

Secure Complaint-Enabled Source-Tracking for Encrypted Messaging

Charlotte Peale
Stanford University

Saba Eskandarian
UNC Chapel Hill

Dan Boneh
Stanford University

ABSTRACT

While the end-to-end encryption properties of popular messaging schemes such as Whatsapp, Messenger, and Signal guarantee privacy for users, these properties also make it very difficult for messaging platforms to enforce any sort of content moderation. This can lead to the unchecked spread of malicious content such as misinformation on such platforms. In 2019, Tyagi et al. developed *message traceback*, which addresses this issue by allowing a messaging platform to recover the path of a forwarded message after a user reports it for malicious content. This paper presents an alternative to message traceback that offers more privacy to users and requires less platform-side storage. We term this approach *source-tracking* for encrypted messaging schemes. Source-tracking enables messaging platforms to provide the privacy guarantees expected from standard end-to-end encryption, but also helps hold the sources of malicious messages accountable: if malicious content is reported by a user, the source can be identified. We formalize security goals for source-tracking schemes and design and implement two source-tracking schemes with different security and performance tradeoffs.

CCS CONCEPTS

• **Security and privacy** → **Cryptography**; **Security services**; **Privacy-preserving protocols**; **Pseudonymity, anonymity and untraceability**;

KEYWORDS

tracing; abuse reporting; message franking; end-to-end encrypted messaging; source-tracking;

ACM Reference Format:

Charlotte Peale, Saba Eskandarian, and Dan Boneh. 2021. Secure Complaint-Enabled Source-Tracking for Encrypted Messaging. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security (CCS '21), November 15–19, 2021, Virtual Event, Republic of Korea*. ACM, New York, NY, USA, 23 pages. <https://doi.org/10.1145/3460120.3484539>

1 INTRODUCTION

End-to-end encrypted messaging apps like WhatsApp and Signal provide users with strong privacy guarantees and deliver billions of messages a day [29]. Concerns about abuse and misinformation on such platforms have led to work on technical solutions for verifiable

reporting of messages with no impact on the privacy guarantees for messages that are not reported [12, 14, 31, 32]. Techniques for enabling some degree of moderation with minimal impact on privacy are all the more important as Brazil and India consider laws that would require messaging providers to be able to identify the sources of misinformation [26–28, 30]. In the absence of strong privacy-preserving tools to minimize the impact of such policies, these requirements would cripple the hard-won privacy victories of end-to-end encrypted messaging.

The majority of work on verifiably reporting abusive messages has been in the context of *message franking* [12, 14, 31], which focuses on reporting the immediate sender of a message and not on identifying the sources of viral misinformation campaigns. While message franking techniques are important for reporting abusive messages, they do not help find the origin of a message that may be forwarded many times before someone reports it. Although it would be possible to find the source of a message by combining a message franking scheme with extensive metadata collection and retention, e.g., a graph of all message sources and destinations with forwarded messages labeled, this approach requires a great deal of storage and forces the platform to collect and keep lots of metadata, which may compromise user privacy in its own right. Solutions relying on having users sign their messages are also unsatisfactory, as this damages deniability and may reveal the identity of the original sender of a message even before it is reported.

To our knowledge, the only work to consider tracing the source of a reported message is the *traceback* scheme of Tyagi et al. [32]. Traceback boasts extremely low client-side costs and reveals the set of users who have received a reported message without revealing any additional information about a message's source before it is reported. Unfortunately, revealing all the recipients of a message may be leaking too much if the goal of reporting is only to find the source of the message. Moreover, tracing a message back to its source takes time linear in the length of the forwarding chain and requires a 32 Byte token to be stored by the server for each message sent, regardless of whether that message is ever reported. The token must be kept for as long as the message can be traced. For a high-volume messaging platform like WhatsApp, this translates into terabytes of additional long-term storage each day. Thus traceback removes the need for a platform to collect message metadata, but it does not reduce the large quantity of data the platform must store.

This paper introduces the notion of *source-tracking* for end-to-end encrypted messaging schemes. Like traceback schemes, source-tracking schemes reveal the source of a reported message, but they do not reveal the identities of users who received a message between when it was first sent and when it was reported. Moreover, our schemes do not require the server to store any data for each message sent, and the time to trace the source of a reported message does not depend on the length of the forwarding chain.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CCS '21, November 15–19, 2021, Virtual Event, Republic of Korea.

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-8454-4/21/11...\$15.00

<https://doi.org/10.1145/3460120.3484539>

We begin by formalizing the notion of a source-tracking scheme and providing security definitions that capture and exceed the privacy requirements of prior message franking and traceback schemes. We consider confidentiality of an unreported message's origin from both the platform and other users who may be forwarded the message. In particular, we offer two notions of confidentiality, *tree-linkable* and *tree-unlinkable* confidentiality. In a tree-linkable scheme, the recipient of two messages with identical plaintexts can tell if the two messages have the same origin. A tree-unlinkable scheme hides even this information.

Our definitions also require accountability, unforgeability, and deniability. These properties require that every successfully delivered message must be able to be traced to its source, that nobody can be framed for originating a message, and that only the platform can verify the origin of a reported message.

Next, we provide two constructions of secure source-tracking schemes, one with tree-linkable security and one with tree-unlinkable security. We implement and evaluate our schemes, comparing performance both to the traceback scheme of Tyagi et al. [32] and an implementation of end-to-end encryption based on the double ratchet protocol [25, 38] without any additional functionality for finding the source of a reported message.

In exchange for its slightly weaker security guarantee, the tree-linkable scheme offers exceptional performance, introducing very reasonable overheads compared to an end-to-end encrypted messaging service with no source-tracking capability. For example, it takes under $60\mu\text{s}$ to find the source of a reported message. The additional server-side overhead of using our tree-linkable scheme is only $20\mu\text{s}$ of cryptographic operations per message delivered. The scheme only relies on standard cryptographic primitives such as signatures and commitments, re-using the security of the underlying end-to-end encrypted messaging service to provide additional security properties at low cost.

Our tree-unlinkable construction, while slower than our first construction, requires only milliseconds to send and receive messages. At a high level, we unlinkably hide the source of a message using a heavily modified version of the algebraic MAC and anonymous credential techniques used in the new Signal group messaging system [4, 5].

Our implementation and raw evaluation data are free and open source and can be accessed at <https://github.com/cpeale/srctracking>.

In summary, this paper makes the following contributions:

- Introduces the notion of *source-tracking* for end-to-end encrypted messaging systems and defines strong security requirements for source-tracking schemes.
- Constructs a *tree-linkable* source-tracking scheme with exceptional performance that can easily be deployed on top of existing end-to-end encrypted messaging services.
- Constructs a *tree-unlinkable* source-tracking scheme which requires milliseconds to process and deliver messages.
- Implements and evaluates both schemes, including comparisons to both prior work and a baseline messaging implementation without source-tracking functionality.

2 DESIGN GOALS

A source-tracking scheme builds on top of an end-to-end encrypted messaging system and allows a user who receives an abusive message to report that message to the messaging platform. This is achieved by augmenting the usual protocol for sending and receiving messages as well as by adding a new Report feature. Given the contents of a report, the messaging platform should be able to uncover the true author of the message, regardless of how many times the message has been forwarded. Source-tracking schemes do not identify objectionable content on their own (as in, e.g., [16]), they only deal with the process of identifying the source of content that a user has reported.

Security goals. Informally, we require the following security properties from a secure source-tracking scheme. We formalize these properties in Section 3.

- **Confidentiality:** Neither the platform nor any user not involved in sending or receiving a given message should learn anything about that message's forwarding history or contents. The scheme should also preserve unlinkability: a user who receives two different forwarded messages should not be able to tell if they were authored by the same user. After a report, the messaging platform should not learn anything new about the reported message's forwarding history other than the identity of the source.
- **Accountability:** No malicious user should be able to send a message which cannot later be traced back to them.
- **Unforgeability:** No malicious user should be able to frame another user for sending a message it did not send.
- **Deniability:** Only the platform can verify a report, so if the contents of a report are leaked, it would be impossible for a third party to verify that the report implicated a particular user.

Tree-linkable and unlinkable confidentiality. We consider two levels of confidentiality: *tree-linkable* and *tree-unlinkable*. When a message is authored and then forwarded, a *forwarding tree* is created consisting of all the users who have received and forwarded that message, rooted at the author of the message. In the strictest form of user confidentiality, a recipient shouldn't be able to tell which forwarding tree a message came from. If a user receives a forward consisting of plaintext m , and then later receives another forward of the same plaintext m , the user shouldn't be able to tell whether the two messages are from the same forwarding tree or from different trees whose messages happen to have the same plaintext. If a scheme satisfies this form of confidentiality, we will say it has *tree-unlinkable* confidentiality.

In some situations, this notion of confidentiality may be unnecessarily strong, as whether or not two identical messages are on the same forwarding tree may not be sensitive information. Additionally, depending on the exact implementation of a messaging scheme, other metadata associated with the message (such as a timestamp) would mean that no two forwarding trees could ever correspond to the same message plaintext. If a scheme is such that a recipient can determine whether two forwards came from the same tree, but can still learn nothing about the tree structure or other messaging activity, we say that it satisfies *tree-linkable* user confidentiality.

It is worth noting, however, that there are circumstances where a stronger tree-unlinkable security guarantee may be preferable.

For example, a malicious insider working for a messaging platform can send the same message to multiple recipients. Later, if one of those recipients sends the message to someone who reports it, the sender will know, from tree-linkability, which person's version of the message was reported. Although we believe tree-linkability suffices for most scenarios, we design and evaluate constructions of both tree-linkable and tree-unlinkable schemes.

Accommodating report metadata. In order to properly handle a report, a messaging platform may wish to receive additional information about a reported message beyond the identity of its original author, e.g., an exact or approximate timestamp of when the message was written. Our schemes support this functionality by allowing messaging platforms to optionally add report metadata to messages such that the metadata associated with a message when it is originally authored will be revealed if and only if that message is reported. We note that this feature is intended for including timestamps or other associated data, and, like any cryptosystem, it is possible for it to be abused in ways that damage security. For example, the platform could damage deniability by including a unique id in each message's metadata and keeping a public database of hashes of ids used in authentic messages.

Limitations. Any system aimed at finding the source of a forwarded message faces certain non-technical limitations. Both our scheme and prior work on traceback [32] track the source of messages that are forwarded using a messaging app's "forward" feature. However, a user could also re-type a message instead of forwarding it using the app, at which point the message appears to be a fresh message which happens to have the same content. Thus the actual guarantee of a scheme that finds the source of a message is that it finds the first person to send the message that was later forwarded, not its real-world author.

As discussed in [15], a more significant issue raised by schemes that find the source of a message is that they allow a user to reveal information about someone other than their immediate conversation partner. This is a fundamental component of the functionality, but it should be weighed carefully whether the potential moderation benefit of such a feature outweighs the privacy risk. Part of our goal in designing source-tracking is to minimize the risk posed compared to traceback. For example, consider the case where a journalist reports a piece of misinformation to a platform. A traceback scheme would, in the process of revealing the source of the misinformation, also reveal the identity of the user who forwarded it to the journalist. Source-tracking does not reveal the path taken by a message, so only the original sender of the message would be revealed.

Finally, both our work and traceback rely on the platform being aware of the identity of the parties sending and receiving a message at the time the message is sent. This information is currently required by most end-to-end encrypted messaging systems in order to deliver messages, so our schemes do not impose additional metadata-collection or retention requirements on these systems. However, it remains an important problem for future work to develop source-tracking schemes that do not even require this minimal metadata at the time of message delivery.

3 FORMALIZING SOURCE-TRACKING

We now formalize the syntax and security definitions for a source-tracking scheme. We begin by describing some notation we will use throughout the rest of the paper.

Let $x \leftarrow F(y)$ denote the assignment of the output of $F(y)$ to x , and let $x \xleftarrow{\mathcal{R}} S$ denote assignment to x of an element sampled uniformly random from set S . We use \mathcal{A}^H to denote that \mathcal{A} has oracle access to some function H . A function $\text{negl}(x)$ is *negligible* if for all $c > 0$, there is a x_0 such that for all $x > x_0$, $\text{negl}(x) < \frac{1}{x^c}$. We omit x if the parameter is implicit. Throughout the paper we also omit an implicit security parameter λ . Finally, we use \perp to indicate an empty message or special character indicating failure.

Interactive protocols. We define an interaction between two parties using the notation

$$(out_1, out_2) \leftarrow \langle P_1(\text{secret params}), P_2(\text{secret params}) \rangle(\text{public params})$$

Where the first party acts according to the protocol defined by P_1 and has access to the secret parameters of P_1 as well as all public parameters, and the second party acts according to P_2 with access to the public parameters as well as its analogous secrets.

Tables and sets. Our security definitions use tables to keep track of important information about adversary queries. Tables are denoted with a capital T and a subscript name, and store key/value pairs. To add a key/value pair to a table, we use the notation $T[key] \leftarrow value$. To retrieve a value corresponding to a particular key from a table, we use the notation $value \rightarrow T[key]$. We use standard set notation to check if a key is included in a table ($key \in T$). Sets are similar to tables, but only store a set of values. We add values v_1, \dots, v_k to a set \mathcal{S} with $\mathcal{S}.add(v_1, \dots, v_k)$.

3.1 Source-Tracking Syntax and Correctness

A source-tracking scheme ST consists of six algorithms: KGen, NewUser, AuthMsg, FwdMsg, RecMsg, and Report, where the latter five are interactive protocols between a user and the messaging platform. Source-tracking involves users U_1, \dots, U_n , each represented by a unique identifier, e.g., a username or long-term public key taken from some set \mathcal{U} , and a platform P .

In contrast with prior work, which operates independently of the underlying end-to-end encrypted messaging scheme [32], we allow source-tracking schemes to make use of the underlying messaging scheme as a black box. We see this as striking a happy medium where source-tracking can take advantage of the security benefits already provided by the messaging scheme while still being easy to deploy on top of existing applications. Although modern end-to-end encrypted messaging protocols provide many useful security properties [1, 7, 25], our constructions will only rely on them to provide authenticated encryption and protect against replay attacks. To abstract away the details of the underlying messaging scheme, we will assume the existence of a messaging oracle

$$\mathcal{E} = (\text{send}(m, U_s, U_r), \text{receive}(ct, U_s, U_r))$$

that sends and authenticates messages between users using the underlying end-to-end encryption scheme the source-tracking scheme is associated with.

The syntax of a source-tracking scheme is defined as follows:

- $KGen(pp) \rightarrow (pk, sk)$: The platform runs this algorithm at system setup. It takes in public parameters pp and outputs a platform key pair (pk, sk) .
- $NewUser$ is an interactive protocol between a new user and the messaging platform to register that user in the system. It is represented by the pair of algorithms U_{new}, P_{new} :

$$(ad, \mathcal{U}') \leftarrow \langle U_{new}, P_{new}(sk, \mathcal{U}) \rangle (U_n, pk)$$

Where \mathcal{U} is the set of currently registered users, (pk, sk) is the platform's key pair, and U_n is the user information or id associated with the new user. We assume that the user runs a separate protocol to register with the underlying messaging scheme as needed. Optionally, U_{new} can add additional private inputs, e.g., a secret used to prove its identity.

On success, the protocol should return some authoring data ad to the user that can be used to send messages and the updated membership set \mathcal{U}' to the platform.

While we define this function interactively in our security definitions for maximum flexibility, a non-interactive function run by the platform suffices for both of our schemes. For ease of notation in the first scheme, which uses no authoring data, we will denote the non-interactive version of this protocol as the function $newUser(U_i, sk)$, run by the platform, which returns only the new set of users \mathcal{U} or \perp on failure.

- $AuthMsg$ is an interactive protocol between a user authoring a message and the platform, represented by the pair of interactive algorithms

$$((ad', e), (pd, e)) \leftarrow \langle U_{auth}(msg), P_{send}(sk, md) \rangle (U_s, U_r, pk)$$

where msg is a tuple (m, ad) , m is the plaintext being sent, and ad ("authoring data") is any associated data required to send a message.

The platform has secret inputs sk , its secret key, and report metadata md that it wants to recover if the message is later reported. Both algorithms receive the identities of the sender and receiver as well as the platform public key. Upon successful completion of the protocol, the sender gets updated authoring data ad' that it can store for authoring future messages, and the platform receives platform data pd to be used when the message is delivered to its recipient. Both parties receive the message identifier e for the message sent using the underlying messaging platform \mathcal{E} .

- $FwdMsg$ is an interactive protocol between a user forwarding a message and the platform, represented by interactive algorithms

$$((fd', e), (pd, e)) \leftarrow \langle U_{fwd}(msg), P_{send}(sk, md) \rangle (U_s, U_r, pk)$$

where each algorithm's inputs and outputs are identical to $AuthMsg$ except that $msg = (m, fd)$, where fd ("forwarding data") represents associated data required to forward a message, and the user's output is updated forwarding data fd' to be used if the user wishes to forward the same message again.

We note that because the platform shouldn't be able to distinguish between an authored and a forwarded message, the platform protocol P_{send} is identical to the platform's protocol for an authored message. The report metadata (md) is only included if the message is new, so in this case it is passed to the forwarding function but remains unused.

