

# Orca: Blocklisting in Sender-Anonymous Messaging

Nirvan Tyagi  
*Cornell University*

Julia Len  
*Cornell University*

Ian Miers  
*University of Maryland*

Thomas Ristenpart  
*Cornell Tech*

**Abstract.** Sender-anonymous end-to-end encrypted messaging allows sending messages to a recipient without revealing the sender’s identity to the messaging platform. Signal recently introduced a sender anonymity feature that includes an abuse mitigation mechanism meant to allow the platform to block malicious senders on behalf of a recipient.

We explore the tension between sender anonymity and abuse mitigation. We start by showing limitations of Signal’s deployed mechanism, observing that it results in relatively weak anonymity properties and showing a new griefing attack that allows a malicious sender to drain a victim’s battery. We therefore design a new protocol, called Orca, that allows recipients to register a privacy-preserving blocklist with the platform. Without learning the sender’s identity, the platform can check that the sender is not on the blocklist and that the sender can be identified by the recipient. We construct Orca using a new type of group signature scheme, for which we give formal security notions. Our prototype implementation showcases Orca’s practicality.

## 1 Introduction

End-to-end (E2E) encrypted messaging, now relied upon by billions of people due to products like Signal, WhatsApp, Facebook Messenger, and more, provides strong E2E confidentiality and integrity guarantees [5, 23]: the messaging platform itself cannot read or modify user messages. The E2E encryption protocols used [52] do not, however, attempt to ensure anonymity, so the platform learns the sender and recipient of every message sent over the network. While academic systems [4, 6, 24, 25, 42, 45, 47, 57, 60, 61] have developed protocols that hide the identity of senders and receivers from platforms, they introduce expensive overheads.

A recent suggestion for pragmatic privacy improvements is to aim solely for sender anonymity. Introduced by Signal in a feature called “sealed sender” [48], sender anonymity ensures that the sender’s identity is never revealed via messages to the platform, e.g., the sender does not authenticate with an account password or digital signature; messages reveal only the intended recipient. While sealed sender does not hide network-level identifiers such as IP addresses, one can do so by composing it with Tor [27] or an anonymous broadcast [24, 41, 47, 53, 61].

In this work, we explore a key tension in sender-anonymous systems: mitigating abuse by malicious senders. Already E2E

encryption makes some kinds of abuse mitigations, such as content-based moderation, more challenging (c.f., [28, 30, 35, 58]). Sender anonymity complicates the setting further because the lack of sender authentication means that the platform cannot block unwanted messages on behalf of a recipient in a conventional way.

To enable platform blocking, Signal’s sealed sender has a user distribute an access key to their contacts that senders must show to the platform when sending the user a sender-anonymous message. If a sender cannot provide an access key, the platform drops the message. A user that blocks a sender in their client triggers a rotation of this key and a redistribution to the (remaining) contacts. Future messages from the blocked sender will be dropped by the platform.

We observe two deficiencies with this approach. First, access keys must be distributed over non-sender-anonymous channels, meaning the platform learns the identities of users who can send sender-anonymous messages to a particular recipient. This significantly lowers the anonymity guarantee—in the limit of having only a single contact, there is no anonymity at all.

Second, we show a simple “griefing” attack that works despite the anti-abuse mechanism. By design, the sender is hidden from the platform, and only the recipient can identify the sender of a sender-anonymous message. However, a malicious sender can trivially craft malformed messages that even the recipient will not be able to identify. The recipient’s client rejects these messages, but not before processing them. This is particularly problematic for mobile clients as it uses up battery life; we experimentally verify that an attacker can easily drain a target’s battery in a short period of time. To make matters worse, neither victim nor platform can identify the attacker, and so the victim will not know who to block.

We design a new abuse mitigation mechanism for privacy-preserving blocklisting in sender-anonymous messaging. Our protocol, called Orca, allows recipients to register a blocklist with the platform. The blocklist is privacy-preserving, meaning it does not reveal the identities of the blocked users. Senders construct messages that are anonymous to the platform, but can be verified by the platform as being attributable to a sender not present on the blocklist. If the sender is on the blocklist or if the message is malformed, then the platform rejects the message; if the message is delivered, the recipient is guaranteed to be able to identify the sender.

Importantly, Orca provides a new non-interactive initialization functionality that allows a user to initiate sender-anonymous messages without having previously communicated with the recipient. This significantly enhances the anonymity guarantees, because it expands the anonymity set to be as large as all registered users of the system.

