, Stamp, Inspect, Verify, Forge_{mod}, Forge_{rec}) is a tuple of algorithms called by different parties in the messaging ecosystem. We assume that each party has a unique identifier id provided by the underlying EEMS, and we define a state variable state for each party containing all keys and tokens generated by the AMF scheme and the underlying EEMS that have not yet been deleted. The algorithms operate as follows.

- $(\text{pk}, \text{sk}) \leftarrow \text{\$KGen}()$: The key generation algorithm accessed by any party in the EEMS and used for creating (potentially multiple) cryptographic keys. The algorithm is at least run once at the beginning of time to setup the long-term key material for each party.
- $\text{TGen}(\text{id}_{\text{src}}, \text{time}_{\text{mod}}, \text{k}_{\text{mod}}) \rightarrow \text{token}$: An algorithm run by the moderator periodically that provides a one-time token for use when sending a message. An honest moderator should only provide tokens to a participant that correspond to their actual identity id_{src} . It is assumed

that the moderator can rely on the EEMS to authenticate a user’s identity id_{src} before running TGen.

- $\text{Frank}(\text{state}_{\text{src}}, m, \text{id}_{\text{rec}}, \text{token}) \rightarrow m_{\text{frank}}$: The message franking algorithm that allows a user with state $\text{state}_{\text{src}}$ to frank a plaintext message m that they wish to send to a receiver id_{rec} , using the token received during preprocessing. The state $\text{state}_{\text{src}}$ contains all key material produced by KGen and the underlying EEMS, although Frank need not use this state. The resulting franked message m_{frank} can be sent to the platform using $\text{send}_{\text{eems}}$.
- $\text{Stamp}(c_{\text{frank}}, \text{sk}_{\text{plat}}, \text{time}) \rightarrow c_{\text{stamp}}$: The stamping procedure run by the platform to authenticate and timestamp a franked cipher c_{frank} . The resulting stamped cipher c_{stamp} can then be delivered to its intended recipient using the $\text{deliver}_{\text{eems}}$ method. Stamp does not have the sender or receiver’s identity, even if $\text{deliver}_{\text{eems}}$ does.
- $\text{Forward}(m_{\text{frank}}, \text{state}_{\text{fwd}}, \text{id}_{\text{rec}}) \rightarrow m_{\text{frank}}'$: Forwarding algorithm that allows a user with franked message m_{frank} to produce a new franked message m_{frank}' intended for a new recipient id_{rec} . The format of m_{frank} and m_{frank}' are identical, so the ciphertexts of new and forwarded messages look indistinguishable to the platform.
- $\text{Verify}(m_{\text{frank}}, \text{state}_{\text{rec}}) \rightarrow (m, \text{report})$ or \perp : Allows a receiver to validate a franked message m_{frank} with respect to its state $\text{state}_{\text{rec}}$. If valid, Verify returns the corresponding plaintext message m along with a string report that the receiver can send to the moderator if they choose to report an abusive message.
- $\text{Inspect}(\text{report}, \text{k}_{\text{mod}}) \rightarrow (\text{id}_{\text{src}}, m, \text{time})$ or \perp : The inspection algorithm that allows a moderator to handle reported message report using their secret key k_{mod} by validating and possibly source tracing them. If the verification step succeeds, the moderator produces the id of the source id_{src} , the message contents m , and a timestamp of the message time.
- $\text{Forge}_{\text{mod}}(\text{id}_{\text{src}}, \text{id}_{\text{rec}}, m, \text{k}_{\text{mod}}) \rightarrow m_{\text{frank}}$: For deniability, this forgery protocol allows a moderator with secret key k_{mod} to forge a franked message with plaintext m on behalf of a user with id id_{src} and with an intended recipient with id id_{rec} .
- $\text{Forge}_{\text{rec}}(\text{id}_{\text{rec}}, m, \text{state}_{\text{rec}}; \text{id}_{\text{src}}) \rightarrow c_{\text{frank}}$: For deniability, this forgery algorithm allows a receiver with id id_{rec} and state $\text{state}_{\text{rec}}$ to forge a franked ciphertext as though the message m was transmitted through the EEMS by the sender id_{src} to the receiver id_{rec} . Note that id_{src} is an optional parameter and may not be needed by systems that support anonymous messaging.

We say that an AMF scheme with preprocessing is secure if all computationally bounded attackers have negligible advantage in winning the deniability, anonymity, confidentiality, accountability, and backward secrecy games. These games are nuanced to describe; rather than doing so here, we defer discussion to the security analysis in §5 and the full version [44].

Command	Actor	KeyGen	Sign	Verify
TGen	mod	1	1	0
Frank	src	0	1	0
Stamp	plat	0	1	0
Forward	fwd	0	0	0
Verify	rec	0	0	3
Inspect	mod	0	0	3

Table 2: The number of public-key digital signature operations required for each of the interactive algorithms within Hecate (except for the one-time KGen at setup). We only count the additional cryptographic operations required for Hecate beyond those already required by the EEMS.

4 Constructing Hecate

In this section, we describe the Hecate construction in detail. As per Def. 1, Hecate has eight algorithms. We describe them within this section, and we provide the full protocol specification of Hecate in Figs. 3-4. Because they are the most expensive of our standard crypto primitives, we also count the number of public key operations in each step here and in Table 2. For context, the prior AMF scheme from TGLMR [75] required at least 11 modular exponentiations per algorithm.

Key generation. KGen initializes a few long-term keys: the moderator samples an authenticated encryption key k_{mod} and both the moderator and platform sample a digital signature key pair $(pk_{\text{mod}}, sk_{\text{mod}})$ and $(pk_{\text{plat}}, sk_{\text{plat}})$. One strength of Hecate is that individual parties do not need any long-term key material besides their existing EEMS keys, which simplifies our analysis of forward and backward security.

Token generation during preprocessing. In TGen, the moderator creates a batch of tokens for users at specific time intervals. Each token provides users with:

- Ephemeral session keys (pk_e, sk_e) that they can use to sign their message. None of the keys tie to users’ long-term key material, thus giving the sender plausible deniability and confidentiality with respect to other users.
- A dual purpose randomized encryption $x_1 := \text{Enc}_{k_{\text{mod}}}(\text{id}_{\text{src}})$ of the user’s identity id_{src} under the moderator’s secret key k_{mod} that enforces accountability with respect to the moderator, confidentiality with respect to other users, and provides token integrity. The latter property is ensured by having the sender create a share x_2 that along with x_1 reconstructs to a hash of the sent message.
- A timestamp t_1 that provides backward security.
- A signature σ_1 by the moderator of the entire token that guarantees integrity and unforgeability of the token.

As shown in Table 2, the moderator requires two public key operations for token generation: 1 keygen operation to produce the ephemeral key pair and 1 signature to sign the public ephemeral key with the identity of the sender.

Message franking. The Frank method is executed every time the source wishes to send a message. The $\text{construct}_{\text{frank}}$ procedure requires an input plaintext message m from the source and consumes a single token at a time, and it produces a franked message m_{frank} . This can be combined with $\text{send}_{\text{eems}}$ to relay a franked ciphertext c_{frank} to the platform.

To produce the franked message, the sender begins by unpacking x_1 from the token and computes x_2 such that these variables constitute a 2-out-of-2 sharing of $H(m)$. Next, x_2 is signed via the ephemeral keys in the original token to produce σ_2 . Collectively, x_2 , σ_2 , and elements of the pre-processing token (excluding the secret ephemeral key) will constitute the payload of the franked message. Then, the sender creates a commitment com of $x_1 || x_2$ using the randomness r . The user then pushes com onto the envelope of the franked message and appends r to its payload. In total, a sender only requires 1 public key operation to sign the second share x_2 . The constructed franked message m_{frank} has several properties: x_2 and com bind the online and preprocessing stages together, the signature allows the receiver to check the well-formedness of the message, and the use of an ephemeral signing key provides deniability with respect to anyone other than the moderator.

Stamping. In Stamp, the platform timestamps and digitally signs the envelope of a franked cipher c_{frank} ; ergo, this procedure requires 1 public key operation. Then, the platform relays the resulting stamped cipher c_{stamp} to its intended recipient. Stamping prevents a preprocessing token from being used indefinitely after a compromise to blame a victim client for unsent messages.