- $RecMsg$ is an interactive protocol between a user receiving a message and the platform, represented by the pair of interactive algorithms

$$((m, fd), \perp) \leftarrow \langle U_{rec}, P_{rec}(sk, pd) \rangle (U_s, U_r, e, pk)$$

The platform has access to its secret key sk and the transaction data pd generated when the message was sent (from either $AuthMsg$ or $FwdMsg$), and both user and platform have access to the sending and receiving users' identities, the message identifier for the message sent by the underlying message scheme, and the platform's public key. Upon successful completion, the receiving user gets the message plaintext m and forwarding data fd that it can use to forward or report the message in the future.

- $Report$: This is an interactive protocol between a user who would like to report a received message and the platform, represented by the pair of interactive algorithms

$$(fd', (U, md)) \leftarrow \langle U_{rep}(fd), P_{rep}(sk) \rangle (m, pk)$$

where the user knows the forwarding data fd for the message it would like to report, the platform knows the platform secret key, and both participants have access to the platform public key and the plaintext m of the message being reported. Upon successful completion, the platform gets the source user identity U and associated metadata md for the reported message. Optionally, the reporting user could get new forwarding data fd' used to forward the message again in the future, but we do not require use of this in either of our schemes.

We define correctness for a source-tracking scheme as follows.

Definition 3.1. A source-tracking scheme ST is *correct* if when all users and the platform are honest (i.e. follow the protocols and don't try to forward or report messages that they have not received), all protocols will fail with zero probability, messages will be delivered as intended, and if a user reports a message to the platform, the platform can recover the identity of the original author of the message as well as the metadata that the platform included at the time of authoring with probability one.

3.2 Confidentiality

Confidentiality guarantees that a source-tracking scheme does not break the privacy of the underlying end-to-end encryption system or leak more forwarding metadata than the underlying messaging system to either the platform or to users. We formalize this by defining *user confidentiality*, which ensures that a malicious user cannot learn more about other users' messaging activity than it could without the source-tracking scheme, and *platform confidentiality*, which makes the same guarantee for the platform before a message is reported and additionally requires that, even after a report, the platform learns only the source and associated metadata of a reported message and nothing more.

We present definitions for both tree-linkable and tree-unlinkable user confidentiality. For the platform, we only define tree-linkable platform confidentiality because the platform can set the report metadata to be unique for each forwarding tree, allowing it to distinguish between two forwarding trees with the same plaintext when messages are reported. Note that tree-linkability/unlinkability for the platform is only relevant when a report is made: the security

guarantee for unreported messages is identical in the two settings. For simplicity, we present our definitions without allowing for adaptive corruption of users.

User Confidentiality. We define a security game for user confidentiality in Figure 1 with respect to a source tracking scheme ST , an adversary \mathcal{A} , and an underlying messaging scheme \mathcal{E} .

Intuitively, the user confidentiality game allows an adversary to assume the identity of a malicious user and simulate messaging activity with access to the view of that user. An adversary is given the ability to forward an arbitrary number of message pairs to adversary-controlled users. Only one message from each pair is actually forwarded, and the adversary wins if it can determine which of the two messages was sent.

To prevent trivial wins, the game requires that within a pair, messages have the same plaintext and have been sent to the adversary by the same honest user, although the forwarding history of the two messages, including their original source, can differ. The tree-linkable version of the game adds the additional restriction that the forwarding trees of the queried messages must either be the same tree, always be previously queried together, or have never before been sent to the malicious user. While the adversary cannot make calls to the challenge with two different message plaintexts, the game still captures basic unlinkability because an adversary can make two or more calls to the challenge, allowing it to receive multiple different messages along paths of its choice.

This game is constructed by giving the adversary access to a number of oracles it can use to simulate activity in a messaging system. The `getUser` function allows the adversary to create an arbitrary number of malicious (adversary-controlled) users, which are stored in the set \mathcal{U}_{mal} , as well as honest users that it does not control, which are stored in \mathcal{U}_{honest} . `goodAuth` and `goodFwd` allow the adversary to send messages between users it does not control, while `malSend` and `malRec` allow an adversary to send messages between an honest user and a malicious user. The adversary is also given oracles `send`(\cdot, U, \cdot) and `receive`(\cdot, \cdot, U) to send and receive messages from the adversary-controlled users $U \in \mathcal{U}_{mal}$ in the underlying encrypted messaging scheme. In all cases, the adversary gets to choose the plaintext m that is sent as well as the metadata md included with authored messages. The same game applies for both tree-linkable and tree-unlinkable security, with the only change being that the *linkable* flag in the `malRec` oracle is set to true for tree-linkable security.

In order to keep track of messages and their associated forwarding trees, each message sent in the game is identified by a unique message id mid , and each forwarding tree is identified by a unique tid , which is the value of mid for the message at the root of the tree. The game uses table T_{auth} to keep track of authoring data ad for each honest user and table T_{rec} to keep track of messages received by honest users. The set \mathcal{S}_{seen} records which forwarding trees the adversary has encountered, and the set $\mathcal{S}_{allowed}$ keeps track of the pairs of trees from which `malRec` can be given messages without enabling a trivial win in the tree-linkable variant of the game.

Definition 3.2. The *user confidentiality* advantage of an adversary \mathcal{A} against a source-tracking scheme ST and messaging scheme \mathcal{E} is defined as

$$\text{Adv}_{ST, \mathcal{E}}^{\text{uconf}}(\mathcal{A}) = \left| \Pr \left[\text{UCONF}_{ST, \mathcal{E}}^{\mathcal{A}, 1, \ell} = 1 \right] - \Pr \left[\text{UCONF}_{ST, \mathcal{E}}^{\mathcal{A}, 0, \ell} = 1 \right] \right|$$

where the adversary is given access to send and receive oracles for all malicious users in \mathcal{U}_{mal} as well as oracles `getUser`, `goodAuth`, `goodFwd`, `malSend`, and `malRec` $_{\ell, b}(\cdot, \cdot, \cdot, \cdot)$. We refer to the advantage as the *tree-linkable* advantage if $\ell = \text{true}$ and as the *tree-unlinkable* advantage otherwise.

Platform confidentiality. We assume that a platform for an encrypted messaging scheme can see that a user U_s sent a message to U_r , but that it cannot learn anything about the content of that message, or whether it was a forward. There has been research into the possibility of metadata-hiding messaging schemes where the sender or receiver of a message can be hidden from the platform (e.g., [8, 17–19, 36]), but we will be aligning our confidentiality goals with the former setting, which corresponds to most deployed messaging schemes, and assuming that the identities of the sending and receiving users are visible to the platform.

The platform confidentiality game operates similarly to the user confidentiality game, except that now the adversary controls the platform and an arbitrary number of malicious users. Similar to user confidentiality, the game gives the adversary oracles `goodSend` and `goodRec` for sending messages between honest users and `malSend` and `malRec` to send a message from or to a malicious user, respectively. We separate the sending and receiving of messages into separate oracles to allow the adversary to launch attacks that involve delivering message data to unintended recipients or otherwise tampering with the message delivery process. Additionally, a new report oracle allows a user to report a message to the platform. Finally, since the security game must ensure that the platform adversary learns nothing about the plaintext of messages, the challenge message in the platform confidentiality game is sent to an honest user instead of a malicious user. We formalize platform confidentiality with the *PCONF* game, described in Appendix A.1.

Definition 3.3. The advantage of an adversary \mathcal{A} in the *PCONF* game against a source-tracking scheme ST and messaging scheme \mathcal{E} is defined as

$$\text{Adv}_{ST, \mathcal{E}}^{\text{pconf}}(\mathcal{A}) = \left| \Pr \left[\text{PCONF}_{ST, \mathcal{E}}^{\mathcal{A}, 1} = 1 \right] - \Pr \left[\text{PCONF}_{ST, \mathcal{E}}^{\mathcal{A}, 0} = 1 \right] \right|$$

where the adversary is given access to the oracles in Figure 7.

Definition 3.4. We say that a source-tracking scheme ST and associated messaging scheme \mathcal{E} satisfy *confidentiality* if for all efficient adversaries \mathcal{A}_u and \mathcal{A}_p , we have that both $\text{Adv}_{ST, \mathcal{E}}^{\text{uconf}}(\mathcal{A})$ and $\text{Adv}_{ST, \mathcal{E}}^{\text{pconf}}(\mathcal{A}_p)$ are negligible. If the advantage used in the *UCONF* game is tree-linkable, then we say the scheme is *tree-linkably confidential*, and likewise *tree-unlinkably confidential* if the advantage is tree-unlinkable.

3.3 Accountability and Unforgeability

Accountability and unforgeability ensure that a source-tracking scheme can provide the platform with useful and accurate information when a message is reported. We present security games for both properties in Figure 2.

The *srcBIND* game addresses accountability by allowing a user to send and receive any messages it would like, and then challenging the adversary to send a message that is received by an honest user, but fails when the honest user tries to report it to the platform. The *unFORGE* game addresses unforgeability by challenging an

<pre> getUser($U, isMal$) $\mathcal{U} \leftarrow \mathcal{U}_{honest} \cup \mathcal{U}_{mal}$ if $U \in \mathcal{U}$: return \perp if $isMal$: $\mathcal{U}' \leftarrow \langle \mathcal{A}, P_{new}(sk, \mathcal{U}) \rangle(U, pk)$ if $\mathcal{U}' = \perp$: return \perp return $\mathcal{U}_{mal}.add(U)$ (ad, \mathcal{U}') $\leftarrow \langle U_{new}, P_{new}(sk, \mathcal{U}) \rangle(U, pk)$ if (ad, \mathcal{U}') = \perp : return \perp $T_{auth}[U] \leftarrow ad$ return $\mathcal{U}_{honest}.add(U)$ goodAuth(U_s, U_r, m, md) if $U_s \notin \mathcal{U}_{honest} \vee U_r \notin \mathcal{U}_{honest}$: return \perp $ad \leftarrow T_{auth}[U_s], msg \leftarrow (m, ad)$ (ad', pd, e) $\leftarrow \langle U_{auth}(msg), P_{send}(sk, md) \rangle(U_s, U_r, pk)$ (m, fd) $\leftarrow \langle U_{rec}, P_{rec}(sk, pd) \rangle(U_s, U_r, e, pk)$ if (m, fd) = \perp : return \perp $T_{auth}[U_s] \leftarrow ad'$ $mid \xleftarrow{\\$} \{0, 1\}^n$ $T_{rec}[mid] \leftarrow (U_s, U_r, m, fd, mid)$ return mid goodFwd(U_s, U_r, mid, md) if $U_s \notin \mathcal{U}_{honest} \vee U_r \notin \mathcal{U}_{honest} \vee mid \notin T_{rec}$: return \perp (U'_s, U'_r, m, fd, tid) $\leftarrow T_{rec}[mid]$ if $U'_r \neq U_s$: return \perp $msg \leftarrow (m, fd)$ (fd_s, pd, e) $\leftarrow \langle U_{fwd}(msg), P_{send}(sk, md) \rangle(U_s, U_r, pk)$ (m, fd_r) $\leftarrow \langle U_{rec}, P_{rec}(sk, pd) \rangle(U_s, U_r, e, pk)$ if $fd_r \vee fd_s = \perp$: return \perp $mid' \xleftarrow{\\$} \{0, 1\}^n$ $T_{rec}[mid] \leftarrow (U'_s, U'_r, m, fd_s, tid)$ $T_{rec}[mid'] \leftarrow (U_s, U_r, m, fd_r, tid)$ return mid' </pre>	<pre> $UCONF_{ST, \mathcal{E}}^{\mathcal{A}, b, \ell}$ (pk, sk) $\leftarrow KGen(params)$ $b' \leftarrow \mathcal{A}^{O_{\ell, b}}(pk)$ return b' malSend(U_s, U_r, md) if $U_s \notin \mathcal{U}_{mal}$: return \perp (pd, e) $\leftarrow \langle \mathcal{A}, P_{send}(sk, md) \rangle(U_s, U_r, pk)$ if $U_r \in \mathcal{U}_{mal}$: return $\langle \mathcal{A}, P_{rec}(sk, pd) \rangle(U_s, U_r, e, pk)$ (m, fd) $\leftarrow \langle U_{rec}, P_{rec}(sk, pd) \rangle(U_s, U_r, e, pk)$ if $fd = \perp$: return \perp $mid \xleftarrow{\\$} \{0, 1\}^n$ $\mathcal{S}_{seen}.add(mid), \mathcal{S}_{allowed}.add((mid, mid))$ $T_{rec}[mid] \leftarrow (U_s, U_r, m, fd, mid)$ return mid malRec_{linkable, b}(U_r, mid_0, mid_1, md) if $U_r \notin \mathcal{U}_{mal}$: return \perp if $mid_0 \vee mid_1 \notin T_{rec}$: return \perp ($U_s^{(0)}, U_r^{(0)}, m_0, fd_0, tid_0$) $\leftarrow T_{rec}[mid_0]$ ($U_s^{(1)}, U_r^{(1)}, m_1, fd_1, tid_1$) $\leftarrow T_{rec}[mid_1]$ if $U_r^{(0)} \neq U_r^{(1)} \vee U_r^{(0)} \in \mathcal{U}_{mal} \vee m_0 \neq m_1$: return \perp if $linkable \wedge (tid_0, tid_1) \notin \mathcal{S}_{allowed} \wedge \{tid_0, tid_1\} \cap \mathcal{S}_{seen} \neq \emptyset$: return \perp $\mathcal{S}_{allowed}.add((tid_0, tid_1)), \mathcal{S}_{seen}.add(tid_0, tid_1)$ $msg \leftarrow (m_b, fd_b)$ (fd_s, pd, e) $\leftarrow \langle U_{fwd}(msg), P_{send}(sk, md) \rangle(U_r^{(b)}, U_r, pk)$ $T_{rec}[mid_b] \leftarrow (U_s^{(b)}, U_r^{(b)}, m_b, fd_s, tid_b)$ return $\langle \mathcal{A}, P_{rec}(sk, pd) \rangle(U_r^{(b)}, U_r, e, pk)$ </pre>
---	---

Figure 1: User confidentiality game and oracles. The oracle $O_{\ell, b}$ gives the adversary access to $getUser(\cdot, \cdot)$, $goodAuth(\cdot, \cdot, \cdot, \cdot)$, $goodFwd(\cdot, \cdot, \cdot, \cdot)$, $malSend(\cdot, \cdot, \cdot)$, and $malRec_{\ell, b}(\cdot, \cdot, \cdot, \cdot)$ as well as the $send(\cdot, U, \cdot)$ and $receive(\cdot, \cdot, U)$ oracles for the underlying encrypted messaging scheme for all $U \in \mathcal{U}_{mal}$. ℓ is a fixed variable that determines whether the $UCONF$ game requires tree-linkability ($\ell = 1$) or tree-unlinkability ($\ell = 0$), and b is a hidden variable that determines which version of the game the adversary sees.

adversary who can see the result of receiving and sending messages to malicious users to create a report that is validated by the platform but implicates an honest user in sending a message it never sent.

Both games give the adversary a set of oracles which allow it to send and receive messages between arbitrary users and control the sending and receiving on a set of users that it controls. The oracle $getUser$ allows an adversary to add honest and malicious users. Oracles $goodAuth$ and $goodFwd$ allow sending messages from honest users, and $malSend$ allows malicious users to send messages to honest users. The adversary is also given access to a reporting oracle $\langle \cdot, P_{rep}(sk) \rangle(\cdot, pk)$ and oracles for sending messages between malicious users in the underlying messaging scheme.

We assume the oracle has an associated \mathcal{M}_{sent} set that keeps track of the messages that have been sent by honest users. As in the confidentiality games (Section 3.2), we use unique message ids (mid) to refer to message transactions between honest users, use a table T_{auth} to keep track of the authoring data for honest users, and use a table T_{rec} indexed by mid to keep track of the data associated with messages received by honest users.

Definition 3.5. We define the advantage of an adversary \mathcal{A} in the $srcBIND$ game for a source-tracking scheme ST and messaging scheme \mathcal{E} as

$$\text{Adv}_{ST, \mathcal{E}}^{\text{src-bind}}(\mathcal{A}) = \Pr[\text{srcBIND}_{ST, \mathcal{E}}^{\mathcal{A}} = 1].$$

$srcBIND_{ST,\mathcal{E}}^{\mathcal{A}}$	$goodAuth(m, md, U_s, U_r)$ if $U_s \notin \mathcal{U}_{honest}$: return \perp $ad \leftarrow T_{auth}[U_s]$ $msg \leftarrow (m, ad)$ $(ad', pd, e) \leftarrow \langle U_{auth}(msg), P_{send}(sk, md) \rangle(U_s, U_r, pk)$ $T_{auth}[U_s] \leftarrow ad', M_{sent}.add((U_s, md, m))$ if $U_r \in \mathcal{U}_{mal}$: return $\langle \mathcal{A}, P_{rec}(sk, pd) \rangle(U_s, U_r, e, pk)$ $(m, fd) \leftarrow \langle U_{rec}, P_{rec}(sk, pd) \rangle(U_s, U_r, e, pk)$ $mid \xleftarrow{\mathbb{R}} \{0, 1\}^n$ $T_{rec}[mid] \leftarrow (U_s, U_r, m, fd)$ return mid	$getUser(U, isMal)$ $\mathcal{U} \leftarrow \mathcal{U}_{honest} \cup \mathcal{U}_{mal}$ if $U \in \mathcal{U}$: return \perp if $isMal$: $\mathcal{U}' \leftarrow \langle \mathcal{A}, P_{new}(sk, \mathcal{U}) \rangle(U, pk)$ if $\mathcal{U}' = \perp$: return \perp return $\mathcal{U}_{mal}.add(U)$ $(ad, \mathcal{U}') \leftarrow \langle U_{new}, P_{new}(sk, \mathcal{U}) \rangle(U, pk)$ if $(ad, \mathcal{U}') = \perp$: return \perp $\mathcal{U}_{honest}.add(U)$ $T_{auth}[U] \leftarrow ad$
$unFORGE_{ST,\mathcal{E}}^{\mathcal{A}}$	$goodFwd(mid, md, U_r)$ if $mid \notin T_{rec}$: return \perp $(U'_s, U_s, m, fd) \leftarrow T_{rec}[mid]$ $msg \leftarrow (m, fd)$ $(fd', pd, e) \leftarrow \langle U_{fwd}(msg), P_{send}(sk, md) \rangle(U_s, U_r, pk)$ $T_{rec}[mid] \leftarrow (U'_s, U_s, m, fd')$ if $U_r \in \mathcal{U}_{mal}$: return $\langle \mathcal{A}, P_{rec}(sk, pd) \rangle(U_s, U_r, e, pk)$ $(m, fd'') \leftarrow \langle U_{rec}, P_{rec}(sk, pd) \rangle(U_s, U_r, e, pk)$ $mid' \xleftarrow{\mathbb{R}} \{0, 1\}^n$ $T_{rec}[mid'] \leftarrow (U_s, U_r, m, fd'')$ return mid'	$malSend(U_s, U_r, md)$ if $U_s \notin \mathcal{U}_{mal}$: return \perp $(pd, e) \leftarrow \langle \mathcal{A}, P_{send}(sk, md) \rangle(U_s, U_r, pk)$ if $U_r \in \mathcal{U}_{mal}$: return $\langle \mathcal{A}, P_{rec}(sk, pd) \rangle(U_s, U_r, e, pk)$ $(m, fd) \leftarrow \langle U_{rec}, P_{rec}(sk, pd) \rangle(U_s, U_r, e, pk)$ if $fd = \perp$: return \perp $mid \xleftarrow{\mathbb{R}} \{0, 1\}^n$ $T_{rec}[mid] \leftarrow (U_s, U_r, m, fd)$ return mid

Figure 2: Unforgeability and accountability games with their accompanying oracles. The adversary is also given access to $send(\cdot, U, \cdot)$ and $receive(\cdot, \cdot, U)$ oracles for the underlying encrypted messaging scheme for all adversary-controlled users $U \in \mathcal{U}_{mal}$, as well a reporting oracle $\langle \cdot, P_{rep}(sk) \rangle(\cdot, pk)$.

We say that ST is *accountable* or has accountability if $\text{Adv}_{ST,\mathcal{E}}^{\text{src-bind}}(\mathcal{A})$ is negligible for all efficient adversaries \mathcal{A} .

Definition 3.6. We define the advantage of an adversary \mathcal{A} in the $unFORGE$ game for a source-tracking scheme ST and messaging scheme \mathcal{E} as

$$\text{Adv}_{ST,\mathcal{E}}^{\text{unforge}}(\mathcal{A}) = \Pr[unFORGE_{ST,\mathcal{E}}^{\mathcal{A}} = 1].$$

We say that ST is *unforgeable* if $\text{Adv}_{ST,\mathcal{E}}^{\text{unforge}}(\mathcal{A})$ is negligible for all efficient adversaries \mathcal{A} .

We note that our unforgeability and accountability requirements do not apply to the platform, only to users. As observed by [31], requiring that a scheme be secure against platform-produced forgeries would conflict with the deniability properties discussed below.

3.4 Deniability

We require that a source-tracking scheme satisfies two types of deniability to protect users if the messaging system is compromised:

- *universal deniability*: Users should be able to deny participating in a forwarding path of a reported message.
- *platform compromise deniability*: If the platform's secret key sk is leaked publicly, a user should still be able to deny participating in a reported forwarding path.

Motivated by the approach of [31], we capture these goals by requiring the existence of efficient protocols $UForge$ and $PForge$ that can be executed by any user (with access to the platform's secret keys in the case of $PForge$). These protocols must successfully create a set of forged transcripts and forwarding data indistinguishable to a third party from the actual transcripts resulting from the path and subsequent report of a forwarded message. In universal deniability, the party tasked with differentiating transcripts has access to the platform's public keys and the secret keys of all users in the system, while for platform compromise deniability, they are given access to the platform's secret keys as well. We discuss the space of potential deniability definitions and formalize security games for universal deniability ($UnivDEN$) and platform-compromise deniability ($PlatDEN$) in Appendix A.2.

Definition 3.7. For a deniability game DEN and forgery algorithm $Forge$, the advantage of an adversary \mathcal{A} against this game, $\text{Adv}_{ST,\mathcal{E},Forge}^{\text{den}}(\mathcal{A})$, is defined as

$$\text{Adv}_{ST,\mathcal{E},Forge}^{\text{den}}(\mathcal{A}) = \left| \Pr \left[DEN_{ST,\mathcal{E},Forge}^{\mathcal{A},1} = 1 \right] - \Pr \left[DEN_{ST,\mathcal{E},Forge}^{\mathcal{A},0} = 1 \right] \right|.$$

Definition 3.8. We say that a source-tracking scheme ST is *deniable* if there exist efficient (possibly interactive) algorithms $UForge$

and PForge such that for any efficient adversaries \mathcal{A}_u and \mathcal{A}_p , $\text{Adv}_{ST, \mathcal{E}, \text{UForge}}^{\text{univden}}(\mathcal{A}_u)$ and $\text{Adv}_{ST, \mathcal{E}, \text{PForge}}^{\text{platen}}(\mathcal{A}_p)$ are both negligible.

4 TREE-LINKABLE SOURCE TRACKING

Straw-man scheme. A first straw-man attempt at tree-linkable source-tracking could have the identity U_S of the original sender of a message included with the message plaintext, so the final message becomes an encryption of (U_S, m) . While this approach would satisfy platform confidentiality and deniability requirements, it does not satisfy user confidentiality, accountability, or unforgeability, as nothing would prevent malicious users from observing or tampering with the sender identity included in the message.

Our scheme. Instead, the starting point for our approach is to have the platform append a signature on $(\text{Enc}(k, U_S), c_m)$ to each message, where k is an encryption key known only to the platform and c_m is a commitment to the message plaintext, provided by the sender. This signature is passed along whenever a message is forwarded, and the signature is checked by each message recipient using the platform's public key and commitment randomness r included by the sender alongside the end-to-end encrypted plaintext. This scheme enforces user confidentiality, unforgeability, and accountability, but it does not yet provide platform confidentiality because the platform can always decrypt the sender identity that is sent alongside each message.

Our final scheme combines the construction thus far with the straw-man construction. When a message is forwarded for the first time, the forwarder includes the signature, commitment, and commitment randomness alongside the message plaintext that is end-to-end encrypted and hidden from the platform, i.e., it encrypts and sends $(m, \text{Enc}(k, U_S), c_m, r)$. Subsequent forwards pass along the same information *inside* the end-to-end encrypted ciphertext, with each recipient checking that the signature and commitment match the provided message and encrypted sender pair. Fresh messages include padding to hide the fact that they do not contain forwarding information, and forwarded messages send the server a commitment to an empty message to hide the fact that they are forwards. Users report messages by sending the message plaintext, commitment, randomness, and signature to the platform, who checks the signature and decrypts the identity of the original sender.

In terms of security, hiding the forwarding information inside the end-to-end encrypted message provides platform confidentiality, while encrypting the original sender provides user confidentiality. The fact that the forwarding information a user sees does not change when a message is forwarded makes the scheme tree-linkable: a user can tell when the same original message has been forwarded to it twice. Accountability and unforgeability are provided by the platform being the one to include the sender identity and signing it, and deniability results from the sender identity being encrypted and therefore hidden from everyone except the platform (and being trivial to forge if the platform keys are known).

Our scheme, formalized in Figure 3, makes use of the underlying messaging scheme $\mathcal{E} = (\text{send}(m, U_S, U_r), \text{receive}(ct, U_S, U_r))$ as well as the following cryptographic tools.

- A symmetric encryption scheme

$$\mathcal{P} = (\text{KGen}_{sym} \rightarrow k, \text{Enc}(k, m) \rightarrow ct, \text{Dec}(k, ct) \rightarrow m).$$

KGen(params)	Report(m, pk)
$k \leftarrow \text{KGen}_{sym}$	$U_{rep}(fd) \quad P_{rep}(sk = (k, sk_s))$
$(sk_s, vk_s) \leftarrow \text{KGen}_{sig}$	fd
$pk \leftarrow vk_s, sk \leftarrow (k, sk_s)$	\longrightarrow
return (pk, sk)	$(\sigma, src, c_m, r) \leftarrow fd$
newUser(U_i, sk)	if $\neg \text{Open}(c_m, m, r) :$
if $U_i \in \mathcal{U} :$	return \perp
return \perp	if $\neg \forall f(pk, \sigma, (c_m, src)) :$
$\mathcal{U}.add(U_i)$	return \perp
return \mathcal{U}	return $\text{Dec}(k, src)$
AuthMsg(U_s, U_r, pk)	
$U_{auth}(msg)$	$P_{send}(sk = (k, sk_s), md)$
$(m, \perp) \leftarrow msg$	
$(c_m, r) \leftarrow \text{Commit}(m)$	
//send $(m, data)$ via underlying msg scheme	
$e \leftarrow \text{send}((m, \perp, c_m, r), U_s, U_r)$	
$(c_m, e) \longrightarrow$	
return e	$src \leftarrow \text{Enc}(k, (U_s, md))$
	$\sigma \leftarrow \text{Sig}(sk_s, (c_m, src))$
	return $(pd = (\sigma, src), e)$
FwdMsg(U_s, U_r, pk)	
$U_{fwd}(msg)$	$P_{send}(sk = (k, sk_s), md)$
$(m, fd) \leftarrow msg$	
$(c_m, r) \leftarrow \text{Commit}(\perp)$	
$e \leftarrow \text{send}((m, fd, c_m, r), U_s, U_r)$	
$(c_m, e) \longrightarrow$	
	$src \leftarrow \text{Enc}(k, (U_s, md))$
	$\sigma \leftarrow \text{Sig}(sk_s, (c_m, src))$
	return $(pd = (\sigma, src), e)$
RecMsg(U_s, U_r, e, pk)	
U_{rec}	$P_{rec}(sk, pd)$
	$pd \longleftarrow$
$(\sigma, src) \leftarrow pd$	
$(m, fd', c_m, r) \leftarrow \text{receive}(e, U_s, U_r)$	
//verify platform signature	
if $\neg \forall f(pk, \sigma, (c_m, src)) :$	return \perp
if $fd' = \perp :$ //If message is new	
if $\neg \text{Open}(c_m, m, r) :$	return \perp
return $m, fd \leftarrow (\sigma, src, c_m, r)$	
else :	$(\sigma', src', c'_m, r') \leftarrow fd'$
if $\neg \text{Open}(c_m, \perp, r) :$	return \perp
if $\neg \text{Open}(c'_m, m, r') :$	return \perp
if $\neg \forall f(pk, \sigma', (c'_m, src')) :$	return \perp //verify platform signature
return (m, fd')	//msg = (m, fd') will be used to forward the message

Figure 3: Tree-linkable construction.

- A signature scheme

$$\begin{aligned} S &= (\text{KGen}_{\text{sig}} \rightarrow (\text{vk}_s, \text{sk}_s), \text{Sig}(\text{sk}_s, m) \rightarrow \sigma, \\ &\quad \text{Vf}(\text{vk}_s, \sigma, m) \rightarrow \{0, 1\}). \end{aligned}$$

- A commitment scheme

$$C = (\text{Commit}(m) \rightarrow (c_m, r), \text{Open}(c_m, m, r) \rightarrow \{0, 1\}).$$

One of the merits of the scheme in Figure 3 is that it is non-interactive. The communication for each protocol only consists of a single message sent from the user to the platform in the case of a send, or from the platform to the user in the case of a receipt, though optional status messages could be added to let a platform or user know if an operation succeeds.

We prove the following security theorems in Appendix B. Definitions of CPA- and AE-security can be found in [2].

Theorem 4.1. *Assuming that the platform encryption scheme \mathcal{P} is CPA-secure, the commitment scheme C is hiding, and that the messaging encryption scheme \mathcal{E} is AE-secure, then Scheme 1 is tree-linkably confidential (Definition 3.4).*

Theorem 4.2. *Assuming that the platform's signature scheme S is secure against any efficient adversary \mathcal{A} , then Scheme 1 satisfies accountability (Definition 3.5).*

Theorem 4.3. *Assuming that the signature scheme S is unforgeable and the commitment scheme C is binding against any efficient adversary \mathcal{A} , then Scheme 1 is unforgeable (Definition 3.6).*

Theorem 4.4. *If for any efficient adversary, the messaging system's encryption scheme \mathcal{E} satisfies deniability¹ and the platform's encryption scheme \mathcal{P} is CPA-secure, then the scheme is deniable (Definition 3.8).*

5 TREE-UNLINKABLE SOURCE-TRACKING

Intuitively, we could achieve tree-unlinkability in the scheme described in Section 4 if there was a way to re-randomize the forwarding data fd each time a message is forwarded. However, such a re-randomization is not easily achieved because re-randomizing the platform's signature is not sufficient to make forwarding data unlinkable. Both the signature and the signed contents need to be re-randomized without revealing the contents to the platform.

Our tree-unlinkable construction preserves the basic framework of the tree-linkable scheme in that forwarding data consists of a "signature" (in this case a MAC) on a commitment to the message and an encryption of the source user and optional metadata. Users re-randomize this forwarding data by re-randomizing the contents of the signature and proving the validity of the re-randomized values in zero-knowledge to the platform, which then issues the user a fresh signature on those values.

Our scheme can be thought of as similar to a keyed-verification anonymous credential scheme [4–6] that we modify to allow for anonymous and unlinkable credential delegation via forwarding. The attributes included in each credential are an encryption of the source identity and metadata as well as a hash of the message content. Since our attributes are a mix of group elements and scalars and we need users to efficiently prove properties about these

attributes to the platform in zero-knowledge, we use the keyed-verification anonymous credential scheme presented in [5]. This construction relies on an algebraic MAC rather than a signature for the credential. This is fine for us because only the platform will need to distribute and verify forwarding credentials.

5.1 Tools

Zero Knowledge Proofs. We use zero knowledge proofs, denoted with the standard Camenisch-Stadler notation ([3]), i.e.,

$$\pi \leftarrow PK\{(\text{secrets}) : \text{expressions with secret and public vals}\}$$

When verifying a proof π , we use the notation $\text{Vf}(\pi, [\text{optional}: P_1, P_2, \dots])$, where P_1, P_2, \dots are any public values used in the proof.

Algebraic MAC. The main building block of our scheme is an algebraic MAC construction from Chase et al. [5]. Our particular instantiation requires a MAC on three attributes, two of which are group elements encoding the source information about a message (E_1 and E_2), and the other a scalar corresponding to a hash of the message plaintext m , which we refer to by $d \leftarrow H(m)$.

This MAC requires a group \mathbb{G} of prime order q where the discrete log problem is hard and 10 (public) fixed group elements. Because we use most of the MAC functionality as a black-box in our scheme, we will only deal with five of these parameters explicitly, and refer to them as $G, G_d, G_{y_1}, G_{y_2}, G_{y_3}$. These parameters correspond to the parameters of the same name described in Section 3.1 of [5].

In an effort to reduce redundancy and keep our scheme simple, we use this MAC as a black-box via the following functions:

- $\text{KGen}_{\text{MAC}}(\text{params}) \rightarrow \text{sk}_{\text{MAC}}$: Generates a MAC key.
- $\text{issue}(d, E_1, E_2, \text{sk}_{\text{MAC}}) \rightarrow \sigma, \pi_{\text{issue}}$: This function issues a MAC, σ , on the attributes d, E_1 , and E_2 , which are provided in the clear, as well as a proof that the MAC has been computed correctly, π_{issue} (See [5], Section 3.2). When verifying the unblinded MAC on attributes $\text{info} = (d, E_1, E_2)$, we use the verification function $\text{Vf}(\sigma, \pi_{\text{issue}}, \text{info})$.
- $\text{blindIssue}(ct_d, ct_{E_1}, ct_{E_2}, \text{pk}, \text{sk}_{\text{MAC}}) \rightarrow ct_\sigma, \pi_{\text{issue}}$: Given ElGamal ciphertexts of each attribute encrypted under pk , this function blindly issues an encrypted MAC ct_σ encrypted under the same key. π_{issue} is a proof that the MAC is well-formed. (See [5], Section 5.10).
- $\text{Vf}_{\text{issue}}(\pi_{\text{issue}}, ct_\sigma, ct_d, ct_{E_1}, ct_{E_2}, \text{sk}) \rightarrow \sigma \text{ or } \perp$: This is called by a user to verify that a blindly issued MAC is well-formed. If the proof is correct, the function returns the decrypted MAC. (See [5], Section 5.10).
- $\text{prepPresent}(\sigma, d, E_1, E_2) \rightarrow C_d, C_{E_1}, C_{E_2}, z, \pi_{\text{present}}$: When a user wants to prove to a platform that it has a valid MAC, it uses this function to commit to its attributes and MAC with opening z , and a proof π_{present} that it has a valid MAC for these commitments. The commitments have the form $C_{E_1} = G_{y_1}^z E_1, C_{E_2} = G_{y_2}^z E_2, C_d = G_{y_3}^z d$. We include these commitments separately because they are later used for a second proof of the MAC contents, but assume all necessary commitments are included in π_{present} . (See [5], Section 3.2).
- $\text{blindVf}(\pi_{\text{present}}, C_F, \text{sk}_{\text{MAC}}) \rightarrow \text{valid or } \perp$: Verifies that π_{present} is a valid proof that the user holds a MAC on the committed attributes C_F . (See [5], Section 3.2).