Verification and reporting. Upon reception, the receiver executes Verify to validate the signatures, timestamps, packet integrity, and envelope commitments. If a message fails the verification check, then the receiver drops the packet and the application never displays the plaintext message. Otherwise, Verify generates a plaintext message m that can be displayed on the receiver’s phone, and a report that can be sent out to the moderator if the receiver so chooses.

In Hecate, the report solely consists of the franked message m_{frank} . When the moderator receives a report, they locally run the Inspect method which performs the same verification procedure as the recipient, and if successful, decrypts the source’s identity from the ciphertext x_1 within the token. Both the receiver of a message and a moderator who receives a report require 3 signature verifications to check that the two shares are not tampered with and have the right timestamps.

Message forwarding. Verified messages can alternatively be forwarded using the optional Forward method. There are two differences between Forward and Frank: the forwarder creates a dummy commitment outside the Signal envelope, and it moves the true commitment and signed timestamp into the payload of the franked message. Because it reuses the prior signature, the forwarder doesn’t perform any public key operations of its own. Note that Frank and Forward payloads are distinguishable by the receiver but indistinguishable to the

Frank [src → plat]

```

1: mfrank := constructfrank(token, m)
2: cfrank := sendeems(mfrank)
3: return cfrank

```

Forward [fwd → rec]

```

1: mfrank' := constructfwd(mfrank)
2: cfrank' := sendeems(mfrank')
3: return cfrank'

```

Stamp [plat → rec]

```

1: cstamp := stamptime(cfrank)
2: cstamp' := sendeems(cstamp)
3: return cstamp'

```

KGen [mod and plat, separately]

```

1: kmod ← $EncKGen(1n)
2: (pkmod, skmod) ← $SigKGen(1n)
3: (pkplat, skplat) ← $SigKGen(1n)

```

TGen [mod → src]

```

1: token := constructtoken(skmod, idsrc)
2: return token

```

Verify [rec → mod/rec]

```

1: mfrank := delivereems(cstamp)
2: mfrank := movestamp(mfrank)
3: (m, report) := vfRec(mfrank)
4: if (m, report) ? = ⊥ :
5:   return ⊥
6: return (m, report)

```

Inspect [mod]

```

1: if vfMsg(report) :
2:   return (Dec(report.x1), report.t2)
3: return ⊥

```

Figure 3: Hecate’s construction along with the transmissions using the encrypted messenger. The notation $[a \rightarrow b]$ means that party a executes the method and sends the returned value to b . The receiver of a message may elect to forward or report it; we assume that Forward is preceded by a successful invocation of Verify. See Fig. 4 for more details on the methods used here.

construct_{token}(k_{mod}, sk_{mod}, id_{src})

```

1: (pke, ske) ← $SigKGen(1n)
2: t1 := time()
3: x1 := Enckmod(idsrc)
4: σ1 := Signskmod(x1 || pke || t1)
5: token := (x1, t1, σ1, (pke, ske))
6: return token

```

construct_{frank}(m, token)

```

1: r ← $ {0, 1}n
2: (x1, t1, σ1, (pke, ske)) = token
3: x2 := split(x1, H(m))
4: σ2 := Signske(x2)
5: com := comr(x1 || x2)
6: envelope := com
7: payload := (x1, x2, pke, r, t1, σ1, σ2)
8: fwdpayload := ∅ // reserved for forwarder
9: mfrank := (m, payload, fwdpayload, envelope)
10: return mfrank

```

construct_{fwd}(m_{frank})

```

1: mfrank := movestamp(mfrank)
2: mfrank.envelope ← $ {0, 1}n
3: return mfrank

```

stamp_{time}(c_{frank}, sk_{plat})

```

1: t2 = time()
2: σ3 := Signskplat(com || t2)
3: cstamp.envelope := (com || σ3 || t2)
4: cstamp.payload := cfrank.payload
5: return cstamp

```

move_{stamp}(m_{frank})

```

1: if mfrank.fwdpayload ? = ∅ :
2:   mfrank.fwdpayload := mfrank.envelope
3: return mfrank

```

vfRec(m_{frank})

```

1: if vfMsg(mfrank) :
2:   report := mfrank
3:   return (mfrank.m, report)
4: return ⊥

```

vfMsg(report)

```

1: b1 := vfToken(report)
2: b2 := vfCom(report)
3: b3 := vfExp(report)
4: return b1 ∧ b2 ∧ b3

```

vfToken(report)

```

1: reveal := open(x1, x2)
2: b1 := (reveal ? = H(m))
3: b2 := Vfpkmod(x1 || pke || t1, σ1)
4: b3 := Vfpke(x2, σ2)
5: return b1 ∧ b2 ∧ b3

```

vfExp(report)

```

1: b := |t1 - t2| ? < expiry
2: return b

```

vfCom(report)

```

1: (com, σ3, t2) ← report.fwdpayload
2: b1 := Vf(x1 || x2, com, r)
3: b2 := Vfpkplat(com || t2, σ3)
4: return b1 ∧ b2

```

Figure 4: Hecate’s subroutines. See Appendix A for specifications and security guarantees of the crypto primitives used. We omit writing out attribute access notation when it is obvious from the context (i.e. com for instance is a shorthand for c_{frank}.com).

Source’s Payload								Forwarder’s Payload		Envelope		
x ₁	x ₂	nonce	pk _e	r	t ₁	σ ₁	σ ₂	envelope of source		com	σ ₃	t ₂
32B	32B	12B	32B	32B	8B	64B	64B	104B		32B	64B	8B

Table 3: The format of a franked message delivered to the receiver, along with sizes in bytes for the implementation in §6. A franked message sent by the source is similar, except the envelope does not yet contain σ₃ or t₂.

platform (meaning that it will stamp a forwarded cipher); see Table 3 for the format of m_{frank} . The receiver of a forwarded message executes Verify using the commitment, signature, and timestamp inside the payload (ignoring the envelope).

We defer discussion of Hecate’s forgery algorithms to §5 with more details provided in the full version of this work [44], since these are proof artifacts of the deniability property rather than actual elements of the construction.

5 Security Analysis

In this section, we formally define the security properties of asymmetric message franking (AMF) with preprocessing, and we prove that Hecate guarantees them. All of our definitions are written as indistinguishability games $\text{GAME}_b^{\mathcal{A}}$ for $b \in \{0, 1\}$, and we want to show that the adversary’s advantage

$$\text{Adv}_{\text{Hecate}}^{\text{game}}(\mathcal{A}) = \left| \Pr \left[\text{GAME}_1^{\mathcal{A}} = 1 \right] - \Pr \left[\text{GAME}_0^{\mathcal{A}} = 1 \right] \right|$$

is negligible for each game, if the adversary \mathcal{A} is computationally bounded to probabilistically polynomial time (PPT).

Due to space constraints, we include here an informal description of each game and concrete theorem statements. We defer full game descriptions and security proofs to [44].

Deniability. Deniability states that *a sender should always be able to deny that they sent a particular message to anyone, except to the moderator when a message is reported*. Deniability could hold with respect to a colluding moderator and receivers, or it can hold against malicious receivers who are not colluding with an honest moderator.

Theorem 5.1. *Hecate is deniable against a moderator. Any adversary \mathcal{A} has advantage $\text{Adv}_{\text{Hecate}}^{\text{denm}}(\mathcal{A}) = 0$.*

The essence of Thm. 5.1 is the claim that Hecate’s real send routine is indistinguishable from the moderator’s forgery. Intuitively, Hecate achieves moderator deniability because Hecate implements algorithms TGen, Frank and Forward without ever using the user’s long term key materials and instead relies on ephemeral keys that the moderator generated. Additionally, the preprocessing token relies on an encryption and signature by the moderator in a way that is not directly bound to the message. This claim holds even against a distinguisher who also has the moderator’s secret key – that is, if the moderator chooses to leak their own keys in an attempt to convince the rest of the world about the actions of a sender. For a rigorous proof of this theorem (and all theorems in this section), see the full version of this work [44].