¹See Appendix A.2 for a formal definition of deniability for encrypted messaging schemes.

```

present( $d, E_1, E_2, \sigma, \text{pk} = Y$ )
( $C_d, C_{E_1}, C_{E_2}, z, \pi_{\text{present}}$ )  $\leftarrow$  prepPresent( $\sigma, d, E_1, E_2$ )
//Additional commitments for proving re-randomization
( $z', \text{rnd}$ )  $\xleftarrow{\$} \mathbb{Z}_q^2$ 
 $C'_d \leftarrow C_d G_{y_3}^{z'}$ 
 $C'_{E_1} \leftarrow G_{y_1}^{z'} G^{\text{rnd}} C_{E_1}, C'_{E_2} \leftarrow G_{y_2}^{z'} Y^{\text{rnd}} C_{E_2}$ 
 $\pi_{\text{rerand}} = \text{PK}\{(z', \text{rnd}) : C'_d / C_d = G_{y_3}^{z'}$ 
 $\wedge C'_{E_1} / C_{E_1} = G_{y_1}^{z'} G^{\text{rnd}} \wedge C'_{E_2} / C_{E_2} = G_{y_2}^{z'} Y^{\text{rnd}}\}$ 
 $C_F \leftarrow (C_d, C_{E_1}, C_{E_2})$ 
 $C'_F \leftarrow (C'_d, C'_{E_1}, C'_{E_2})$ 
 $o_f \leftarrow (z + z', d, (E_1 G^{\text{rnd}}, E_2 Y^{\text{rnd}}))$ 
return ( $C_f, C'_f, o_f, \pi_{\text{present}}, \pi_{\text{rerand}}$ )

```

Figure 4: Construction of a sender's proof in our tree-unlinkable scheme. This function constructs a proof that a user has a valid forwarding credential, and additionally that C'_F is a commitment to the same message and a re-randomization of the source of the forwarding credential in question.

Protocols using El-Gamal Encryption. Chase et al.'s algebraic MAC (Section 5.1) relies on a number of protocols based on manipulating El-Gamal ciphertexts, variants of which will also be used in our scheme. These protocols all rely on the homomorphic properties of El-Gamal encryption, which allow users to re-randomize ciphertexts without knowledge of the platform's secret key.

Similarly to the Chase et al. scheme, re-randomizing El-Gamal ciphertexts will enable blind issuance and unlinkable presentation of a MAC on hidden attributes in our scheme. We will use their techniques to prove that a user possesses valid credentials to forward a message. Whereas the original Chase et al. paper shows how to prove that a commitment commits to the same value as a given ciphertext, our protocol will augment theirs and show that a commitment commits to a *re-randomization* of a given ciphertext. We need the re-randomization to render the encryption of a message author's identity unlinkable to past forwards of the same message.

5.2 Our Construction

Our scheme requires a group \mathbb{G} of prime order q in which the discrete log problem is hard and the 10 fixed group elements required by the MAC described in Section 5.1.

We assume that each message m can be hashed to a representation in \mathbb{Z}_q , which we denote as $d \leftarrow H(m)$. We also assume that each (src, md) pair corresponding to a source user and metadata value have a reversible representation in \mathbb{G} . Our scheme can easily be extended if more than one element is needed to represent these values.

Keys. Our scheme assumes that a platform has an El-Gamal key pair $(y, Y = G^y)$ for $y \in \mathbb{Z}_q$. This is used for encrypting source ciphertexts. The platform also has a MAC key sk_{MAC} .

Authoring and forwarding data structure. Forwarding data fd held by a user consists of a tuple $(m, (E_1, E_2), \sigma)$, where m corresponds to the message plaintext, (E_1, E_2) is an ElGamal encryption of some (src, md) pair $S \in \mathbb{G}$, and σ is a MAC on attributes $d \leftarrow H(m), E_1, E_2$.

```

redeem( $C_A, src, o_A, C_F, o_F, type, \text{pk} = Y$ )
( $C_d, C_{E_1}, C_{E_2}$ )  $\leftarrow C_F, (E_1^{(A)}, E_2^{(A)}) \leftarrow src$ 
( $z_F, d_F, E_1^{(F)}, E_2^{(F)}$ )  $\leftarrow o_F, (z_A, d_A) \leftarrow o_A$ 
 $h, r_1, r_2, r_3, \text{rnd} \xleftarrow{\$} \mathbb{Z}_q, H \leftarrow G^h$ 
 $(A_1, A_2) \leftarrow (G^{r_1}, H^{r_1} G_d^{type})$ 
 $(B_1, B_2) \leftarrow (G^{r_2}, H^{r_2} E_1^{(type)} G^{\text{rnd}})$ 
 $(C_1, C_2) \leftarrow (G^{r_3}, H^{r_3} E_2^{(type)} Y^{\text{rnd}})$ 
 $P_1 \leftarrow H = G^h \wedge A_1 = G^{r_1} \wedge B_1 = G^{r_2} \wedge C_1 = G^{r_3}$ 
 $\wedge C_A = G_d^{d_A} G_{y_3}^{z_A} \wedge C_d = G_d^{d_F} G_{y_3}^{z_F}$ 
 $P_A \leftarrow A_2 = H^{r_1} G_d^{d_A} \wedge B_2 / E_1^{(A)} = H^{r_2} G^{\text{rnd}}$ 
 $\wedge C_2 / E_2^{(A)} = H^{r_3} Y^{\text{rnd}}$ 
 $P_F \leftarrow A_2 = H^{r_1} G_d^{d_F} \wedge B_2 / C_{E_1} = H^{r_2} G^{\text{rnd}} / G_{y_1}^{z_F}$ 
 $\wedge C_2 / C_{E_2} = H^{r_3} Y^{\text{rnd}} / G_{y_2}^{z_F}$ 
 $\pi \leftarrow \text{PK}\{(h, r_1, r_2, r_3, \text{rnd}, d_A, z_A, d_F, z_F) :$ 
 $P_1 \wedge (P_A \vee P_F)\}$ 
info  $\leftarrow (H, (A_1, A_2), (B_1, B_2), (C_1, C_2))$ 
return ( $\pi, \text{info}, h, E_1^{(type)} G^{\text{rnd}}, E_2^{(type)} Y^{\text{rnd}}$ )

```

Figure 5: Proof to receive a message in our tree-unlinkable scheme. This function constructs a proof that the provided ciphertexts to be used for blind issuance of forwarding credentials encrypt either 1) the same message and re-randomized source as are committed to in C_F , or 2) commit to the same message as committed to in C_A and a re-randomization of src .

Users also store an authoring credential ad , which has the same form as the forwarding data, but for an unused message \perp . The authoring data is given to the user by the platform when the new user is created. The authoring data is used by the user to prove credentials on a new message so that new messages are indistinguishable from forwarded messages to the platform.

Once a user gets forwarding data after receiving a message, the forwarding data stays constant for the rest of its sending activity for that message. The MAC is presented blindly during future forwards, and so the possession of a valid forwarding credential can be proved to a platform multiple times unlinkably.

Sending a message. To send a message, the sender provides authoring information C_A and forwarding information $(C_F, C'_F, \pi_p, \pi_r)$. Only one of these is filled with useful information, depending on the type of message that the user wants to send. If the message is an authored message, the authoring information C_A is for the actual message plaintext, and the forwarding information is for an unused message \perp (created from the user's authoring data ad). If the message is a forward, the message contents are swapped, and the relevant forwarding data fd is used to create the forwarding information.

The authoring information is a commitment, C_A , to the \mathbb{Z}_q representation of the message, while the forwarding information is constructed in the present sub-protocol (Figure 4) and consists of C_F , which is a commitment to the forwarded message's \mathbb{Z}_q representation and source ciphertext, as well as C'_F , which is a new commitment to the same message and a re-randomization of the source ciphertexts. The forwarding information also includes π_p , a

KGen(<i>params</i>) $(y, Y \leftarrow G^y) \xleftarrow{\$} \text{KGen}_{\mathcal{P}}(\text{params})$ $\text{sk}_{MAC} \leftarrow \text{KGen}_{MAC}(\text{params})$ return ($\text{pk} = Y, \text{sk} = (y, \text{sk}_{MAC})$) <hr/> newUser(<i>U_i</i>, <i>sk</i>, <i>pk</i> = <i>Y</i>) U_{new} $P_{new}(\text{sk}, \mathcal{U})$ if $U_i \in \mathcal{U}$: return \perp $r \xleftarrow{\$} \mathbb{Z}_q, (E_1, E_2) \leftarrow (G^r, U_i Y^r)$ $(\sigma, \pi_{issue}) \leftarrow \text{issue}(\perp, E_1, E_2, \text{sk}_{MAC})$ <hr/> $\text{info} \leftarrow (\perp, E_1, E_2)$ $\sigma, \pi_{issue}, (E_1, E_2)$ if $\forall (\sigma, \pi_{issue}, \text{info})$: $ad \leftarrow (H(\perp), \sigma, (E_1, E_2))$ return ad $\mathcal{U}.add(U_i)$ return \perp return \mathcal{U} <hr/> AuthMsg(<i>U_s</i>, <i>U_r</i>, <i>pk</i> = <i>Y</i>) $U_{auth}(\text{msg})$ $P_{send}(\text{sk}, md)$ $m, ad \leftarrow \text{msg}$ //create proof sender has valid authoring data $out \leftarrow \text{present}(ad, \text{pk})$ $(C_F, C'_F, o_F, \pi_p, \pi_r) \leftarrow out$ $z \xleftarrow{\$} \mathbb{Z}_q, d \leftarrow H(m)$ $o_A \leftarrow (z, d)$ $C_A \leftarrow G_{y_3}^z G_d^d$ //commit to new message //send via underlying msg scheme $e \leftarrow \text{send}((m, o_A, o_F), U_s, U_r)$ <hr/> $e, C_A, C_F, C'_F, \pi_p, \pi_r$ if $\neg \text{blindVf}(\pi_p, C_F, \text{sk}_{MAC})$ $\quad \quad \quad \vee \neg \text{Vf}(\pi_r, C_F, C'_F)$: return \perp $r \xleftarrow{\$} \mathbb{Z}_q, S \leftarrow (U_s, md)$ $src \leftarrow (G^r, SY^r)$ return $(pd = (C_A, src, C'_F), e)$ <hr/> FwdMsg(<i>U_s</i>, <i>U_r</i>, <i>pk</i> = <i>Y</i>) $U_{fwd}(\text{msg})$ $P_{send}(\text{sk}, md)$ $(m, fd) \leftarrow \text{msg}$ $out \leftarrow \text{present}(fd, \text{pk})$ $(C_F, C'_F, o_F, \pi_p, \pi_r) \leftarrow out$ $z \xleftarrow{\$} \mathbb{Z}_q, C_A \leftarrow G_{y_3}^z G_d^\perp$ $o_A \leftarrow (z, \perp)$ $e \leftarrow \text{send}((m, o_A, o_F), U_s, U_r)$ <hr/> $e, C_A, C_F, C'_F, \pi_p, \pi_r$ if $\neg \text{blindVf}(\pi_p, C_F, \text{sk}_{MAC})$ $\quad \quad \quad \vee \neg \text{Vf}(\pi_r, C_F, C'_F)$: return \perp $r \xleftarrow{\$} \mathbb{Z}_q, S \leftarrow (U_s, md)$ $src \leftarrow (G^r, SY^r)$ return $(pd = (C_A, src, C'_F), e)$		RecMsg(<i>U_s</i>, <i>U_r</i>, <i>e</i>, <i>pk</i> = <i>Y</i>) U_{rec} $P_{rec}(\text{sk} = (y, \text{sk}_{MAC}), pd)$ $\quad \quad \quad \xleftarrow{\quad} pd$ $(C_A, src, C_F) \leftarrow pd$ $(m, o_A, o_F) \leftarrow \text{receive}(e, U_s, U_r)$ $(z_A, d_A) \leftarrow o_A$ if $C_A \neq G_{y_3}^{z_A} G_{d_A}^{d_A}$: return \perp $(z_F, d_F, E_1, E_2) \leftarrow o_F$ $(C_d, C_{E_1}, C_{E_2}) \leftarrow C_F$ if $C_d \neq G_{y_3}^{z_F} G_{d_F}^{d_F} \vee C_{E_1} \neq G_{y_1}^{z_F} E_1$ $\quad \vee C_{E_2} \neq G_{y_2}^{z_F} E_2$: return \perp if $d = d_A \neq \perp \wedge d_F = \perp$: $type \leftarrow A$ //new message elseif $d_A = \perp \wedge d_F = d \neq \perp$: $type \leftarrow F$ //forwarded message else : return \perp //prove rerandomizations commit to o_A or o_F $o \leftarrow \text{redeem}(C_A, src, o_A, C_F, o_F, type, \text{pk})$ $(\pi, info, h, E'_1, E'_2) \leftarrow o$ <hr/> $\pi, info$ $(C_A, src, C_F) \leftarrow pd$ $\quad \quad \quad \xrightarrow{\quad}$ //verify proof that rerandomizations are valid if $\neg \text{Vf}(\pi, info, C_A, src, C_F)$: return \perp //blindly issue MAC on rerandomized values $(ct_\sigma, \pi) \leftarrow \text{blindIssue}(info, \text{sk}_{MAC})$ <hr/> $\quad \quad \quad \xleftarrow{\quad} ct_\sigma, \pi$ //verify proof that the new MAC is valid $\sigma \leftarrow \text{Vf}_{issue}(\pi, ct_\sigma, info, h)$ if $\sigma == \perp$: return \perp return $(m, fd = (H(m), (E'_1, E'_2), \sigma))$ //(m, fd) used to forward the message <hr/> Report(<i>m</i>, <i>pk</i> = <i>Y</i>) $U_{rep}(fd)$ $P_{rep}(\text{sk} = (y, \text{sk}_{MAC}))$ //prove knowledge of valid fd $out \leftarrow \text{present}(fd, \text{pk})$ <hr/> $\quad \quad \quad \xrightarrow{\quad} out$ $C_F, (C'_F, o_F, \pi_p, \pi_r) \leftarrow out$ if $\neg (\text{blindVf}(\pi_p, C_F, \text{sk}_{MAC}) \wedge \text{Vf}(\pi_r, C_F, C'_F))$: return \perp $(z, d', E_1, E_2) \leftarrow o_F$ $d \leftarrow H(m)$ if $C'_F \neq (G_{y_3}^z G_{d'}^d, G_{y_1}^z E_1, G_{y_2}^z E_2)$ $\quad \vee d \neq d' \vee d = \perp$: return \perp return E_2/E_1^y //decrypt to recover (src, md)
---	--	--

Figure 6: Protocols for Scheme 2.

proof that the user holds a valid MAC on the values stored in C_F , and π_r , which proves that C'_F commits to the same message and (re-randomized) source as C_F . We note that these two proofs can easily be combined into a single statement, but we've written them out separately to distinguish between scheme-specific statements and the standard MAC issuance proof.

The authoring commitment to a hash of a message $d \in \mathbb{Z}_q$ has the form of a standard Pedersen commitment [24] with bases G_{y_3} and G_d : $C_A = G_{y_3}^r G_d^d$ for a random $r \leftarrow \mathbb{Z}_q$. A forwarding commitment takes the form of a multi-attribute Pedersen commitment on d and two group elements E_1 and E_2 with bases $G_{y_3}, G_{y_1}, G_{y_2}, G_d$: $C_F = (G_{y_1}^r E_1, G_{y_2}^r E_2, G_{y_3}^r G_d^d)$. Hashing the plaintext to an element in \mathbb{Z}_q ensures that these schemes are binding without the need for an additional group element to fix the opening to a particular value of r . We note that this structure can easily be extended to commit to multiple additional attributes that may be required to store additional metadata.

The sending user encrypts openings to C_A and C'_F alongside the message plaintext and passes that information on to the receiver.

Processing a message. To process a sent message, the platform first checks that the proofs π_p and π_r are valid for the provided commitments. If this is the case, it creates a new encryption src of the sending user's identity and some associated metadata to be used if the message is new. It then passes along the platform data $pd = C_A, C'_F, src$ to the receiver.

Receiving a message. During the receipt of a message, the receiver presents encryptions of a message plaintext and src ciphertext and proves that these are either (re-randomizations of) the forwarding information committed to in C'_F or of the authored message committed to in C_A and the new src ciphertext. The platform then blindly issues the receiver a new MAC on those attributes. The helper function `redeem` (Figure 5), is used by the receiver to construct a proof to the platform that the values it would like to get a MAC on are valid.

Reporting a message. To report a message, a user presents the platform with the message plaintext, a re-randomization of the source ciphertext it would like to report, and a proof that it has a valid credential on those values. This proof is created in the same manner as when sending a message by calling the helper function `present` (Figure 4). After verifying the proof, the platform can decrypt the ciphertext to reveal the identity of the source.

Security. We prove the following theorems in Appendix C.

Theorem 5.1. *Assuming that the MAC of [5] satisfies blind issuance and anonymity as defined in [5], the platform's El Gamal encryption scheme is CPA-secure, the proof system employed is zero-knowledge, and \mathcal{E} is an AE-secure encryption scheme, then Scheme 2 satisfies tree-unlinkable confidentiality.*

Theorem 5.2. *Assuming the correctness of the MAC's blind presentation protocol [5] and the completeness of the zero-knowledge proof system, Scheme 2 satisfies accountability.*

Theorem 5.3. *Assuming that the MAC presentation and issuance protocols of [5] satisfy unforgeability, the proof of knowledge system satisfies a strong extractability property, and the discrete log problem is hard in \mathbb{G} , then Scheme 2 is unforgeable.*

Theorem 5.4. *Assuming that the El Gamal encryption scheme \mathcal{P} is CPA-secure and \mathcal{E} is deniable, Scheme 2 is deniable.*

6 IMPLEMENTATION AND EVALUATION

We implemented both our constructions in Rust using the double-ratchet crate, an implementation of Signal's Double Ratchet protocol, as a baseline implementation of a messaging scheme [38]. The implementation of our tree-linkable scheme relies on the ed25519-dalek crate, which provides an implementation of ed25519 signing and verification [20]. For the tree-unlinkable scheme, we used curve25519-dalek's implementation of curve25519 [21] to implement the algebraic MAC of Chase et al. [5], El-Gamal encryption and verification, and the necessary zero-knowledge proofs, which we implemented using the zkp crate [9]. The zkp crate does not currently provide functionality to prove disjunctions of statements, so to implement the receive proof (Figure 5), we modified the zkp prover and verifier to support the OR-protocol for Schnorr proofs [2]. Only our tree-linkable implementation includes report metadata, which we set to be a UNIX timestamp.

We evaluated our implementations on a variety of message lengths, ranging from 10 to 8000 bytes. Because we found that increased message lengths impacted runtime by <1% over this range, we only provide data for 1KB messages. Evaluation was performed using Criterion on an Intel i7-6700 processor @ 3.40GHz running Ubuntu Linux, and all estimates are computed from the average of at least 300 trials, where the number of iterations was determined by the number sufficient to achieve a <1% margin of error ($p = 0.05$). We conclude that our tree-linkable scheme does not significantly increase computation or communication costs in comparison to schemes deployed today that do not support source-tracking.

We compare our schemes to a baseline messaging scheme without source-tracking in Table 1. The times for server-side computation and reporting reflect the additional computations incurred by our scheme. We do not include a baseline time for a scheme with no tracking for these values, as there are no server-side operations to support source-tracking in the baseline.

While our tree-linkable scheme results in increased client-side costs, the overall client-side cost is still on the same order of magnitude as a scheme without source-tracking, and costs remain concretely very low. The costs on the server-side, where cost increases are much more sensitive, are minimal by comparison. The majority of the cost of sending messages in this scheme is the necessary cost of delivering a message in the underlying messaging scheme. Note that we separately report receiving times for fresh and forwarded messages. The different times do not open a timing side channel, as the components of each protocol that involve interaction with the server are identical. A constant-time implementation would incur the greater of the two costs for each message.

Our tree-unlinkable scheme has costs that are an order of magnitude higher than the linkable scheme, but concrete costs still remain on the order of milliseconds for each operation.

Table 2 Shows the additional communication we incur compared to a standard double ratchet end-to-end encrypted messaging scheme that does not use source-tracking. The increase in communication costs across our two schemes mirrors the increase in computation costs. Although the greatest communication increase

	No Tracking	Tree-Linkable	Tree-Unlinkable
Send	175 μ s	201 μ s	2.7 ms
Rec. (fresh)	116 μ s	181 μ s	6.4 ms
Rec. (forward)	116 μ s	237 μ s	6.3 ms
Server	N/A	20 μ s	4.1 ms
Report	N/A	57 μ s	2.6 ms

Table 1: Computation time for using our constructions to send 1KB messages compared to a standard double ratchet end-to-end encrypted message. Although the cost of sending a message increases with message size, we found that the increase was always less than 1% when going from 10B to 1KB messages.

	Tree-Linkable	Tree-Unlinkable
Send	256B	712B
Rec.	320B	1688B
Report	160B	648B

Table 2: Additional communication incurred to send a message using our constructions compared to sending the message via a standard double ratchet end-to-end encrypted message.

is in the process of receiving a message, the data that clients must keep in order to report a message is only 128 Bytes in the tree-linkable scheme and 136 Bytes in the tree-unlinkable scheme.

Comparison to prior work. The only prior work that considers finding the source of a forwarded message is the traceback scheme of Tyagi et al. [32]. Whereas their work considers the setting where the system wishes to reveal all the users who received a particular message, we explicitly aim to prevent leaking this information. Despite the difference in security properties, it is worthwhile to compare our tree-linkable scheme to the path-lookup of Tyagi et al. (their more efficient scheme) to understand the performance trade-offs between the two approaches.

Traceback incurs lower client overhead than our scheme, requiring less than 8 μ s overhead to send a message and less than 2 μ s to receive one, meaning their client-side sending and receiving costs almost match the performance of a messaging scheme with no tracking as shown in Table 1. Per-message communication costs for traceback are also lower. Sending and receiving a message require an additional 64 and 48 bytes of communication, respectively (compared to 256 and 320 bytes for tree-linkable tracing), while a report consists of only 16 bytes in addition to the message plaintext. A report in our tree-linkable scheme requires 160 additional bytes.

The lower client costs to support message traceback come at the cost of increased server side storage and report verification costs. First, finding the source of a reported message using traceback requires computation linear in the length of the forwarding chain. At approximately 90 μ s per forward, the report verification cost is comparable to ours (57 μ s on a slightly faster processor) for a fresh message but becomes worse if a message is forwarded several times. The time to find the source of a reported message in our schemes is constant, regardless of the number of forwards.

More importantly, the traceback scheme requires the server to store 36 Bytes per message sent. The server cannot know when a client has read or deleted a message, so it must conservatively store tracing information for as long as the message should be traceable. At the scale of a messaging system like WhatsApp that delivers

on the order of 100 billion messages a day [29], this requires an additional 3.6 Terabytes of storage per day. On the other hand, our source-tracking schemes require no persistent server-side storage.

ACKNOWLEDGMENTS

This work was funded by NSF, DARPA, a grant from ONR, and the Simons Foundation. Opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of DARPA.

REFERENCES

- [1] Joël Alwen, Sandro Coretti, and Yevgeniy Dodis. 2019. The double ratchet: Security notions, proofs, and modularization for the signal protocol. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, Springer, Cham, Darmstadt, 129–158.
- [2] Dan Boneh and Victor Shoup. 2020. A Graduate Course in Applied Cryptography.
- [3] Jan Camenisch and Markus Stadler. 1997. Proof systems for general statements about discrete logarithms. *Technical Report/ETH Zurich, Department of Computer Science* 260 (1997).
- [4] Melissa Chase, Sarah Meiklejohn, and Greg Zaverucha. 2014. Algebraic MACs and keyed-verification anonymous credentials. In *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security*. ACM, Scottsdale, 1205–1216.
- [5] Melissa Chase, Trevor Perrin, and Greg Zaverucha. 2020. The Signal Private Group System and Anonymous Credentials Supporting Efficient Verifiable Encryption. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security (Virtual Event, USA) (CCS '20)*. Association for Computing Machinery, New York, NY, USA, 1445–1459. <https://doi.org/10.1145/3372297.3417887>
- [6] David Chaum. 1985. Security without identification: Transaction systems to make big brother obsolete. *Commun. ACM* 28, 10 (1985), 1030–1044.
- [7] Katriel Cohn-Gordon, Cas Cremers, Benjamin Dowling, Luke Garratt, and Douglas Stebila. 2020. A formal security analysis of the signal messaging protocol. *Journal of Cryptology* 33, 4 (2020), 1914–1983.
- [8] Henry Corrigan-Gibbs, Dan Boneh, and David Mazières. 2015. Riposte: An Anonymous Messaging System Handling Millions of Users. In *2015 IEEE Symposium on Security and Privacy, SP 2015, San Jose, CA, USA, May 17-21, 2015*. IEEE Computer Society, San Jose, 321–338.
- [9] Henry de Valence. 2019. zkp: Experimental zero-knowledge proof compiler in Rust macros. <https://github.com/hdevalence/zkp>.
- [10] Mario Di Raimondo and Rosario Gennaro. 2009. New Approaches for Deniable Authentication. *Journal of Cryptology* 22, 4 (Oct. 2009), 572–615. <https://doi.org/10.1007/s00145-009-9044-3>
- [11] Mario Di Raimondo, Rosario Gennaro, and Hugo Krawczyk. 2005. Secure Off-the-Record Messaging. In *Proceedings of the 2005 ACM Workshop on Privacy in the Electronic Society*. Association for Computing Machinery, New York, NY, United States, 81–89.
- [12] Yevgeniy Dodis, Peul Grubbs, Thomas Ristenpart, and Joanne Woodage. 2018. Fast message franking: From invisible salamanders to encryption. In *Annual International Cryptology Conference*. Springer, Springer, Cham, Santa Barbara, 155–186.
- [13] Marc Fischlin and Sogol Mazaheri. 2015. Notions of Deniable Message Authentication. In *Proceedings of the 14th ACM Workshop on Privacy in the Electronic Society (WPES '15)*. Association for Computing Machinery, New York, NY, USA, 55–64. <https://doi.org/10.1145/2808138.2808143>
- [14] Paul Grubbs, Jiahui Lu, and Thomas Ristenpart. 2017. Message Franking via Committing Authenticated Encryption. *IACR Cryptol. ePrint Arch.* 2017 (2017), 664.
- [15] Seny Kamara, Mallory Knodel, Emma Llansó, Greg Nojeim, Lucy Qin, Dhanaraj Thakur, and Caitlin Vogus. 2021. Outside looking in: Approaches to content moderation in end-to-end encrypted systems. <https://cdt.org/insights/outside-looking-in-approaches-to-content-moderation-in-end-to-end-encrypted-systems/>
- [16] Anunay Kulshrestha and Jonathan Mayer. 2021. Identifying Harmful Media in End-to-End Encrypted Communication: Efficient Private Membership Computation. In *USENIX Security*. USENIX, Virtual Event.
- [17] Albert Kwon, Henry Corrigan-Gibbs, Srinivas Devadas, and Bryan Ford. 2017. Atom: Horizontally Scaling Strong Anonymity. In *Proceedings of the 26th Symposium on Operating Systems Principles, Shanghai, China, October 28-31, 2017*. Association for Computing Machinery, New York, NY, United States, 406–422.
- [18] Albert Kwon, David Lu, and Srinivas Devadas. 2020. XRD: Scalable Messaging System with Cryptographic Privacy. In *17th USENIX Symposium on Networked*

- Systems Design and Implementation, NSDI 2020, Santa Clara, CA, USA, February 25–27, 2020*. USENIX, Santa Clara, 759–776.
- [19] David Lazar, Yossi Gilad, and Nickolai Zeldovich. 2018. Karaoke: Distributed Private Messaging Immune to Passive Traffic Analysis. In *13th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2018, Carlsbad, CA, USA, October 8–10, 2018*. USENIX, Carlsbad, 711–725.
 - [20] Isis Lovercruft. 2020. ed25519-dalek: Fast and efficient ed25519 signing and verification in Rust. <https://github.com/dalek-cryptography/ed25519-dalek>.
 - [21] Isis Lovercruft and Henry de Valence. 2020. curve25519-dalek: A pure-Rust implementation of group operations on Ristretto and Curve25519. <https://github.com/dalek-cryptography/curve25519-dalek>.
 - [22] Moxie Marlinspike. 2013. Simplifying OTR Deniability.
 - [23] Moxie Marlinspike and Trevor (ed.) Perrin. 2016. Signal » Specifications » The X3DH Key Agreement Protocol. <https://signal.org/docs/specifications/x3dh/>. (Accessed on 09/03/2020).
 - [24] Torben Pryds Pedersen. 1992. Non-Interactive and Information-Theoretic Secure Verifiable Secret Sharing. In *Advances in Cryptology — CRYPTO '91 (Lecture Notes in Computer Science)*, Joan Feigenbaum (Ed.). Springer, Berlin, Heidelberg, 129–140. https://doi.org/10.1007/3-540-46766-1_9
 - [25] Trevor Perrin and Moxie Marlinspike. 2016. The Double Ratchet Algorithm. <https://signal.org/docs/specifications/doublerratchet/>.
 - [26] Katitza Rodriguez and Seth Schoen. 2020. FAQ: Why Brazil's Plan to Mandate Traceability in Private Messaging Apps Will Break User's Expectation of Privacy and Security. <https://www EFF.org/deeplinks/2020/08/faq-why-brazils-plan-mandate-traceability-private-messaging-apps-will-break-users>.
 - [27] Prasanto K Roy. 2019. Why India wants to track WhatsApp messages. <https://www.bbc.com/news/world-asia-india-50167569>.
 - [28] Manish Singh. 2020. India likely to force Facebook, WhatsApp to identify the originator of messages. <https://techcrunch.com/2020/01/21/india-likely-to-force-facebook-whatsapp-to-identify-the-originator-of-messages/>.
 - [29] Manish Singh. 2020. 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/>.
 - [30] Udbhav Tiwari and Jochai Ben-Avie. 2020. Mozilla's analysis: Brazil's fake news law harms privacy, security, and free expression. <https://blog.mozilla.org/netpolicy/2020/06/29/brazils-fake-news-law-harms-privacy-security-and-free-expression/>.
 - [31] Nirvan Tyagi, Paul Grubbs, Julia Len, Ian Miers, and Thomas Ristenpart. 2019. Asymmetric Message Franking: Content Moderation for Metadata-Private End-to-End Encryption. *IACR Cryptol. ePrint Arch.* 2019 (2019), 565.
 - [32] Nirvan Tyagi, Ian Miers, and Thomas Ristenpart. 2019. Traceback for End-to-End Encrypted Messaging. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security (London, United Kingdom) (CCS '19)*. Association for Computing Machinery, New York, NY, USA, 413–430. <https://doi.org/10.1145/3319535.3354243>
 - [33] Nik Unger, Sergej Dechand, Joseph Bonneau, Sascha Fahl, Perl Henning, Ian Goldberg, and Matthew Smith. 2015. SoK: Secure Messaging. In *2015 IEEE Symposium on Security and Privacy*. IEEE, San Jose, 232–249.
 - [34] Nik Unger and Ian Goldberg. 2015. Deniable Key Exchanges for Secure Messaging. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security (CCS '15)*. Association for Computing Machinery, New York, NY, USA, 1211–1223. <https://doi.org/10.1145/2810103.2813616>
 - [35] Nik Unger and Ian Goldberg. 2018. Improved Strongly Deniable Authenticated Key Exchanges for Secure Messaging. *Proceedings on Privacy Enhancing Technologies* 2018, 1 (Jan. 2018), 21–66. <https://doi.org/10.1515/popets-2018-0003>
 - [36] Jelle van den Hooff, David Lazar, Matei Zaharia, and Nickolai Zeldovich. 2015. Vuvuzela: scalable private messaging resistant to traffic analysis. In *Proceedings of the 25th Symposium on Operating Systems Principles, SOSP 2015, Monterey, CA, USA, October 4–7, 2015*. Association for Computing Machinery, New York, NY, United States, 137–152.
 - [37] Nihal Vatasdas, Rosario Gennaro, Bertrand Ithurburn, and Hugo Krawczyk. 2020. On the Cryptographic Deniability of the Signal Protocol. In *ACNS 2020*. Springer, Rome, Italy, 188–209.
 - [38] S.R. Verschoor. 2020. double-ratchet: Pure Rust implementation of the Double Ratchet algorithm. <https://github.com/sebastianv89/double-ratchet>.

A DEFERRED DEFINITIONS

A.1 Platform Confidentiality

The platform confidentiality game gives the adversary oracles `goodSend` and `goodRec` for sending messages between honest users and `malSend` and `malRec` to send a malicious from or to a malicious user, respectively. We separate the sending and receiving of messages into separate oracles to allow the adversary to launch

attacks that involve delivering message data to unintended recipients or otherwise tampering with the message delivery process. Additionally, the report oracle allows a user to report a message to the platform.

To further strengthen the adversary, we allow for a platform to create a set \mathcal{U}_{mal} of platform-controlled users that can collude with the platform using the `getUser` oracle passed with `malUser` set to true. We define the set of honest users \mathcal{U}_{honest} as the set of users not controlled by the adversary, each created by a call to `getUser` with `malUser` set to false. These honest users are created by running the `newUser` protocol with the adversary as the platform, and adding the new user to the set of users only if the protocol is successful. The adversary is given access to oracles `send(·, U, ·)` and `receive(·, ·, U)` for $U \in \mathcal{U}_{mal}$ to send and receive messages from malicious users in the underlying messaging scheme.