Theorem 5.2. *Hecate is deniable against a malicious receiver. Concretely, for any PPT adversary \mathcal{A} , there exist PPT adversaries \mathcal{A}' and \mathcal{A}'' such that:*

$$\text{Adv}_{\text{Hecate}}^{\text{denr}}(\mathcal{A}) \leq \text{Adv}_{\mathcal{E}}^{\text{eemsdeniability}}(\mathcal{A}') + \text{Adv}_{\mathcal{E}}^{\text{enccpa}}(\mathcal{A}'').$$

That is: Hecate’s deniability reduces to the deniability and CPA security properties of the underlying EEMS scheme \mathcal{E} .

Anonymity. In this work, anonymity considers the amount of metadata that each entity can learn through the abuse reporting system alone, if the network provides perfect anonymity.

First, anonymity with respect to the receiver guarantees that *receivers should not be able to learn any other member of the forwarding path of a message beyond their direct neighbors*.

Theorem 5.3. *Hecate is anonymous with respect to the receiver. For any PPT adversary \mathcal{A} , there exists an adversary \mathcal{A}' that can win the chosen plaintext attack game with advantage $\text{Adv}_{\text{Hecate}}^{\text{anonr}}(\mathcal{A}) \leq \text{Adv}_{\text{Hecate}}^{\text{enccpa}}(\mathcal{A}')$.*

Informally, this theorem holds because the preprocessing tokens in Hecate only contain any information about the original sender’s identity in encrypted form; without access to the moderator’s secret key, a receiver can’t distinguish between tokens that originate from different senders. Additionally, Hecate stores no information about forwarders of a message at all, thereby guaranteeing their anonymity as well.

Second, anonymity with respect to the moderator ensures that *the moderator should not be able to learn members of the forwarding path of a reported message beyond the neighbors of colluding receivers and the reported source*. Here, honest forwarders want to be assured that, when their direct contacts are honest, only their neighboring recipients know that they forwarded a specific message.

Theorem 5.4. *Hecate is anonymous with respect to the moderator. Any adversary \mathcal{A} has advantage $\text{Adv}_{\text{Hecate}}^{\text{anonm}}(\mathcal{A}) = 0$.*

Confidentiality and Forward Secrecy. Message confidentiality dictates that *any party not involved in the creation, reception or reporting of a message should not be able to learn anything about the message*. Moreover, forward secrecy guarantees that *corrupted users should be guaranteed confidentiality of all their messages and interactions prior to the time of compromise*. In this work, we consider the state of users to consist entirely of their key material and their tokens; ergo, Hecate can only guarantee confidentiality for messages that have been securely deleted from the local device prior to the compromise event.

Theorem 5.5. *Our scheme Hecate is message confidential and forward secure. Concretely, for any PPT adversary \mathcal{A} , there exist PPT adversaries \mathcal{A}' and \mathcal{A}'' such that:*

$$\text{Adv}_{\text{Hecate}}^{\text{conf-fs}}(\mathcal{A}) \leq \text{Adv}_{\text{Hecate}}^{\text{hidingcom}}(\mathcal{A}') + \text{Adv}_{\text{Hecate}}^{\text{enccca}}(\mathcal{A}'').$$

Informally, the theorem holds because Hecate constructs franked messages by appending tokens to the payload of the message, and by adding a commitment and timestamp to its envelope (see $\text{construct}_{\text{frank}}$ and $\text{stamp}_{\text{time}}$ in Figure 4). The

tokens are encrypted alongside the plaintext message. The identifying content of the commitment is encrypted, and we rely on the hiding properties of the commitment scheme and the security of the connection that exists between parties and the platform. No part of a Hecate franked message can therefore break this security property.

Unforgeability and Accountability. These properties describe a scheme’s ability to *bind senders to well-formed messages while guaranteeing that no user can be accused of sending a message that they did not send*. They go hand-in-hand because well-formed messages are necessarily bound to their original sender and cannot be attributed to anyone else.

Theorem 5.6. *Hecate holds users accountable. For any PPT adversary \mathcal{A} that makes at most q queries to its $\mathcal{O}^{\text{send}}$ oracle, there exist PPT adversaries \mathcal{A}' and \mathcal{A}'' such that:*

$$\text{Adv}_{\text{Hecate}}^{\text{acc}}(\mathcal{A}) \leq (q + 1) \cdot \text{Adv}_{\mathcal{S}}^{\text{sig}_{\text{geu-cma}}}(\mathcal{A}') + \text{Adv}_{\mathcal{H}}^{\text{hash}_{\text{coll}}}(\mathcal{A}'').$$

The description of the $\mathcal{O}^{\text{send}}$ oracle and the proof of this theorem are deferred to the full version of this work [44].

Backward Security. Backward Security requires that *an adversary who controlled the state and keys of a device pre-compromise should not be able to benefit from them after device recovery*. In particular, the adversary should be unable to craft new messages from a recovered user or claim that during-compromise messages were sent out (not forwarded) after the compromise period.

Theorem 5.7. *Hecate is backward secure. For any PPT adversary \mathcal{A} , there exist PPT \mathcal{A}' , \mathcal{A}'' , and \mathcal{A}''' such that:*

$$\begin{aligned} \text{Adv}_{\text{Hecate}}^{\text{bac}}(\mathcal{A}) &\leq \text{Adv}_{\text{Hecate}}^{\text{acc}}(\mathcal{A}') \\ &+ \text{Adv}^{\text{sig}_{\text{geu-cma}}}(\mathcal{A}'') + \text{Adv}^{\text{binding}_{\text{com}}}(\mathcal{A}'''). \end{aligned}$$

Informally, this theorem holds because Hecate requires the moderator and platform timestamps to be close to one another. Hecate also binds each of these timestamps to the franked message and the source’s identity, but in a confidential way, to ensure that the adversary cannot evade backward secrecy by using unexpired timestamps from other messages.

6 Implementation and Evaluation

We implemented Hecate as a Rust library that can be used as a back-end by other systems. Our implementation uses Signal’s official platform agnostic API library libsignal-client [69] for our encryption, commitment and hashing building blocks. To that effect, we use libsignal-client’s implementation of AES-256 GCM for symmetric encryption and HMAC with SHA-256 for commitments. We use the ed25519-dalek [27] crate for ed25519 signatures and their associated functions and SHA256 from the sha2 crate for our hash functions. Our implementation is open source and available at [41].

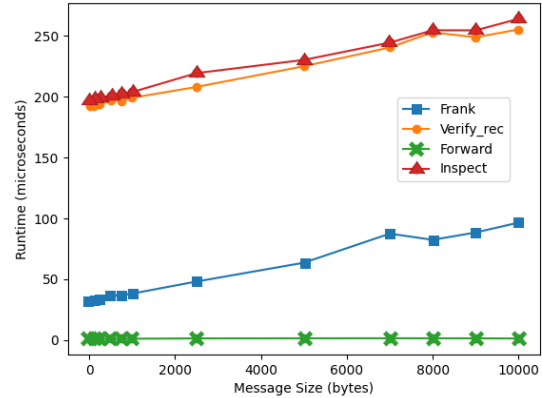


Figure 5: Online runtime of Hecate’s components as a function of message size in bytes

In this section, we show experimental results when executing each component of Hecate in isolation. Then, we measure the overhead of Hecate when integrated into a Signal client.

6.1 Performance Cost and Comparison

In this section, we measure the runtimes and transmission sizes for each component within Hecate using Criterion, a Rust benchmarking suite. We also evaluate prior open-source message franking systems on the same machine and compare them to Hecate.

Experimental setup. We ran all experiments on Amazon Web Services in the US East-Ohio Region, using a t3.small EC2 virtual machine running Ubuntu 20.04 LTS with 2GB of RAM on a 3.1 GHz Intel Xeon Platinum Processor. Each data point shown is the average of 300 trials, with outliers removed. We chose this machine and number of experiment trials to align with prior work [59, 75].

Hecate communication costs. We list the size of Hecate’s franked ciphers in Table 3. The sizes in Table 3 stem from the fact that our libraries yield 32 byte commitments with 32 byte long commitment randomness, 32 byte ciphertexts, 64 byte signatures, 32 byte symmetric and public keys, 12 byte nonces for symmetric encryption, and 8 byte Unix timestamps. We remark that there exist more compact instantiations of these primitives; our choices were motivated by ease of implementation on top of libsignal-client.