As in the user confidentiality game, messages are identified by a unique identifier `mid`, trees are identified by `tid`, the identifier of the message at their root. When set to \perp , the `tid` value acts as a flag for messages that can't be revealed to be the adversary without allowing it to trivially win the confidentiality game.

The game keeps track of messages sent and received in tables T_{send} and T_{rec} , respectively. The table T_{auth} additionally keeps track of authoring data (`ad`) received by honest users when they author a new message. The `goodSend` oracle represents the challenge portion of the game. It allows the adversary to choose two messages c_0 and c_1 , one of which will actually be sent between two honest users, and the adversary must guess which message was sent. The c_i s include an associated value $c_i[\text{type}]$ that determines whether they are new messages or forwarded messages. If the type is new, then c_i is simply a message plaintext to be sent. Otherwise, c_i is the `mid` value for a message to be forwarded.

A.2 Deniability

Here, we present an extended discussion of our deniability definitions discussed in Section 3.4, and present formalizations of deniability for source-tracking schemes.

Deniability, or the guarantee that only the messaging platform can prove that a sender sent a message, is a common goal of many secure messaging applications [22]. Because source-tracking is designed to be applied on top of these existing messaging applications, it is important that source-tracking schemes preserve the messaging application's deniability guarantees.

One of the key settings where deniability is important is protection for whistleblowers who use the messaging application to anonymously relay sensitive information. From this perspective, deniability for source-tracking schemes becomes particularly important because we want to guarantee that not only the source of the forwarded message can deny sending it, but also that intermediate users along the path who may have reported or forwarded the message to an authority can deny this action if their messages are compromised.

Prior works [10, 11, 13, 31, 33] have investigated the large space of potential deniability definitions and their associated tradeoffs with notions of authentication and unforgeability.

$\frac{PCONF_{ST, \mathcal{E}}^{\mathcal{A}, b}}{(\text{pk}, \text{sk}) \leftarrow \text{KGen}(\text{params})}$ $b' \leftarrow \mathcal{A}^O(\text{pk}, \text{sk})$ $\text{return } b'$ $\frac{\text{getUser}(U, \text{isMal})}{\text{if } U \in \mathcal{U}_{\text{honest}} \cup \mathcal{U}_{\text{mal}} :}$ $\quad \text{return } \perp$ $\text{if isMal} :$ $\quad \text{return } U_{\text{mal}}.\text{add}(U)$ $ad \leftarrow \langle U_{\text{new}}, \mathcal{A} \rangle(U, \text{pk})$ $\text{if } ad = \perp : \text{return } \perp$ $T_{\text{auth}}[U] \leftarrow ad$ $\text{return } U_{\text{honest}}.\text{add}(U)$ $\frac{\text{report}(mid)}{\text{if } mid \notin T_{\text{rec}} : \text{return } \perp}$ $(U_s, U_r, m, fd, tid) \leftarrow T_{\text{rec}}[mid]$ $// \text{check if report contents can be revealed}$ $\text{if } tid = \perp : \text{return } \perp$ $\text{return } \langle U_{\text{rep}}(fd), \mathcal{A} \rangle(m, \text{pk})$ $\frac{\text{goodRec}(mid, e)}{\text{if } mid \notin T_{\text{send}} : \text{return } \perp}$ $(U_s, U_r, m, e', tid) \leftarrow T_{\text{send}}[mid]$ $\text{if } U_s \vee U_r \notin \mathcal{U}_{\text{honest}} : \text{return } \perp$ $(m', fd) \leftarrow \langle U_{\text{rec}}, \mathcal{A} \rangle(U_s, U_r, e, \text{pk})$ $\text{if } fd = \perp : \text{return } \perp$ $mid' \leftarrow \{0, 1\}^n, \text{ if } e \neq e' : tid \leftarrow \perp$ $T_{\text{rec}}[mid'] \leftarrow (U_s, U_r, m, fd, tid)$ $\text{return } mid'$	$\text{goodSend}(U_s, U_r, c_0, c_1)$ $\text{if } U_s \vee U_r \notin \mathcal{U}_{\text{honest}} : \text{return } \perp$ $mid \leftarrow \{0, 1\}^n$ $\text{for } i \in \{0, 1\} :$ $\quad \text{if } c_i[\text{type}] = \text{new} :$ $\quad \quad m_i \leftarrow c_i$ $\quad \quad tid_i \leftarrow mid$ $\quad \text{else} :$ $\quad \quad mid_i \leftarrow c_i$ $\quad \quad \text{if } mid_i \notin T_{\text{rec}} : \text{return } \perp$ $\quad \quad (U_s^{(i)}, U_r^{(i)}, m_i, fd_i, tid_i) \leftarrow T_{\text{rec}}[mid_i]$ $\quad \quad \text{if } U_r^{(i)} \neq U_s : \text{return } \perp$ $\text{if } m_0 \neq m_1 \vee tid_0 \neq tid_1 \vee tid_0 = \perp \vee tid_1 = \perp :$ $\quad tid \leftarrow \perp$ $\text{if } c_b[\text{type}] = \text{new} :$ $\quad ad \leftarrow T_{\text{auth}}[U_s], msg \leftarrow (m_b, ad)$ $\quad (ad', e) \leftarrow \langle U_{\text{auth}}(msg), \mathcal{A} \rangle(U_s, U_r, \text{pk})$ $\quad T_{\text{auth}}[U_s] \leftarrow ad'$ $\quad \text{else} :$ $\quad \quad msg \leftarrow (m_b, fd_b)$ $\quad \quad (fd, e) \leftarrow \langle U_{\text{fwd}}(msg), \mathcal{A} \rangle(U_s, U_r, \text{pk})$ $\quad \quad T_{\text{rec}}[mid_b] \leftarrow (U_s^{(b)}, U_s, m_b, fd, tid_b)$ $T_{\text{send}}[mid] \leftarrow (U_s, U_r, e, m_b, tid)$ $\text{return } mid$	$\frac{\text{malSend}(U_s, U_r, e)}{\text{if } U_s \notin \mathcal{U}_{\text{mal}} \vee U_r \notin \mathcal{U}_{\text{honest}} : \text{return } \perp}$ $(m, fd) \leftarrow \langle U_{\text{rec}}, \mathcal{A} \rangle(U_s, U_r, e, \text{pk})$ $\text{if } fd = \perp : \text{return } \perp$ $mid \leftarrow \{0, 1\}^n$ $T_{\text{rec}}[mid] \leftarrow (U_s, U_r, m, fd, mid)$ $\text{return } mid$ $\frac{\text{malRec}(U_s, U_r, c)}{\text{if } U_s \notin \mathcal{U}_{\text{honest}} \vee U_r \notin \mathcal{U}_{\text{mal}} : \text{return } \perp}$ $// \text{case 1: } c \text{ corresponds to a new message}$ $\text{if } c[\text{type}] = \text{new} : \quad m \leftarrow c$ $\quad ad \leftarrow T_{\text{auth}}[U_s], msg \leftarrow (m, ad)$ $\quad (ad', e) \leftarrow \langle U_{\text{auth}}(msg), \mathcal{A} \rangle(U_s, U_r, \text{pk})$ $\quad T_{\text{auth}}[U_s] \leftarrow ad'$ $// \text{case 2: } c \text{ corresponds to a forward}$ $\text{else} : \quad mid \leftarrow c$ $\quad \text{if } mid \notin T_{\text{rec}} : \text{return } \perp$ $\quad (U'_s, U'_r, m, fd, tid) \leftarrow T_{\text{rec}}[mid]$ $\quad \text{if } U'_r \neq U_s \vee tid = \perp : \text{return } \perp$ $\quad msg \leftarrow (m, fd)$ $\quad fd', e \leftarrow \langle U_{\text{fwd}}(msg), \mathcal{A} \rangle(U_s, U_r, \text{pk})$ $T_{\text{rec}}[mid] \leftarrow (U'_s, U'_r, m, fd', tid)$
--	--	--

Figure 7: Platform confidentiality game and oracles. The adversary is also given access to $\text{send}(\cdot, U, \cdot)$ and $\text{receive}(\cdot, \cdot, U)$ oracles for the underlying encrypted messaging scheme for all adversary-controlled users $U \in \mathcal{U}_{\text{mal}}$.

While there are many different flavors of deniability definitions that could be extended to source-tracking, we chose to focus on definitions that were reasonable with respect to the current deniability guarantees offered by secure messaging systems. *Online deniability*, which guarantees a user can maintain deniability even when a third party interacts with it during the protocol to gain evidence of the user's participation, is a strong form of deniability that has been shown not to hold for the Signal messaging protocol [34, 35]. Instead, we focus on *offline deniability*, which guarantees a user can deny participating in a conversation in the event that a transcript of the exchange is provided to a third party. This form of deniability has been shown to hold for the Signal protocol under strong assumptions [33, 37]. We show that our source-tracking schemes satisfy the strongest form of offline deniability using the simulated transcript approach of [11], meaning that any user can forge a

transcript of a forwarding path and report that is indistinguishable from a real transcript to a third party who holds the long-term keys for all involved users, assuming that the underlying messaging scheme satisfies this guarantee as well. This is stated formally in the following definition:

Definition A.1 (Deniability of an Encrypted Messaging Scheme). We say an encrypted messaging scheme

$$\mathcal{E} = (\text{send}(m, U_s, U_r), \text{receive}(ct, U_s, U_r))$$

satisfies deniability if there exists an efficient simulation algorithm $\text{Sim}_{\mathcal{E}}(U_s, U_r, m)$ that produces a transcript of U_s sending the message m using \mathcal{E} to U_r without access to the private keys of U_s or U_r such that this simulated transcript is indistinguishable from a real transcript to any efficient distinguisher with access to U_s and U_r 's secret keys.

Substituting in a weaker assumption on the messaging scheme would just result in the same weakened guarantee for the overall system, so using this approach means that we can guarantee a source-tracking scheme preserves the offline deniability of any messaging scheme it is built on top of.

Following the approach of [31], we additionally show that in addition to *universal deniability*, where the forger has access to an arbitrary third party user and the discerning party has access to the long-term keys of all involved users, our schemes also satisfy *platform-compromise deniability*, which gives the forger and discerning party additional access to the platform's keys. This second form of deniability ensures that even in the event that the platform was compromised and its secret keys exposed, deniability can be preserved.

While not required in our definition for deniable source-tracking schemes, it is worth noting that our scheme constructions provide strong deniability protections for *reporters* of malicious messages, regardless of the deniability guarantees of the underlying messaging scheme. Neither scheme makes use of the messaging oracles or users' secret keys while making a report, and additionally anyone with knowledge of the forwarding data for a particular message can report it even if they themselves did not receive it, meaning that the pool of users who could have reported a message is not restricted to the users on the forwarding path of the message in question.

We present two security games to address each of these types of deniability: *UnioDEN* and *PlatDen* (Figure 8). Each game accepts the corresponding forgery algorithm that simulates a forwarding path's transcripts as a parameter. The games get access to a challenge oracle $\text{Chal}(\cdot, \cdot, \cdot, \text{type})$, where the value of *type* is either *u* (universal) or *p* (platform) depending on the type of game, as well as the ability to create new users of their choice.

The deniability challenge consists of the adversary presenting a query that consists of a message *m*, a path of users *p*, a forging user U_D , and a list of metadata *mds*. The challenge will either output the actual transcripts (Tr_r , Tr_f , and Tr_{rep}) and forwarding data resulting from sending the message along the path and then being reported by the last user, or a forged version constructed by inputting the query into the forgery algorithm, which is given access to the sending and receiving capabilities of U_D and a simulator for the underlying messaging scheme, $\text{Sim}_{\mathcal{E}}$. The functions of this oracle are presented in Figure 8 as interactive protocols where \mathcal{F} is the portion of the interaction that the forger can control.

B SECURITY PROOFS FOR TREE-LINKABLE SOURCE-TRACKING

B.1 Confidentiality – Proof of Theorem 4.1

PROOF. The proof of this statement follows immediately from combining the results of Lemmas B.1 and B.2, which prove user and platform confidentiality, respectively. \square

Lemma B.1. *Assuming that \mathcal{P} is CPA-secure, and the commitment and signatures schemes \mathcal{C} and \mathcal{S} are correct, the advantage of any efficient adversary against the *IUCONF* game for Scheme 1 is negligible.*

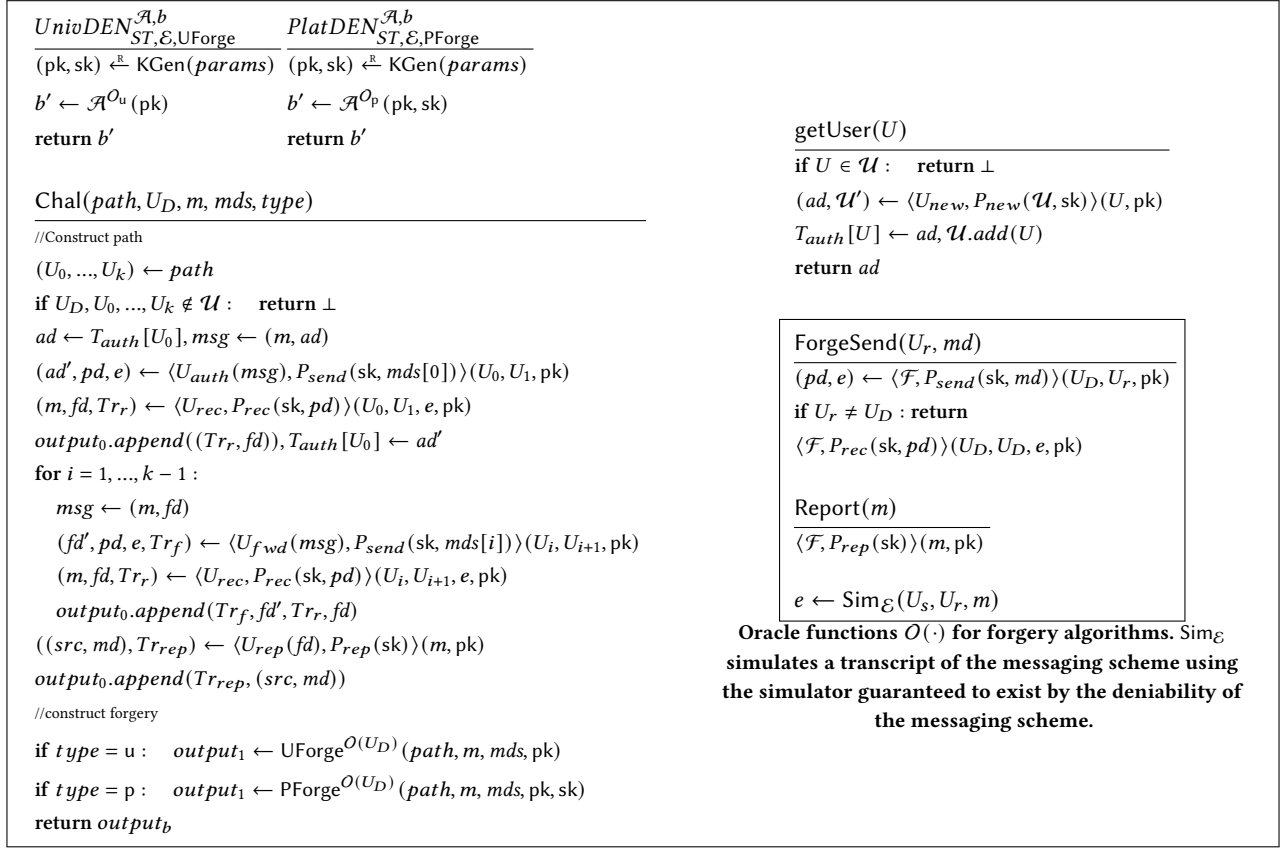
PROOF. We construct a series of hybrid games to show that $IUCONF^1$ is indistinguishable from $IUCONF^0$ to any efficient adversary. This proves that the scheme satisfies tree-linkable user confidentiality.

- G0: This game is identical to the standard *IUCONF* game when $b = 0$.
- G1: We modify *goodAuth* (oracle function for authoring a message between two honest users) and *malSend* (oracle function for sending a message from a malicious user) so that the *src* the platform computes and includes in the platform data is an encryption of some default user and metadata $U_D || md_D$ rather than the actual author of the message (the default data can, for example, be an all-zero string). Because the adversary does not have the secret key, and this value is never decrypted by the platform, this version of the game is indistinguishable from Game 0 by the CPA-security of the platform's encryption scheme, \mathcal{P} .
- G2: We modify *goodFwd* (oracle function for forwarding a message between two honest users) to immediately abort if the user provides an $mid \in T_{rec}$ such that the signature or commitment included in the forwarding data for the associated message is invalid. By the definition of the scheme, any message stored in T_{rec} that's associated with an honest receiver was received successfully by that honest user who checked the correctly formed signature and commitment, so this will never happen, and this game is indistinguishable from Game 1.
- G3: We again modify *goodFwd* so that after checking that the users associated with the given *mid* are correct and that $mid \in T_{rec}$, if m, fd, tid is the data associated with the provided *mid*, the oracle doesn't actually forward the message but instead just adds a new entry $T_{rec}[mid'] = (U_s, U_r, m, fd, tid)$ to the table. By definition of the scheme, the resulting entry in T_{rec} is identical to the entry if the message had actually been forwarded, and so this game is indistinguishable from Game 2 as the oracle does not modify any other persistent state outside of the table T_{rec} .

We note that after these modifications, the output of *malRec* (oracle function for forwarding a message to a malicious receiver, the function gets two potential forwards and uses the value of *b* to determine which one the adversary sees), which is the only oracle function that uses *b* to decide what to output, no longer depends on *b* because the forwarding data of all honestly sent messages have the same contents. This means that applying an identical series of hybrid steps to those described in G0 to G3 starting at $IUCONF^1$ instead of $IUCONF^0$ gets us to a game identical to Game 3 that is indistinguishable from $IUCONF^1$ to an efficient adversary.

Thus we have shown that we can go from $IUCONF^0$ to $IUCONF^1$ through a series of indistinguishable hybrids, and we can conclude that any adversary's advantage against the *IUCONF* game for Scheme 1 must be negligible. \square

Lemma B.2. *Assuming that the commitment scheme is hiding and the messaging encryption scheme \mathcal{E} is AE-secure, the advantage of any efficient adversary \mathcal{A} against the *PCONF* game for Scheme 1 is negligible.*



Oracle functions $\mathcal{O}(\cdot)$ for forgery algorithms. $\text{Sim}_{\mathcal{E}}$ simulates a transcript of the messaging scheme using the simulator guaranteed to exist by the deniability of the messaging scheme.

Figure 8: Deniability games and oracle functions. Oracles $\mathcal{O}_{\text{type}}$ for $\text{type} = u$ or p give the adversary access to the functions $\text{getUser}(\cdot, \cdot)$, $\text{Chal}(\cdot, \cdot, \cdot, \cdot, \text{type})$, and $\text{send}(\cdot, U, \cdot)$ and receive (\cdot, \cdot, U) oracles for the underlying messaging scheme for all users $U \in \mathcal{U}$. Oracle $\mathcal{O}(U_D)$ gives the forgery algorithms access to the Send and Report functions.

Note that the lemma requires that the messaging scheme guarantee authenticated encryption. This sort of guarantee is common among most encrypted messaging schemes, such as Signal's Double Ratchet Protocol [23].

PROOF. We first present a series of hybrid games that show an adversary cannot gain advantage by tampering with messages between honest users. By definition of the scheme, the platform is tasked with passing on the platform data (σ, src) to the receiving user, and the security game additionally gives the option for the platform to choose the message identifier e that the receiver gets as input. The platform could stray from the correct protocol by calling goodRec (oracle function for receiving a message between honest users) with an e that wasn't the same as the one outputted by the sending interaction associated with the mid , or by signing a value other than the src value chosen by the platform and the commitment presented by the sender.

- G0: This initial game is identical to $\text{PCONF}_{ST, \mathcal{E}}^{\mathcal{A}, 0}$, the game when $b = 0$.
- G1: We add the additional condition to goodRec to abort immediately if the inputted e value has never been outputted by the sender U_s as part of a message to the receiver U_r at any

point in the game. By definition, this can only happen during a call to goodSend because U_s and U_r are honest. This is indistinguishable from G0 because finding a valid e that was not sent from U_s to U_r contradicts the authenticated encryption properties of the messaging scheme, specifically its ciphertext integrity, and can therefore happen with only negligible probability.

We note that if the adversary presents a different e but the game does not abort due to the above condition, e must have resulted from a call to goodSend between U_s and U_r earlier in the game, and so there exists some mid'' in T_{send} that stores e . The only difference between calling goodRec with the correct pair (mid'', e) instead of (mid, e) is that the tid will always be \perp , which only weakens the adversary's power because it can no longer see the message contents by having the message reported and so we can assume that $e = e'$ in all other cases.

- G2: In this game, if during goodRec the platform presents an honest user with platform data (σ, src) and a message identifier e such that σ is a signature on the pair (c_m, src) , but (c_m, e) was never sent by the sending user, the game aborts. By the properties of Game 1, we know that the e must have

been created by the sending user and so contains an encryption of the commitment c_m that the sender linked to e , and therefore the receipt of the message will fail when the receiving user compares the commitment that is signed by the platform to the value encrypted in e . Therefore, this game is indistinguishable from Game 3 because an adversary cannot create an otherwise valid message that fails for this reason. Similarly, the src value is chosen by the platform and sent over in the clear, so the platform gains no advantage by signing a value other than the correct src , because it already knows it will cause the signature verification to fail.

We conclude that a platform gains no advantage by passing a receiver something other than the platform data and message identifier outputted by following the honest protocol. We now present two hybrids that show forwarded messages are indistinguishable from authored messages.

- G3: Game 3 is identical to Game 2, except that during `goodSend`, a forwarded message m with forwarding data σ, src, c_m, r that would normally be forwarded as e, c_\perp where c_\perp is a commitment to \perp is instead forwarded with e, c'_m where c'_m is a new commitment to m . The honest receiver evaluates the message as if c'_m was the original commitment to \perp (tracked by the oracle). By the hiding properties of the commitment, this is indistinguishable to an efficient adversary.
- G4: Game 4 is identical to Game 3, but during `goodSend`, e is encrypted as $(m_0, \perp, c'_{m_0}, r)$ rather than (m_0, fd, c_\perp, r) . This is indistinguishable to an efficient adversary because of the authenticated encryption properties of \mathcal{E} .

The messages and associated data of Game 4 are all identically distributed to authored messages. We now show that message contents are hidden from the adversary as well.

- G5: In this game, modify Game 4 so that when messages are sent via `goodSend`, the associated message identifier corresponds to an encryption of (m_1, \perp, c_{m_1}, r) and the commitment is still the original commitment to m_0 . The oracle keeps a table of the corresponding (m_0, \perp, c_{m_0}, r) that should have been sent and evaluates the results based on that value. By the authenticated encryption properties of \mathcal{E} , this is indistinguishable from Game 4 because all that changes is the decryption of e .
- G6: In this game, we switch the commitment sent in the clear, c_{m_0} , to instead be c_{m_1} , the commitment contained in e that commits to m_1 . The message is evaluated in the standard way using the decryption of e . This is indistinguishable from Game 5 by the hiding properties of the commitment.

We now note that applying the same game hops 1 - 4 starting from the game where $b = 1$ instead of 0 are each indistinguishable by the same arguments, and result in a game identically distributed to Game 6. Therefore we have shown that both $PCONF_{ST, \mathcal{E}}^1$ and $PCONF_{ST, \mathcal{E}}^0$ are indistinguishable from Game 7 to an efficient adversary, and so they are also indistinguishable from each other.

This means an efficient adversary can gain at most negligible advantage in the $PCONF$ game, and so Scheme 1 satisfies platform confidentiality. \square

B.2 Accountability – Proof of Theorem 4.2

PROOF. Suppose for purposes of contradiction we have an adversary that can win the $srcBIND$ game against Scheme 1 with non-negligible probability. This means that an honest user successfully received a message, but then the report of that message failed.

First, we consider the conditions that would cause a report to fail, and then show that each can happen with only negligible probability assuming that the receipt was successful:

- (1) The opening for the provided commitment was invalid.
This is checked by the receive protocol, and the protocol fails if the opening is invalid, so this will happen with probability zero on a message that was received successfully.
- (2) The signature verification on the source and commitment fails.
This exact check is also performed by the receive protocol, so again there is zero probability that this check will fail on a report of a validly received message.
- (3) The decryption of the source ciphertext fails.
By definition of the scheme, the platform only ever signs ciphertexts that it created, and can therefore be sure of their validity. Therefore, the probability that this occurs is upper bounded by the probability that an adversary could create a valid signature without the keys of the platform. By the unforgeability properties of the signature scheme, this happens with at most negligible probability.

We therefore conclude that the scheme satisfies accountability because if a message is received successfully, all three possibilities for a report failure happen with negligible probability. \square

B.3 Unforgeability – Proof of Theorem 4.3

PROOF. First, we remind ourselves of the conditions that must be met for an adversary to win the $unFORGE$ game. After sending some amount of messages between honest and malicious users, the adversary must present forwarding data (σ, src, c_m, r) and a message m to the platform such that:

- (1) c_m is a valid commitment with opening (m, r)
- (2) σ is a valid signature on src and c_m
- (3) src decrypts to an honest user U and metadata md such that U never sent m with metadata md during the adversary's initial queries.

We show that the probability an adversary can construct such an output is negligible via a series of hybrid games.

- G0: This is identical to the original $unFORGE$ game.
- G1: We modify G0 to abort immediately and return 0 if σ was not created by the platform oracle during the query portion of the game. By the unforgeability of the signature, the probability that the adversary can create a new valid signature without the platform's secret keys is negligible, and so this game is indistinguishable from G0 to efficient adversaries.
- G2: Suppose that the src value decrypts to U, md . We modify G1 to abort immediately if U did not send some message m' with commitment c_m and metadata md . By definition of the scheme, the platform would only create a valid signature linking U and md to the commitment c_m that U provided.

By G1, the signature must have been created by the platform, so this additional condition happens with zero probability, and we conclude G2 is indistinguishable from G1.

- G3: We modify G2 to immediately abort if m was not the plaintext of the message sent when U sent c_m and md . This can only happen if an adversary can find a second opening (m', r') for c_m that is different from the original opening created by U . By the binding property of the commitment scheme, this can happen with only negligible probability, and so G3 is indistinguishable from G2 to an efficient adversary.

In a series of three hops, we have reached a game where if the first two requirements described above are met and the U revealed is honest, then U must have sent m with metadata md at some point during the adversary's initial queries. This violates the third requirement, and so we conclude that there is no way to win Game 3. We have shown that *unFORGE* is indistinguishable from Game 3 to any efficient adversary, and so we can conclude that any efficient adversary can gain at most negligible advantage against the *unFORGE* game for Scheme 1, and therefore the scheme is unforgeable. \square

B.4 Deniability – Proof of Theorem 4.4

PROOF. The full proof that this scheme satisfies the deniability requirements of a secure source-tracking scheme appears in the full version of this paper. \square

C SECURITY PROOFS FOR TREE-UNLINKABLE SOURCE-TRACKING

C.1 Confidentiality – Proof of Theorem 5.1

PROOF. The proof follows immediately from the results of Lemmas C.1 and C.2. \square

C.1.1 User Confidentiality.

Lemma C.1. *Assuming the CPA-security of the platform's El Gamal scheme, any efficient adversary has negligible tree-unlinkable advantage against the *UNCONF* game.*

PROOF. We note that conditioned on the *src* ciphertext values (E_1, E_2) used in the forwarding data of messages sent via *malRec*, the adversary's view in *UNCONF*⁰ and *UNCONF*¹ is identical.

It therefore suffices to show that the difference in ciphertexts doesn't provide the adversary with a non-negligible advantage. We do this below by constructing a series of hybrid games to show that *UNCONF*⁰ and *UNCONF*¹ are indistinguishable to an efficient adversary.

- G0: This game is identical to the standard *UNCONF*⁰ game (when $b = 0$).
- G1: We modify *malRec* (oracle function that sends one of two messages from an honest sender to a malicious receiver) so that if the message would normally be sent with forwarding data that includes the source ciphertext (E_1, E_2) , which is an encryption of the source information $S_0 \in \mathbb{G}$, we instead replace (E_1, E_2) with a fresh encryption of S_0 , (E'_1, E'_2) , before sending.

This means that the only difference between G0 and G1 is that instead of seeing a re-randomization of the original source ciphertext upon receipt, the adversary sees a re-randomization of a fresh encryption of the same value. These two ciphertexts follow the same distribution, and are therefore indistinguishable to the adversary.

- G2: We further modify *malRec* so that instead of a fresh encryption of the original source S_0 , we make E'_1, E'_2 a fresh encryption of S_1 , the user/metadata pair that would have been used in the case that $b = 1$. Because the adversary does not know the platform's secret key, this is indistinguishable from Game 1 by the CPA security of El Gamal encryption.
- G3: This is the standard *UNCONF*¹ game. The only difference in the adversary's view between Games 2 and 3 is that the source stored in the forwarding data used to send the message that the adversary receives during a call to *malRec* is a fresh encryption of S_1 in the case of Game 2, and a re-randomized encryption of S_1 in the case of Game 3. By the same argument as the hop from Game 0 to Game 1, this difference is indistinguishable to an efficient adversary.

We have shown that in each hop, the two consecutive games are indistinguishable to an efficient adversary. We can therefore conclude that *UNCONF*⁰ is indistinguishable from *UNCONF*¹ to an efficient adversary, and therefore Scheme 2 satisfies tree-unlinkable user confidentiality. \square

C.1.2 Platform Confidentiality.

Lemma C.2. *Assuming that the messaging system's encryption scheme \mathcal{E} is a secure authenticated encryption scheme, the proof system used to prove attribute values on messages is zero-knowledge, and the algebraic MAC presented in [5] satisfies blind issuance and anonymous presentation properties, then the advantage of any efficient platform-confidentiality adversary for Scheme 2 is negligible.*

Note that the proof relies on the blind issuance and anonymity properties of the algebraic MAC presented in [5]. These are formally defined in [4] and discussed in [5]. Briefly, blind issuance guarantees that when the platform and user interact to issue a MAC on attributes known only to the user, the interaction is a secure two-party computation against malicious adversaries in which the user gets a valid MAC on their attributes and the platform learns nothing about the user's hidden attributes or the MAC that they receive. Anonymity guarantees that a user can prove to the platform they hold a valid MAC on particular attributes without revealing the MAC or hidden attributes to the platform.

PROOF. We show that the two *PCONF* games are indistinguishable to an efficient adversary via a series of game hybrids.

The first few hybrids show that the adversary can't gain a non-negligible advantage by not following the correct protocol. We begin by using the same approach as in the proof of Lemma B.2 to show that the adversary gains no advantage from calling *goodRec* (oracle function to make an honest user receive a message from another honest user) with a message identifier e that is not associated with the inputted *mid* or correct platform data *pd*.

- G0: This is identical to the standard *PCONF*⁰ game for $b = 0$.

G1: We add the additional condition to `goodRec` to abort immediately if the inputted e value has never been outputted by the sender U_s as part of a message to the receiver U_r at any point in the game. This is indistinguishable from G0 because finding a valid e that was not sent from U_s to U_r contradicts the authenticated encryption properties of the messaging scheme, and can therefore happen with only negligible probability.

We note that otherwise, e must have resulted from a call to `goodSend` (oracle function to make an honest user send a message to another honest user) between U_s and U_r earlier in the game, and so there exists some mid'' in T_{send} that stores e . The only difference between calling `goodRec` with the correct pair (mid'', e) instead of (mid, e) is that the tid will always be \perp , which only weakens the adversary's power because it can no longer see the message contents by having the message reported, so we can assume that $e = e'$ in all other cases.

G2: Consider the contents of the platform data (C_A, src, C'_F) that the platform passes to the receiver. In the correct protocol, the platform would pass over identical C_A and C'_F values to the ones that were presented by the sender (The src ciphertext is chosen by the platform and is allowed to be unique for every new message, so it doesn't matter what this is chosen to be). We modify G1 to immediately abort if C_A and C'_F are not identical to the commitments presented by the sender. This game is indistinguishable from Game 1, because the openings of C_A and C'_F are encrypted with the message plaintext in e , and so the receive function will always fail when the receiver checks the openings against the commitments that the platform provides.

We have ensured that the platform must act honestly and pass the correct information between protocols. We now show that it cannot act dishonestly *during* either of the interactive protocols `getUser` or `RecMsg`.

G3: We modify G2 to immediately abort if the platform gives the user an invalid MAC for the expected attributes in `goodRec`, `malSend`, or `getUser`.

By the blind issuance properties of the MAC (with no hidden attributes in the case of `NewUser`), the adversary can achieve this and still create a correct issuance proof with only negligible probability, and so this game is indistinguishable from Game 2.

We've now ensured that the platform interacts with the user according to the expected protocol.

G4: We modify Game 3 so that each e sent between honest users is just an encryption of a default value of all zeroes, and the oracle passes the actual information between the two users. Since the adversary has none of the secret keys for these interactions, this game is indistinguishable from Game 3 by the authenticated encryption properties of \mathcal{E} .

G5: We modify G4 so that during calls to `goodSend` and `goodRec` on an unrevealable message, i.e. one with tid set to \perp , we replace the receiving proof output by `redeem` and the re-randomization proof π_r output by `present` with the outputs

of the zero-knowledge simulator for the proofs. By the zero-knowledge properties of the proof system, this is indistinguishable from a real proof, and so the game is indistinguishable from G4.

G6: We alter each call to `goodSend` with $tid = \perp$ to use the authoring data for the currently sending user as the MAC that gets presented to the platform when the message is sent rather than the correct forwarding data for the message (if it is a forward). The commitments C_A and C'_F still commit to the same values that would have been used in the standard game. By the anonymity of the MAC, the platform cannot distinguish between a presentation of the authoring data and the forwarding data, because both are valid MACs. Therefore this game is indistinguishable from G5.

G7: We modify the ciphertexts $(A_1, A_2), (B_1, B_2), (C_1, C_2)$ to just be encryptions of the default values \perp, G^\perp, G^\perp when `goodRec` is called on a message with $tid = \perp$. By the blind-issuance property of the keyed-verification anonymous credentials scheme [5], this is indistinguishable to the platform from issuing a MAC on the expected values.