Hecate’s online runtime. Fig. 5 shows the performance of each component of Hecate for message sizes ranging from 10 bytes to 10 kB. Overall, the runtime costs remain low especially when compared to the cryptography already required within an end-to-end messaging system (cf. §6.2).

Most components require executing a SHA-256 hash func-

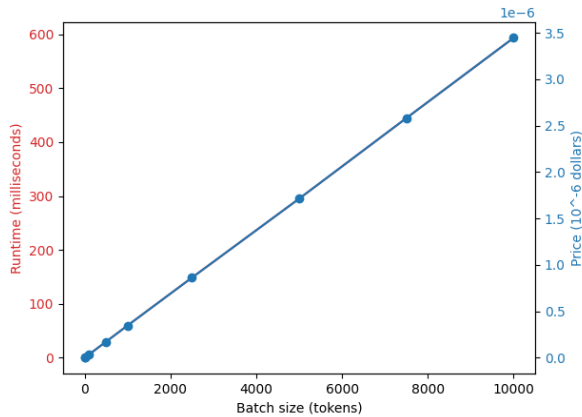


Figure 6: Runtime and dollar pricing of pre-processing token generation (TGen) as function of the token batch size.

	Sent	Received	Report
Tyagi et al. [75]	489 B	489 B	489 B
Peale et al. [59]	256 B	320 B	160 B
Hecate	380 B	484 B	380 B

Table 4: Communication overhead of Hecate and other message franking schemes, in bytes.

tion (to calculate x_2) whose runtime is linear in the message size, along with 0-3 digital signature operations whose cost is independent of message size (cf. Table 2). As a result, the signature(s) dominate the cost for small message sizes and the hash function dominates the cost for large messages. The two costs are balanced at a message size of 7.5 kB, where hashing and digital signing each take about $33\mu\text{s}$. Verify and Inspect are slower because they require about $192\mu\text{s}$ to verify 3 signatures. On the other hand, Forward, Stamp and TGen all have fast runtimes that are independent of message size. We remark that a forwarder is assumed already to have verified a message at reception time, so its only work during Forward is to move the envelope contents into the payload.

Hecate’s preprocessing cost. We also measure TGen over various batch sizes from 1 to 10,000 tokens. Fig. 6 shows the runtime and computational cost based on a rate of 2.09¢ per hour for a t3.small AWS instance at the time of this writing. Our measurements show that the price of generating a batch of 10^4 tokens is 3.45×10^{-6} USD. Extrapolating to the scale of 10^{11} tokens (the approximate number of messages sent through WhatsApp daily [56]), we estimate the cost of token generation to be 35 cents per day. We also highlight that the moderator does not need to remember these ephemeral signing keys in between preprocessing and reporting; in fact the moderator doesn’t require any storage cost at all.

Comparison with prior work. In this section, we compare

our Rust implementation with the open-source software by TGLMR [75] and Peale et al. [59]. To ensure a level comparison: we re-ran the benchmarks from prior work [59,75] on our t3.small AWS instance, we only considered the tree-linkable version of Peale et al., and we removed the double ratchet encryption within the benchmarks of Peale et al. in order to measure only the overhead of their message franking scheme.

We show a comparison of communication overhead in Table 4, and we compare computation overhead in Table 5 for a message size of 1 kB. The benchmarks of TGLMR [75] were orders of magnitude slower than the other works because their construction of designated verifier signature performs more group operations; their communication overhead was also the highest. The comparison between Hecate and Peale et al. [59] is more nuanced. We stress that Hecate achieves additional security properties like anonymity and backward security. As a consequence, senders perform more work in Hecate and transmit more data, whereas Peale et al. leverage a non-anonymous network so that the platform can tag the originator of a message. On the other hand, Peale et al. require a forwarder to generate a commitment, whereas Hecate only requires generating a random 32 byte string (which could even be sampled beforehand).

6.2 End-to-End Prototype Deployment

In this section, we integrate Hecate into an existing Signal client and show that Hecate adds minimal overhead.

Implementation. To test the end-to-end overhead of sending and receiving messages, we integrated our Rust Hecate library into the Java tools signal-cli [67] and libsignal-client [69]. The sender’s Frank procedure adds the Hecate payload to a message before encrypting it using the EEMS, and then appends the Hecate envelope. The receiver decrypts the franked message and runs Verify. Our modified libraries are available as open source repositories [42,43]. The sender gets tokens by running the Rust library prior to the start of the experiment.

We deployed the sender and receiver signal-cli instances on one machine with a 1.90GHz Intel i7-8650U CPU and 16GB of RAM running Ubuntu 20.04 LTS. They were connected over a wide-area network to an instance of signal-server [71].

Evaluation. We measured the client side overhead of running signal-cli with and without Hecate on messages of size 1 kB. In both cases, we measured the average of 10,000 trials of running local signal-cli operations and 600 trials of end-to-end (E2E) latency, with outliers removed. Our timer for the end-to-end latency test starts as soon as the source’s signal-cli begins franking a message, and it ends when the receiver’s signal-cli completes processing the franked ciphertext and outputs the message. These sample sizes are larger than in §6.1 to overcome the noise added by the network latency, the polling rate of the receiver, and the warm-up time of the jvm instance of signal-cli for each of the parties.

	TGen	Frank	Verify	Inspect	Forward	Stamp
Tyagi et al. [75]	–	6339 μ s	5461 μ s	5939 μ s	–	–
Peale et al. [59]	–	8.98 μ s	69.56-138.19 μ s	73.64 μ s	8.46 μ s	24.53 μ s
Hecate	58.4 μ s	38.24 μ s	199.15 μ s	203.87 μ s	1.16 μ s	29.17 μ s

Table 5: Runtimes of message franking schemes, in microseconds, for a message size of 1 kB. The runtime of Verify within Peale et al. [59] differs based on whether the message is authored (left) or forwarded (right).

	Hecate	No Hecate	Diff
Send	2.55 ms	2.38 ms	0.16 ms
Receive	3.19 ms	2.52 ms	0.67 ms
Total	5.74 ms	4.91 ms	0.83 ms
E2E Latency	37.28 ms	36.3 ms	0.98 ms

Table 6: Computation and communication costs for signal-cli with and without Hecate, for a message size of 1 kB. Send and Receive (10^4 trials) correspond to local computation prior to sending or after receiving the message over the network. E2E Latency (600 trials) starts at the beginning of Send and stops at the end of Receive, with network latency included.

Our results are shown in Table 6. They showcase how the findings from Table 5 translate to imperceptible overheads in an actual deployment of Hecate on a Signal client. The inclusion of Hecate adds less than a millisecond of runtime locally and over the network, on average. Moreover, this difference is dwarfed by the sample variance of signal-cli due to the sources of measurement uncertainty.

7 Conclusion and Discussion

In this work, we constructed the first abuse reporting protocol that combines asymmetric message franking and source tracing. We integrated this construction into a Signal client and showed that its performance impact was imperceptible. Along the way, we generalized the AMF model to accommodate pre-processing, and we formalized security properties that hadn’t previously been considered by message franking schemes.

In this final section, we discuss some extensions of Hecate, known limitations, and opportunities for future work.

Extensions. We extend Hecate’s communication from the two-device setting to more realistic flows supported by Signal.

Group Messaging. Hecate’s definitions and constructions can be ported in a straightforward manner to Signal’s group messaging protocol, in which broadcasts to a group of size N are implemented via N individual point-to-point messages, after a server-assisted consensus protocol to determine the group [22]. We note that there exist recent works and an IETF standardization effort on sub-linear ends-to-ends encrypted group chats [5, 6, 12, 24, 66]; it remains an open problem to design abuse reporting mechanisms for these protocols.

Multiple Devices. Hecate can easily support multiple de-

vices for the same user (e.g., a phone and laptop) by giving each device its own independent set of tokens. The moderator can use the same id_{src} for both sets of tokens, so that reports only name an identity rather than a device.

Limitations. We discuss a few limits of our approach.

Reporting Benign Messages. Our construction allows receivers to report messages that may later be deemed to be non-abusive. While it might be possible to require the receiver to prove to an honest moderator that the message they are reporting is actually abusive, this question is incredibly delicate and is therefore out of scope for this and all prior works on end-to-end abuse reporting.

Distinguishing Forwarded vs. Original Messages. In our construction, receivers can distinguish between sent and forwarded messages. While this may be a desirable feature in a messaging app, it is still a leakage in our system.

Forwarding Cycle Linkability. If the forwarding path of a message contains a cycle, i.e. a receiver receives the same forwarded message multiple times, then they can tell that these messages originate from the same source. This is an inherent weakness of Hecate as a result of forwarding the same tokens per message that we do not attempt to protect against. It may be possible to combine Hecate with Peale et al.’s techniques to remove this leakage [59].

Future Work. Looking ahead, we identify several avenues of future research into privacy-respecting content moderation.

Content Censorship. Content moderation systems can be misused for censorship purposes. Questions surrounding what constitutes misinformation or a “bad” message fall outside the scope of this work and into the realm of policy making and social media regulation. We believe however that it may be interesting to federate the role of the moderator in: (1) defining bad messages, (2) verifying reports, (3) taking actions with respect to flagged contents and users.

Super Spreaders. A recent line of work [72, 81] on misinformation spread in social media distinguishes between honest users who forward misinformation and malicious actors that act as super spreaders of misinformation. Honest users can mistakenly forward or send misinformation content without ever realizing it. Super spreaders on the other are adversarially creating or spreading bad content. Future work could examine aggregate behavior in order to distinguish malicious vs. mistaken users.

Partial opening. Known AMF constructions like Hecate

only allow the receiver to report all or none of a message. It should be possible to achieve partial opening to the moderator by extending the message franking techniques of Leontiadis and Vaudenay [54] to the asymmetric setting.

Stronger Notions of Backward Security. Backward security makes no guarantees with respect to anything created during the time of compromise. In the context of content moderation, this implies that the adversary can blame users for old compromised messages. We encourage future research into ways to limit the damage of adversarial moderation reports or to allow honest parties to correct the record post-recovery.

Ensuring System Security. Finally, we emphasize that our study of abuse reporting has been primarily through a cryptographic lens, and as a result does not capture all aspects of security. For example, many of our crypto definitions assume that clients already have sufficient preprocessing tokens in hand. When implementing Hecate, careful attention is required to ensure that adversaries cannot obtain a side channel by, e.g., influencing when preprocessing is run.

Acknowledgments

The authors thank Ran Canetti, Hoda Maleki, Leo Reyzin, Sarah Scheffler, and the anonymous reviewers for their insightful comments and valuable feedback. This material is based upon work supported by the National Science Foundation under Grants No. 1718135, 1739000, 1801564, 1915763, and 1931714, by the DARPA SIEVE program under Agreement No. HR00112020021, and by DARPA and the Naval Information Warfare Center (NIWC) under Contract No. N66001-15-C-4071. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of NSF, DARPA, or NIWC.

References

- [1] Michel Abdalla, Mihir Bellare, and Gregory Neven. Robust encryption. In Daniele Micciancio, editor, *TCC 2010*, volume 5978 of *LNCS*, pages 480–497. Springer, Heidelberg, February 2010.
- [2] Surabhi Agarwal. India proposes alpha-numeric hash to track WhatsApp chat. *India Times*, 2021.
- [3] Nikolaos Alexopoulos, Aggelos Kiayias, Riivo Talviste, and Thomas Zacharias. MCMix: Anonymous messaging via secure multiparty computation. In Engin Kirda and Thomas Ristenpart, editors, *USENIX Security 2017*, pages 1217–1234. USENIX Association, August 2017.
- [4] 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, *EUROCRYPT 2019, Part I*, volume 11476 of *LNCS*, pages 129–158. Springer, Heidelberg, May 2019.
- [5] Joël Alwen, Sandro Coretti, Yevgeniy Dodis, and Yiannis Tselekounis. Security analysis and improvements for the IETF MLS standard for group messaging. In Daniele Micciancio and Thomas Ristenpart, editors, *CRYPTO 2020, Part I*, volume 12170 of *LNCS*, pages 248–277. Springer, Heidelberg, August 2020.
- [6] Joël Alwen, Sandro Coretti, Daniel Jost, and Marta Mularczyk. Continuous group key agreement with active security. In Rafael Pass and Krzysztof Pietrzak, editors, *TCC 2020, Part II*, volume 12551 of *LNCS*, pages 261–290. Springer, Heidelberg, November 2020.
- [7] Donald Beaver. Efficient multiparty protocols using circuit randomization. In Joan Feigenbaum, editor, *CRYPTO '91*, volume 576 of *LNCS*, pages 420–432. Springer, Heidelberg, August 1992.
- [8] Amos Beimel, Yuval Ishai, and Tal Malkin. Reducing the servers' computation in private information retrieval: PIR with preprocessing. *Journal of Cryptology*, 17(2):125–151, March 2004.
- [9] 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, *CRYPTO 2017, Part III*, volume 10403 of *LNCS*, pages 619–650. Springer, Heidelberg, August 2017.
- [10] Luca Belli. Whatsapp skewed brazilian election, proving social media's danger to democracy. <https://theconversation.com/whatsapp-skewed-brazilian-election-proving-social-medias-danger-to-democracy-106476>, 2018.
- [11] Abhishek Bhowmick, Dan Boneh, Steve Myers, Kunal Talwar, and Karl Tarbe. The Apple PSI system. https://www.apple.com/child-safety/pdf/Apple_PSI_System_Security_Protocol_and_Analysis.pdf, 2021.
- [12] Alexander Bienstock, Yevgeniy Dodis, and Paul Rösler. On the price of concurrency in group ratcheting protocols. In Rafael Pass and Krzysztof Pietrzak, editors, *TCC 2020, Part II*, volume 12551 of *LNCS*, pages 198–228. Springer, Heidelberg, November 2020.
- [13] Alexander Bienstock, Jaiden Fairoze, Sanjam Garg, Pratyay Mukherjee, and Srinivasan Raghuraman. A more complete analysis of the signal double ratchet algorithm. *Cryptology ePrint Archive*, Report 2022/355, 2022. <https://eprint.iacr.org/2022/355>.
- [14] Dan Boneh and Victor Shoup. A graduate course in applied cryptography. <https://toc.cryptobook.us/book.pdf>, 2020.
- [15] Nikita Borisov, Ian Goldberg, and Eric A. Brewer. Off-the-record communication, or, why not to use PGP. In *WPES*, pages 77–84. ACM, 2004.
- [16] Colin Boyd, Anish Mathuria, and Douglas Stebila. *Protocols for Authentication and Key Establishment, Second Edition*. Information Security and Cryptography. Springer, 2020.
- [17] Brazilian fake news draft bill no. 2.630, of 2020. <https://docs.google.com/document/d/1MHMDHsVJB45PI1R51AyoLmZvZk8eULHisYFqGy9X2s>, 2020.