G8: Finally, we change C_A and C'_F in calls to `goodSend` with $tid = \perp$ so that they commit to the default message \perp and default source G^\perp, G^\perp . By the hiding properties of the commitment, this is indistinguishable from commitments for the actual message and source values.

We now note that applying the same hybrids described in Games 1-8 starting from the *PCONF* game where $b = 1$ instead of 0 are each indistinguishable by the same arguments, and result in a game identically distributed to Game 8, because all message interactions with $tid \neq \perp$ are identically distributed in both games by definition, and we have altered the $tid = \perp$ case to use the same authoring data, commit to the same attributes, and receive a MAC on the same values. Therefore we have shown that both $PCONF_{ST, \mathcal{E}}^1$ and $PCONF_{ST, \mathcal{E}}^0$ are indistinguishable from Game 8 to an efficient adversary, and so they are also indistinguishable from each other.

This means an efficient adversary can gain at most negligible advantage in the *PCONF* game, and so Scheme 2 satisfies platform confidentiality. \square

C.2 Accountability – Proof of Theorem 5.2

Intuitively, because the receiver of a message has complete control over what commitments and attributes it receives a MAC on, our scheme satisfies perfect accountability.

PROOF. As a reminder, the *srcBIND* game challenges an adversary with the ability to send messages with content of their choice between any two users and completely control some set of malicious users to create a message and associated content that is successfully received by an honest user, but then cannot be reported.

The tree-unlinkable scheme's definition of `Report` provides five opportunities for failure. We list these in order of occurrence:

- (1) The platform's call to `blindVf` fails to verify the validity of π_p , the proof that the reporter has a valid MAC for the message that they are reporting.

- (2) The platform's call to V_f fails to verify the validity of π_r , the proof that the commitments provided by the reporter contain correct re-randomizations of the original credentials.
- (3) The first equality check fails ($C'_f \neq (G_{y_3}^z G_d^d, G_{y_1}^z E_1, G_{y_2}^z E_2)$). In other words, the opening o_f is not a valid opening for the commitment C'_f .
- (4) The second equality check fails ($d \neq d'$).
- (5) The third and final equality check fails ($d = \perp$).

We'll show that assuming the platform and reporter honestly received the message being reported, each of these failures can happen with zero probability.

π_p fails to verify. Because the platform and user honestly followed the receipt protocol, a successful receipt guarantees that the reporter has a valid MAC on the attributes being reported. If a user was unable to create a valid presentation proof for this MAC, it would contradict the correctness of the MAC presentation of [5]. We conclude that this failure happens with zero probability.

π_r fails to verify. All of the relationships proved in π_r are constructed by the user during the present protocol. Because the user is honest and follows the protocol, they ensure that all of these relationships are valid, and so an incorrect proof would violate the completeness of the proof scheme, so this failure also happens with zero probability.

The opening to the commitment C'_f is invalid. By definition of the scheme, the reporting user must have received a valid opening (z_F, d, E_1, E_2) for C_F , meaning that

$$C_F = (C_d, C_{E_1}, C_{E_2}) = (G_{y_3}^{z_F} G_d^d, G_{y_1}^{z_F} E_1, G_{y_2}^{z_F} E_2).$$

When the honest reporter follows the present protocol, C'_F and o_F are computed as

$$\begin{aligned} C'_F &= (C_d G_{y_3}^{z'}, C_{E_1} G_{y_1}^{z'} G^{rnd}, C_{E_2} G_{y_2}^{z'} Y^{rnd}) \\ &= (G_{y_3}^{z_F} G_d^d G_{y_3}^{z'}, G_{y_1}^{z_F} E_1 G_{y_1}^{z'} G^{rnd}, G_{y_2}^{z_F} E_2 G_{y_2}^{z'} Y^{rnd}) \\ &= (G_{y_3}^{z_F+z'} G_d^d, G_{y_1}^{z_F+z'} E_1 G^{rnd}, G_{y_2}^{z_F+z'} E_2 Y^{rnd}) \end{aligned}$$

and $o_F = (z_F + z', d, (E_1 G^{rnd}, E_2 Y^{rnd}))$, so o_F is guaranteed to be the correct opening for C'_F , so this check fails with zero probability because the reporter honestly follows the protocol and successfully received the original commitment.

$d \neq d'$. This condition only occurs if the reporter presents a credential for a message other than the one being reported. This will never happen because the reporter is honest.

$d = \perp$. During a message receipt, the user checks that the message being received is not equal to \perp , and aborts otherwise, so this will never happen because the message being reported was successfully received.

In conclusion, we've shown that assuming a message was received successfully by an honest user, a report of that message can fail with zero probability, so the scheme satisfies perfect accountability. \square

C.3 Unforgeability – Proof of Theorem 5.2

PROOF. The unforgeability of our scheme will depend on the extractability of the proofs used. Note that the blind issuance properties and unforgeability of the anonymous credentials used in [5] assume that the system used for proofs of knowledge satisfies a strong extractability property, in particular, that we can extract the opening (z, d, E_1, E_2) of the commitments C_d, C_{E_1}, C_{E_2} used in the MAC presentation proof $\pi_{present}$, which can be used to re-compute the attributes of the MAC being presented. Since we use their scheme, whose security relies on the random oracle and generic group models, we inherit their use of these models in our scheme. Other possibilities for instantiating such a proof system are discussed in Appendix D of [4].

The proof is a bit long, so we will first lay out the high-level intuition. Suppose an adversary wins the game and presents a report to the platform on an (m, U, md) tuple that was never actually sent by the honest user U .

We have two possibilities. On one hand, the adversary could have achieved this by acting dishonestly during the last report interaction in order to report a message that it didn't actually have forwarding credentials for. On the other hand, the adversary may have correctly followed the standard report protocol, but was able to succeed because it had valid credentials that were created dishonestly earlier in the game.

In this second case, we can then look at the point in time when the adversary received these credentials earlier in the game. Once again, these could have been obtained by the adversary acting dishonestly to create credentials that it didn't actually have, or this receipt could have been performed honestly because the necessary credentials were created earlier in the game.

In this way, we can trace back through the adversary's interactions with the oracle to identify the first interaction that resulted in a set of forwarding credentials with an honest source user that didn't originate from a call to `goodAuth`.

We consider all the places where this could have happened, and then in the actual proof, remove these opportunities in a series of game hops:

- (1) *While the adversary was sending a message.*

If the incorrect credentials are created when an adversary is sending a message (i.e., during a call to `malSend`), this could happen because either the adversary is able to prove ownership of a MAC that was never created by the platform (addressed in G1 in the proof below), fake a proof that the new commitments stored in C'_F are correct re-randomizations of the original attributes (G3), or find an alternate opening of the commitments it used to prove the validity of its forwarding credentials, C_F and C'_F (G4).

- (2) *While the adversary was receiving a message.*

If the incorrect credentials are created when an adversary receives a message, this means that either the adversary found alternate openings to the commitments of the attributes it can receive credentials from, C_F and C_A (Games G4 and G5, respectively), it faked a proof that the attribute ciphertexts were a valid re-randomization of C_A or C_F (G2), or it received a valid MAC on different values by acting dishonestly during the blind issuance protocol for the MAC (G1).

(3) *While the adversary was reporting the message.*

A valid report follows the same approach as a send, so it has the same opportunities for dishonest action, which we cover above.

(4) *While the adversary was creating a new adversary-controlled user.*

The last potential place where incorrect credentials could have been created is during a call to `getUser`. However, the `newUser` function is non-interactive, so the adversary cannot act dishonestly in order to alter the credentials they receive.

In each game, we remove one of these opportunities for dishonest action and show that the resulting game is indistinguishable from the previous game to an efficient adversary. This means that in the final game, the adversary has no opportunity to create fake credentials, and so cannot win the game.

We present the games below.

G0: This game is identical to the standard *unFORGE* game.

G1: This game is identical to the *unFORGE* game, but we add a set \mathcal{M}_{MAC} , that keeps track of all attributes that the platform has issued a MAC on. To do this, we modify the receiving interaction so that when a proof π provided by the receiver is verified by the platform, we run the knowledge extractor to extract witnesses $(h, r_1, r_2, r_3, rnd, d_A, z_A, d_F, z_F)$ and abort if extraction fails. We use these witnesses to decrypt (A_1, A_2) , (B_1, B_2) , (C_1, C_2) to get attributes M, E_1, E_2 , i.e. $A_2/A_1^h, B_2/B_1^h, C_2/C_1^h$, respectively. We then add (M, E_1, E_2) to \mathcal{M}_{MAC} . The extractor fails with negligible probability, so this game is indistinguishable from Game 0.

G2: We further modify so that every time a message is sent or reported, if the proof output by present verifies, we again run the proof of knowledge extractor on the MAC presentation proof $\pi_{present}$ to extract attributes M, E_1 , and E_2 and abort if extraction fails. Once again, failure happens with negligible probability and this is indistinguishable from G1.

G3: After extracting the attributes M, E_1, E_2 as described in G2, we immediately abort if $(M, E_1, E_2) \notin \mathcal{M}_{MAC}$. If the game aborts due to the above reason, it means that the user proved ownership of a MAC that the platform never created. (Note that we're not requiring that the attributes be in any way correct at this stage, we are just looking directly at the exact attributes of the MACs that the platform created and comparing them to the attributes that the adversary tries to prove it has. This can happen with negligible probability due to the unforgeability of the anonymous credentials of [5], and therefore G3 is indistinguishable from G2.

G4: In this next game, every time an individual presents a valid receiving proof to the platform, we run the extractor to get witnesses $h, r_1, r_2, r_3, rnd, d_A, z_A, d_F, z_F$. We immediately abort if the extractor fails. By the extractability of the proof system, a valid proof fails extraction with negligible probability and so this game is indistinguishable to the forgery adversary. On success, we are guaranteed that the properties included in the proof are guaranteed to hold, i.e. that $P_1 \wedge P_A$ or $P_1 \wedge P_F$ as described in Figure 5 is satisfied.

We also note that by this check, if $P_1 \wedge P_A$ holds, then we are guaranteed that:

$$\begin{aligned} A_1, A_2 &= G^{r_1}, H^{r_1} G_d^{d_A} \\ B_1, B_2 &= G^{r_2}, H^{r_2} G^{rnd} E_1^{(A)} \\ C_1, C_2 &= G^{r_3}, H^{r_3} Y^{rnd} E_2^{(A)} \\ C_A &= G_d^{d_A} G_{y_3}^{z_A} \\ C_d &= G_d^{d_F} G_{y_3}^{z_F} \end{aligned}$$

So, the encryptions are valid encryptions of d_A and a re-randomization of (E_1, E_2) under public key H .

Similarly, if $P_1 \wedge P_F$ holds, then we are guaranteed that

$$\begin{aligned} A_1, A_2 &= G^{r_1}, H^{r_1} G_d^{d_F} \\ B_1, B_2 &= G^{r_2}, H^{r_2} G^{rnd} C_{E_1} / G_{y_1}^{z_F} \\ &= H^{r_2} G^{rnd} G_{y_1}^{z_F} E_1 / G_{y_1}^{z_F} \\ &= H^{r_2} G^{rnd} E_1 \\ C_1, C_2 &= G^{r_3}, H^{r_3} Y^{rnd} C_{E_2} / G_{y_2}^{z_F} \\ &= H^{r_3} Y^{rnd} G_{y_2}^{z_F} E_2 / G_{y_2}^{z_F} \\ &= H^{r_3} Y^{rnd} E_2 \\ C_d &= G_d^{d_F} G_{y_3}^{z_F} \\ C_A &= G_d^{d_A} G_{y_3}^{z_A} \end{aligned}$$

So, this guarantees that in the forwarded case, the encrypted values must be an encryption of d_F and a re-randomization of (E_1, E_2) , where (z_F, d_F, E_1, E_2) is a valid opening for the commitment C_F .

G5: In this game, we check the sending/reporting proofs in the interaction. When given valid $\pi_{present}$ and π_{rerand} proofs associated with commitments $C_F = (C_d, C_{E_1}, C_{E_2})$ and $C'_F = (C'_d, C'_{E_1}, C'_{E_2})$, we apply the proof extractor to extract the witnesses (z', rnd) from π_{rerand} , aborting on failure. By the strong extractability property, this extraction fails with only negligible probability and otherwise we are guaranteed that the statement in π_{rerand} holds true for the provided values of C_F and C'_F .

Moreover, this guarantees that if (z, d, E_1, E_2) is the opening of C_F extracted from the associated presentation proof $\pi_{present}$, then we must have

$$\begin{aligned} C'_d &= C_d G_{y_3}^{z'} = G_{y_3}^{z+z'} G_d^d \\ C'_{E_1} &= G_{y_1}^{z'} G^{rnd} C_{E_1} = G_{y_1}^{z+z'} G^{rnd} E_1 \\ C'_{E_2} &= G_{y_2}^{z'} Y^{rnd} C_{E_2} = G_{y_2}^{z+z'} Y^{rnd} E_2 \end{aligned}$$

and so $(z + z', d, E_1 G^{rnd}, E_2 Y^{rnd})$ must be a valid opening for C'_F .

By the requirements of Game 3, the sending user must have had a valid MAC on (d, E_1, E_2) , and so this ensures that the opening C'_F must be a re-randomization of attributes that the user has a valid MAC for.

G6: We have shown that the attributes used to prove and create MACs must be internally consistent (within a send/receive operation), but we now show that they need to consistent

across transactions too. Let G4 be identical to G3 except for the modification that every time we have a valid pair of sending proofs π_p and π_r , we extract the opening $(z + z', d, E_1 G^{rnd}, E_2 Y^{rnd})$ of C_F from the sending proof. In the same manner as G3. Then, during the corresponding receiving proof that uses the same C_F and some arbitrary C_A , we once again extract the witnesses $(h, r_1, r_2, r_3, rnd, d_A, z_A, d_F, z_F)$, aborting on failure. By the extractability of the proof system, extraction fails with negligible probability, making this game indistinguishable from G5.

- G7: We modify G6 so that after extracting the above witnesses, we use them to compute the new opening $(z_F, d_F, E'_1 = C_{E_1}/G^{z_F}_{y_1}, E'_2 = C_{E_2}/G^{z_F}_{y_2})$.

We abort if these openings are not the same, i.e. if $(z + z', d, E_1 G^{rnd}, E_2 Y^{rnd}) \neq (z_F, d_F, E'_1, E'_2)$.

We first note that E'_1 and E'_2 are completely determined by the value of z_F , so it suffices to show that we must have $(z + z', d) = (z_F, d_F)$. If this is not the case, this means that the adversary was able to find distinct pairs (z_1, d_1) and (z_2, d_2) such that $G^{z_1}_{y_3} G^{d_1}_d = G^{z_2}_{y_3} G^{d_2}_d$ however, this implies that the adversary knows the discrete log of G_{y_3} with respect to G_d or vice versa, both of which we assume are unknown, so we conclude that finding two such pairs will happen with only negligible probability, and so the game is indistinguishable from G6.

- G8: We just showed that the openings of the forwarding commitment created by the sender can't be changed by the receiver. We will now do the same for authoring commitments. We note that for authoring commitments, the source encryption is chosen by the platform and verified in the clear, so we need only ensure that the message stored in C_A created by an honest sender stays constant.

We modify G7 by adding the requirement that each time `goodAuth` or `goodFwd` is called, we save the message m (which could be \perp) that is associated with the authoring commitment C_A that gets presented to the platform. Then, if the receiving proof verifies, we extract witnesses z_A and d_A from the proof, aborting if the extraction fails. This happens with negligible probability, so this game is indistinguishable from G7.

- G9: After extracting witnesses z_A and d_A in G8, we abort if $H(m) \neq d_A$.

As before, if this is the case, then after being provided a valid opening (z, d) by the honest sender, the adversary found some new (z_A, d_A) such that $G^{z_A}_{y_3} G^{d_A}_d = G^{z_A}_{y_3} G^{d_A}_d$. By the same reasoning as the previous game, this can be achieved only with negligible probability, and therefore this new game is indistinguishable from G4 to an efficient adversary.

By the definition of G9, we know that if a message is successfully sent or reported with attributes d, E_1, E_2 , then either the adversary must have previously received a message with attributes $d, E_1 G^r, E_2 Y^r$ for some random r , or in the case of sending, (E_1, E_2) must decrypt to the sender's identity. (If the attributes came from an unused forwarding or authoring commitment, we would have $d = \perp$ and the protocol would automatically fail).

Because the decryption $E_2 Y^r / (E_1 G^r)^y$ will be constant for any value of r , we can trace back through these extracted receipts to the first time a MAC on attributes (d, E_1, E_2) was created. This must either have been the result of a call to `goodAuth` or `malSend`. If it was from `malSend`, then the sender wasn't an honest user, and so such attributes cannot be used to win the forgeability game.

Otherwise, if it was in `goodAuth`, then the tuple must have been added to \mathcal{M}_{sent} , and so similarly can't be used to win the game.

Therefore, we have shown that there is no way for the adversary to win Game 9. Because we showed that Game 9 was indistinguishable from the *unFORGE* game. We can conclude that any efficient adversary will have negligible advantage against *unFORGE*, and therefore Scheme 2 is unforgeable. \square

C.4 Deniability – Proof of Theorem 5.4

PROOF. As in the first scheme, the complete proof that this scheme satisfies the deniability requirements of a secure source-tracking scheme appears in the full version of this paper. \square