- [18] Sébastien Champion, Julien Devigne, Céline Duguey, and Pierre-Alain Fouque. Multi-device for signal. In Mauro Conti, Jianying Zhou, Emiliano Casalichio, and Angelo Spognardi, editors, *ACNS 20, Part II*, volume 12147 of *LNCS*, pages 167–187. Springer, Heidelberg, October 2020.
- [19] Ran Canetti, Palak Jain, Marika Swanberg, and Mayank Varia. Universally composable end-to-end secure messaging. Cryptology ePrint Archive, Report 2022/376, 2022. <https://eprint.iacr.org/2022/376>.
- [20] Ran Canetti and Hugo Krawczyk. Analysis of key-exchange protocols and their use for building secure channels. In Birgit Pfizmann, editor, *EUROCRYPT 2001*, volume 2045 of *LNCS*, pages 453–474. Springer, Heidelberg, May 2001.
- [21] Ran Canetti, Daniel Shahaf, and Margarita Vald. Universally composable authentication and key-exchange with global PKI. In Chen-Mou Cheng, Kai-Min Chung, Giuseppe Persiano, and Bo-Yin Yang, editors, *PKC 2016, Part II*, volume 9615 of *LNCS*, pages 265–296. Springer, Heidelberg, March 2016.
- [22] Melissa Chase, Trevor Perrin, and Greg Zaverucha. The signal private group system and anonymous credentials supporting efficient verifiable encryption. In Jay Ligatti, Xinming Ou, Jonathan Katz, and Giovanni Vigna, editors, *ACM CCS 2020*, pages 1445–1459. ACM Press, November 2020.
- [23] Katriel Cohn-Gordon, Cas Cremers, Benjamin Dowling, Luke Garratt, and Douglas Stebila. A formal security analysis of the signal messaging protocol. *Journal of Cryptology*, 33(4):1914–1983, 2020.
- [24] Katriel Cohn-Gordon, Cas Cremers, Luke Garratt, Jon Millican, and Kevin Milner. On ends-to-ends encryption: Asynchronous group messaging with strong security guarantees. In David Lie, Mohammad Mannan, Michael Backes, and XiaoFeng Wang, editors, *ACM CCS 2018*, pages 1802–1819. ACM Press, October 2018.
- [25] Katriel Cohn-Gordon, Cas J. F. Cremers, and Luke Garratt. On post-compromise security. In *CSF*, pages 164–178. IEEE Computer Society, 2016.
- [26] Henry Corrigan-Gibbs, Dan Boneh, and David Mazières. Riposte: An anonymous messaging system handling millions of users. In *2015 IEEE Symposium on Security and Privacy*, pages 321–338. IEEE Computer Society Press, May 2015.
- [27] dalek-cryptography. ed25519-dalek. <https://github.com/dalek-cryptography/ed25519-dalek>, 2020.
- [28] Roger Dingledine, Nick Mathewson, and Paul F. Syverson. Tor: The second-generation onion router. In Matt Blaze, editor, *USENIX Security 2004*, pages 303–320. USENIX Association, August 2004.
- [29] Yevgeniy Dodis, Paul Grubbs, Thomas Ristenpart, and Joanne Woodage. Fast message franking: From invisible salamanders to encryption. In Hovav Shacham and Alexandra Boldyreva, editors, *CRYPTO 2018, Part I*, volume 10991 of *LNCS*, pages 155–186. Springer, Heidelberg, August 2018.
- [30] Yevgeniy Dodis, Jonathan Katz, Adam Smith, and Shabsi Wal-fish. Composability and on-line deniability of authentication. In Omer Reingold, editor, *TCC 2009*, volume 5444 of *LNCS*, pages 146–162. Springer, Heidelberg, March 2009.
- [31] Elizabeth Dwoskin and Annie Gowen. On WhatsApp, fake news is fast – and can be fatal. https://www.washingtonpost.com/business/economy/on-whatsapp-fake-news-is-fast--and-can-be-fatal/2018/07/23/a2dd7112-8ebf-11e8-bcd5-9d911c784c38_story.html, July 2018. Accessed: 09-28-2020.
- [32] Facebook. Messenger secret conversations: Technical whitepaper (version 2.0). <https://about.fb.com/wp-content/uploads/2016/07/messenger-secret-conversations-technical-whitepaper.pdf>, 2017.
- [33] Pooya Farshim, Benoît Libert, Kenneth G. Paterson, and Elizabeth A. Quaglia. Robust encryption, revisited. In Kaoru Kurosawa and Goichiro Hanaoka, editors, *PKC 2013*, volume 7778 of *LNCS*, pages 352–368. Springer, Heidelberg, February / March 2013.
- [34] Pooya Farshim, Claudio Orlandi, and Răzvan Roşie. Security of symmetric primitives under incorrect usage of keys. *IACR Trans. Symm. Cryptol.*, 2017(1):449–473, 2017.
- [35] Marc Fischlin. Round-optimal composable blind signatures in the common reference string model. In Cynthia Dwork, editor, *CRYPTO 2006*, volume 4117 of *LNCS*, pages 60–77. Springer, Heidelberg, August 2006.
- [36] Marc Fischlin and Dominique Schröder. On the impossibility of three-move blind signature schemes. In Henri Gilbert, editor, *EUROCRYPT 2010*, volume 6110 of *LNCS*, pages 197–215. Springer, Heidelberg, May / June 2010.
- [37] Georg Fuchsbauer, Christian Hanser, Chethan Kamath, and Daniel Slamanig. Practical round-optimal blind signatures in the standard model from weaker assumptions. In Vassilis Zikas and Roberto De Prisco, editors, *SCN 16*, volume 9841 of *LNCS*, pages 391–408. Springer, Heidelberg, August / September 2016.
- [38] Georg Fuchsbauer, Christian Hanser, and Daniel Slamanig. Practical round-optimal blind signatures in the standard model. In Rosario Gennaro and Matthew J. B. Robshaw, editors, *CRYPTO 2015, Part II*, volume 9216 of *LNCS*, pages 233–253. Springer, Heidelberg, August 2015.
- [39] Paul Grubbs, Jiahui Lu, and Thomas Ristenpart. Message franking via committing authenticated encryption. In Jonathan Katz and Hovav Shacham, editors, *CRYPTO 2017, Part III*, volume 10403 of *LNCS*, pages 66–97. Springer, Heidelberg, August 2017.
- [40] Christoph G. Günther. An identity-based key-exchange protocol. In Jean-Jacques Quisquater and Joos Vandewalle, editors, *EUROCRYPT’89*, volume 434 of *LNCS*, pages 29–37. Springer, Heidelberg, April 1990.
- [41] Hecate Rust implementation. <https://github.com/Ralissa/hecate>, 2022.
- [42] Hecate’s modified libsignal-client implementation. <https://github.com/Ralissa/libsignal-client/>, 2022.
- [43] Hecate’s modified signal-cli implementation. <https://github.com/Ralissa/signal-cli>, 2022.
- [44] Rawane Issa, Nicolas AlHaddad, and Mayank Varia. Hecate: Abuse reporting in secure messengers with sealed sender.

- Cryptology ePrint Archive, Report 2021/1686, 2021. <https://eprint.iacr.org/2021/1686>.
- [45] Joseph Jaeger and Igors Stepanovs. Optimal channel security against fine-grained state compromise: The safety of messaging. In Hovav Shacham and Alexandra Boldyreva, editors, *CRYPTO 2018, Part I*, volume 10991 of *LNCS*, pages 33–62. Springer, Heidelberg, August 2018.
- [46] Markus Jakobsson, Kazue Sako, and Russell Impagliazzo. Designated verifier proofs and their applications. In Ueli M. Maurer, editor, *EUROCRYPT'96*, volume 1070 of *LNCS*, pages 143–154. Springer, Heidelberg, May 1996.
- [47] Daniel Jost, Ueli Maurer, and Marta Mularczyk. Efficient ratcheting: Almost-optimal guarantees for secure messaging. In Yuval Ishai and Vincent Rijmen, editors, *EUROCRYPT 2019, Part I*, volume 11476 of *LNCS*, pages 159–188. Springer, Heidelberg, May 2019.
- [48] Daniel Jost, Ueli Maurer, and Marta Mularczyk. A unified and composable take on ratcheting. In Dennis Hofheinz and Alon Rosen, editors, *TCC 2019, Part II*, volume 11892 of *LNCS*, pages 180–210. Springer, Heidelberg, December 2019.
- [49] Seny Kamara, Mallory Knodel, Emma Llansó, Greg Nojeim, Lucy Qin, Dhanaraj Thakur, and Caitlin Vogus. Outside looking in: Approaches to content moderation in end-to-end encrypted systems. <https://cdt.org/wp-content/uploads/2021/08/CDT-Outside-Looking-In-Approaches-to-Content-Moderation-in-End-to-End-Encrypted-Systems.pdf>, 2021.
- [50] Shuichi Katsumata, Ryo Nishimaki, Shota Yamada, and Takashi Yamakawa. Round-optimal blind signatures in the plain model from classical and quantum standard assumptions. In Anne Canteaut and François-Xavier Standaert, editors, *EUROCRYPT 2021, Part I*, volume 12696 of *LNCS*, pages 404–434. Springer, Heidelberg, October 2021.
- [51] Jonathan Katz and Yehuda Lindell. *Introduction to modern cryptography*. CRC press, 2020.
- [52] Anunay Kulshrestha and Jonathan R. Mayer. Identifying harmful media in end-to-end encrypted communication: Efficient private membership computation. In Michael Bailey and Rachel Greenstadt, editors, *USENIX Security 2021*, pages 893–910. USENIX Association, August 2021.
- [53] Brian A. LaMacchia, Kristin Lauter, and Anton Mityagin. Stronger security of authenticated key exchange. In Willy Susilo, Joseph K. Liu, and Yi Mu, editors, *ProvSec 2007*, volume 4784 of *LNCS*, pages 1–16. Springer, Heidelberg, November 2007.
- [54] Iraklis Leontiadis and Serge Vaudenay. Private message franking with after opening privacy. Cryptology ePrint Archive, Report 2018/938, 2018. <https://eprint.iacr.org/2018/938>.
- [55] Linsheng Liu, Daniel S. Roche, Austin Theriault, and Arkady Yerukhimovich. Fighting fake news in encrypted messaging with the fuzzy anonymous complaint tally system (FACTS). In *NDSS*. The Internet Society, 2022.
- [56] Manish Singh. Whatsapp is now delivering roughly 100 billion messages a day. <https://techcrunch.com/2020/10/29/whatsapp-is-now-delivering-roughly-100-billion-messages-a-day/>, 2020.
- [57] Ian Martiny, Gabriel Kaptchuk, Adam Aviv, Dan Roche, and Eric Wustrow. Improving Signal’s sealed sender. In *NDSS*. The Internet Society, 2021.
- [58] Oversight Board. Ensuring respect for free expression, through independent judgment. <https://oversightboard.com/>, 2021.
- [59] Charlotte Peale, Saba Eskandarian, and Dan Boneh. Secure complaint-enabled source-tracking for encrypted messaging. In Giovanni Vigna and Elaine Shi, editors, *ACM CCS 2021*, pages 1484–1506. ACM Press, November 2021.
- [60] Riana Pfefferkorn. Content-oblivious trust and safety techniques: Results from a survey of online service providers. <https://ssrn.com/abstract=3920031>, 2021.
- [61] Bertram Poettering and Paul Rösler. Towards bidirectional ratcheted key exchange. In Hovav Shacham and Alexandra Boldyreva, editors, *CRYPTO 2018, Part I*, volume 10991 of *LNCS*, pages 3–32. Springer, Heidelberg, August 2018.
- [62] Newley Purnell and Jeff Horowitz. WhatsApp says it filed suit in India to prevent tracing of encrypted messages. <https://www.wsj.com/articles/whatsapp-says-it-filed-suit-in-india-to-prevent-tracing-of-encrypted-messages-11622000307>, 2021.
- [63] Paul Rösler, Christian Mainka, and Jörg Schwenk. More is less: On the end-to-end security of group chats in Signal, WhatsApp, and Threema. In *EuroS&P*, pages 415–429. IEEE, 2018.
- [64] Shahrokh Saeednia, Steve Kremer, and Olivier Markowitch. An efficient strong designated verifier signature scheme. In Jong In Lim and Dong Hoon Lee, editors, *ICISC 03*, volume 2971 of *LNCS*, pages 40–54. Springer, Heidelberg, November 2004.
- [65] Sarah Scheffler and Jonathan Mayer. SoK: content moderation in end-to-end encryption, 2022.
- [66] Michael Schliep and Nicholas Hopper. End-to-end secure mobile group messaging with conversation integrity and deniability. In *WPES@CCS*, pages 55–73. ACM, 2019.
- [67] Sebastian Scheibner. signal-cli. <https://github.com/AsamK/signal-cli>.
- [68] Signal. Technology preview: Sealed sender for signal. <https://signal.org/blog/sealed-sender/>, 2018.
- [69] Signal. libsignal-client. <https://github.com/signalapp/libsignal-client>, 2020.
- [70] Signal. Technical information. <https://signal.org/docs/>, 2021.
- [71] Signal. Signal-server. <https://github.com/signalapp/Signal-Server>, 2022.
- [72] Kate Starbird. Online rumors, misinformation and disinformation: The perfect storm of covid-19 and election2020. In *Enigma 2021*. USENIX Association, February 2021.
- [73] Li Q. Tay, Mark J. Hurlstone, Tim Kurz, and Ullrich K. H. Ecker. A comparison of prebunking and debunking interventions for implied versus explicit misinformation. PsyArXiv, 2021.
- [74] Kurt Thomas, Devdatta Akhawe, Michael Bailey, Dan Boneh, Elie Bursztein, Sunny Consolvo, Nicola Dell, Zakir Durumeric, Patrick Gage Kelley, Deepak Kumar, Damon McCoy, Sarah

Meiklejohn, Thomas Ristenpart, and Gianluca Stringhini. SoK: Hate, harassment, and the changing landscape of online abuse. In *2021 IEEE Symposium on Security and Privacy*, pages 247–267. IEEE Computer Society Press, May 2021.

- [75] Nirvan Tyagi, Paul Grubbs, Julia Len, Ian Miers, and Thomas Ristenpart. Asymmetric message franking: Content moderation for metadata-private end-to-end encryption. In Alexandra Boldyreva and Daniele Micciancio, editors, *CRYPTO 2019, Part III*, volume 11694 of *LNCS*, pages 222–250. Springer, Heidelberg, August 2019.
- [76] Nirvan Tyagi, Ian Miers, and Thomas Ristenpart. Traceback for end-to-end encrypted messaging. In Lorenzo Cavallaro, Johannes Kinder, XiaoFeng Wang, and Jonathan Katz, editors, *ACM CCS 2019*, pages 413–430. ACM Press, November 2019.
- [77] Nik Unger, Sergej Dechand, Joseph Bonneau, Sascha Fahl, Henning Perl, Ian Goldberg, and Matthew Smith. SoK: Secure messaging. In *2015 IEEE Symposium on Security and Privacy*, pages 232–249. IEEE Computer Society Press, May 2015.
- [78] Jelle van den Hooff, David Lazar, Matei Zaharia, and Nikolai Zeldovich. Vuvuzela: scalable private messaging resistant to traffic analysis. In *SOSP*, pages 137–152. ACM, 2015.
- [79] Soroush Vosoughi, Deb Roy, and Sinan Aral. The spread of true and false news online. *Science*, 359(6380):1146–1151, 2018.
- [80] WhatsApp. Two billion users – connecting the world privately. <https://blog.whatsapp.com/two-billion-users-connecting-the-world-privately/>, 2020.
- [81] Liang Wu, Fred Morstatter, Kathleen M Carley, and Huan Liu. Misinformation in social media: definition, manipulation, and detection. *ACM SIGKDD Explorations Newsletter*, 21(2):80–90, 2019.

A Definitions of Cryptographic Building Blocks

This work uses four standard cryptographic building blocks that we use and adapt from Boneh-Shoup [14] and Katz-Lindell [51]. In what follows, we define the message space as $\mathcal{M} := \{0, 1\}^*$, the key space as $\mathcal{K} := \{0, 1\}^n$, the ciphertext space as $\mathcal{C} := \{0, 1\}^*$, the randomness space $\mathcal{R} := \{0, 1\}^n$ and the signature space as $\Sigma := \{0, 1\}^n$, where n denotes the security parameter.

Definition 2 (Commitment scheme). *A non-interactive commitment scheme is defined by two algorithms Com and Vf.*

- Com is an algorithm that takes a random string $r \leftarrow \mathcal{R}$, and a plaintext message $m \in \mathcal{M}$ and outputs a commitment $\text{com} := \text{Com}(m, r)$.
- Vf is an algorithm that takes a commitment com, a string r and a plaintext message m and checks if $\text{Vf}(m, \text{com}, r) := (\text{Com}(m, r) \stackrel{?}{=} \text{com})$.

Definition 3 (Binding commitment). *A commitment scheme $\pi = \{\text{Com}, \text{Vf}\}$ is computationally binding if for all probabilistic polynomial time (PPT) adversaries \mathcal{A} , there is a*

negligible function $\text{negl}(n)$ such that:

$$\text{Adv}_{\pi}^{\text{binding-com}}(\mathcal{A}) = \Pr[\text{Com}(m, r, \text{param}) = \text{Com}(m', r', \text{param}) \mid m \neq m'] \leq \text{negl}(n).$$

Definition 4 (Hiding commitment). *Let $\pi = \{\text{Com}, \text{Vf}\}$ be a commitment scheme. Let $\text{Com}_{\text{hiding}}^{\mathcal{A}}$ be defined by the following experiment:*

- The adversary \mathcal{A} outputs a pair of messages $m_0, m_1 \in \mathcal{M}$.
- A uniform bit $b \in \{0, 1\}$ and the randomness $r \leftarrow \{0, 1\}^n$ are chosen.
- The adversary \mathcal{A} is given access to the commitment oracle $\text{O}^{\text{com-hiding}}$ which on messages m_0 and m_1 computes and returns the commitment $\text{com} \leftarrow \text{Com}(m_b, r)$, where $\text{Vf}(\text{Com}(m_b, r), m_b, r) = 1$.
- The output of the experiment is 1 if $b' = b$ and 0 otherwise.

A commitment scheme π is computationally hiding if for all PPT adversaries \mathcal{A} there is a negligible function $\text{negl}(n)$ such that:

$$\text{Adv}_{\pi}^{\text{hiding-com}}(\mathcal{A}) = \Pr[\text{Com}_{\text{hiding}}^{\mathcal{A}}(n) = 1] \leq \frac{1}{2} + \text{negl}(n).$$

Definition 5 (Encryption scheme). *A randomized private key encryption scheme is defined by three algorithms EncKGen, Enc and Dec over a finite message space \mathcal{M} .*

- EncKGen is a probabilistic key generation algorithm that outputs a key sk sampled uniformly at random from \mathcal{K} .
- Enc is the randomized encryption algorithm that takes as an input sk and plaintext message $m \in \mathcal{M}$ and outputs $c := \text{Enc}_{\text{sk}}(m)$ where $c \in \mathcal{C}$.
- Dec is the decryption algorithm that takes as an input sk and a ciphertext c in the ciphertext space \mathcal{C} and outputs a plaintext message $m := \text{Dec}_{\text{sk}}(c)$ such that $c := \text{Enc}_{\text{sk}}(m)$.

Definition 6 (CCA security). *Let $\pi = \{\text{EncKGen}, \text{Enc}, \text{Dec}\}$ be an encryption scheme. Let $\text{ENC}_{\text{cca}, \pi}^{\mathcal{A}}(n)$ denote the following experiment:*

- EncKGen is run to obtain (pk, sk) and a uniform bit $b \in \{0, 1\}$ is chosen. The adversary \mathcal{A} is given pk .
- The adversary \mathcal{A} is given access to the encryption oracle $\text{O}_{\text{cca}}^{\text{enc}}$ which, on messages m_0, m_1 , outputs a ciphertext $c \leftarrow \text{Enc}_{\text{pk}}(m_b)$.
- The adversary \mathcal{A} is given access to the decryption oracle $\text{O}_{\text{cca}}^{\text{decrypt}}$ which outputs the decrypted plaintext message m under sk when handed out a ciphertext c' .
- \mathcal{A} continues to interact with the decryption and encryption oracles, but may not request a decryption of any ciphertext c returned by $\text{O}_{\text{cca}}^{\text{enc}}$.

- Finally \mathcal{A} output a bit b' . The output of the experiment is defined to be 1 if $b = b'$, and 0 otherwise.

We say that π is secure under a chosen-ciphertext attack (CCA) if for all PPT adversaries \mathcal{A} , there exists a negligible function $\text{negl}(n)$ such that:

$$\text{Adv}_{\pi}^{\text{enc}_{\text{cca}}} = \Pr[\text{ENC}_{\text{cca},\pi}^{\mathcal{A}}(n) = 1] \leq \frac{1}{2} + \text{negl}(n).$$

Definition 7 (CPA security). Let $\pi = \{\text{EncKGen}, \text{Enc}, \text{Dec}\}$ be an encryption scheme. Let $\text{ENC}_{\text{cpa},\pi}^{\mathcal{A}}(n)$ denote a similar experiment to $\text{ENC}_{\text{cca},\pi}^{\mathcal{A}}(n)$ where the adversary \mathcal{A} only has access to the encryption oracle that rename as $\mathcal{O}_{\text{cpa}}^{\text{enc}}$. We say that π is secure under a chosen-plaintext attack (CPA) if for all PPT adversaries \mathcal{A} , there exists a negligible function $\text{negl}(n)$ such that:

$$\text{Adv}_{\pi}^{\text{enc}_{\text{cpa}}} = \Pr[\text{ENC}_{\text{cpa},\pi}^{\mathcal{A}}(n) = 1] \leq \frac{1}{2} + \text{negl}(n).$$

Definition 8 (Digital signature scheme). A signature scheme is defined as the triple of algorithms $\text{SigKGen}, \text{Sign}, \text{Vf}$ over the message space \mathcal{M} and the signature space Σ .

- SigKGen is a probabilistic key generation algorithm that output a key pair (pk, sk) sampled from \mathcal{K} , where pk is the public verification key and sk is the secret signing key.
- Sign is the probabilistic signing algorithm that takes the signing key sk and a plaintext message m and outputs a signature $\sigma \leftarrow \text{Sign}_{\text{sk}}(m)$, where $\sigma \in \Sigma$.
- Vf is the deterministic verification algorithm which checks the signature σ against the plaintext message m and public key pk and outputs \perp or 1 such that:

$$\Pr[\text{Vf}(\text{pk}, m, \text{Sign}_{\text{sk}}(m)) = 1] = 1.$$

Definition 9 (EU-CMA security). Let $\pi = \{\text{SigKGen}, \text{Sign}, \text{Vf}\}$ denote a digital signature scheme. Let $\text{Sig}_{\text{eu-cma}}^{\mathcal{A}}$ be the experiment defined as:

- SigKGen is run to obtain (pk, sk) .
- The adversary \mathcal{A} is given pk and access to the signing oracle $\mathcal{O}_{\text{eu-cma}}^{\text{sign}}$ which on message m computes and outputs the signature σ of that message under the secret signing key sk . Let Q denote the set of all queries that \mathcal{A} makes to $\mathcal{O}_{\text{eu-cma}}^{\text{sign}}$.
- The adversary \mathcal{A} then outputs (m', σ') .
- The experiment outputs 1 if and only if $\text{Vf}_{\text{pk}}(m', \sigma') = 1$ and $m' \notin Q$, and 0 otherwise.

We say that π is existentially unforgeable under an adaptive chosen-message attack (EU-CMA) if for all PPT adverbs \mathcal{A} , there is a negligible function $\text{negl}(n)$ such that:

$$\text{Adv}_{\pi}^{\text{sig}_{\text{eu-cma}}}(\mathcal{A}) = \Pr[\text{Sig}_{\text{eu-cma}}^{\mathcal{A}}(n) = 1] \leq \text{negl}(n).$$

Definition 10 (Hash function). A hash function with output length l is defined by two algorithms Gen and H .

- Gen is a probabilistic algorithm which outputs a key $k \in \mathcal{K}$.
- H is an algorithm which takes as input as key k and a string $m \in \mathcal{M}$ and outputs a string $\text{H}_k(m) \in \{0, 1\}^{l(n)}$.

Definition 11 (Collision resistance). Let $\pi = \{\text{Gen}, \text{H}\}$ denote a hash function. Let $\text{Hash}_{\text{coll}}^{\mathcal{A}}$ be the experiment defined as:

- Gen is run to obtain k .
- The adversary \mathcal{A} is given access to the hashing oracle $\mathcal{O}^{\text{hash}}$ which on input m returns $\text{H}_k(m)$.
- The adversary then outputs m_0 and m_1 .
- The experiment outputs 1 if and only if $m_0 \neq m_1$ and $\text{H}_k(m_0) = \text{H}_k(m_1)$.

We say that π is collision resistant if for all PPT adversaries \mathcal{A} there is a negligible function $\text{negl}(n)$ such that:

$$\text{Adv}_{\pi}^{\text{hash}_{\text{coll}}}(\mathcal{A}) = \Pr[\text{Hash}_{\text{coll}}^{\mathcal{A}}(n) = 1] \leq \text{negl}(n).$$