In summary, our contributions are:

- We build a threat model for sender-anonymous messaging and identify limitations in previous approaches, including a new griefing attack against Signal’s sealed sender that we evaluate.
- We construct a new group signature scheme [22] to make up the core of Orca’s functionality. The new primitive is tailored to the needs of our setting and supports multiple openers, keyed verification, and local revocation; see Section 4 for details. We provide new security definitions, building upon ones from prior work [8, 14].
- We show an extension of Orca that integrates mechanisms from anonymous credentials [20] to arrange that the relatively expensive group signature scheme is only used periodically when initiating a new conversation. Initialization will generate a batch of one-time-use sender tokens [43, 44], which can be spent to authenticate messages and replenished at very low cost.
- We implement and evaluate Orca, suggesting that it is sufficiently performant to deploy at scale. In particular, once initialized, the token-based extension incurs only 30B additional bandwidth cost per message and only one extra group exponentiation of computation for clients; the platform need only compute a group exponentiation and check the token against a strikelist. The computational cost for the platform is paid during initialization which incurs work on the order of the size of the recipient’s blocklist ( $\sim 200\text{ms}$  for a blocklist of length 100). We find that a medium-provisioned server can comfortably support a deployment of a million users depending on frequency of conversation initialization.

## 2 Setting: Sender Anonymity for E2EE

This work focuses on sender-anonymous E2E encrypted messaging hosted by a centralized messaging platform. In this section and throughout the body, we will often use Signal as our running example. However, the techniques that we introduce are relevant for any sender-anonymous messaging system in which the platform learns the recipient identity.

### 2.1 Background: Signal and Sealed Sender

**Non-sender-anonymous E2EE messaging.** We first briefly outline Signal’s non-sender-anonymous protocol. For simplicity we restrict attention to one client per user. A user wishing to send a message first registers an account with the platform using a long-lived identity public key  $pk_s$ , retaining the associated secret key  $sk_s$ . The user then must contact the platform

to obtain the long-lived public key  $pk_r$  of their intended recipient. Once this phase is complete, a client can securely send messages via Signal’s *double ratchet* protocol [52]. This provides state-of-the-art message confidentiality guarantees even in the event of key compromise [5, 23].

Signal, like most other E2E encrypted messaging platforms, requires users to authenticate their account when sending and receiving messages. Importantly, this allows for abuse prevention because the platform can block malicious senders, and even block senders from talking to a specific recipient. On the other hand, such account authentication, e.g., via public key signature or unique account password, does not provide cryptographic sender anonymity.

**Sender anonymity with sealed sender.** Sealed sender is Signal’s protocol [48] for cryptographic sender anonymity motivated by their desire to minimize the amount of trust their users must place in the platform. We will now walk through a high level summary of how sealed sender works.

*Initialization and key exchange.* As before, senders must first register a public key  $pk_s$  with the platform. The user is issued a short-lived *sender certificate* from the platform, that we denote by *cert*. The certificate contains a digital signature by the platform in order to attest to the validity of the user’s identity key. These certificates must be periodically updated, requiring the user to rerun the registration protocol.

To receive sealed messages a recipient must generate their long-lived identity key pair  $(pk_r, sk_r)$  as usual, but now additionally generate a 96-bit *access key* that we denote by *ak*. Both  $pk_r$  and *ak* are registered with the platform. Looking ahead, senders will need to show *ak* to the platform to send a sealed message. This means that the recipient must distribute *ak* to whomever they want to grant the ability to send sealed messages. By default, the access key is distributed to all contacts of a user through Signal’s original non-sender-anonymous channel. Additionally, users can opt into accepting sealed messages from anyone, including non-contacts. In this case, senders do not need a recipient’s access key to send them sealed messages.

*Sending a sealed message.* The pseudocode for sending and receiving a message via sealed sender is provided in Figure 1. It is designed to work modularly as a layer on top of any non-sender-anonymous E2E encryption protocol. At a high level, the protocol creates two ciphertexts: (1) an identity ciphertext encrypting the sender’s long-lived public key  $pk_s$  to the recipient, and (2) a content ciphertext encrypting the standard E2E encryption ciphertext along with the sender certificate. The identity ciphertext and content ciphertext cryptographically hide the sender identity even if the underlying E2E encryption ciphertext does not<sup>1</sup>.

More specifically, the protocol encrypts the sender identity

<sup>1</sup> Signal’s use of the double ratchet algorithm produces ciphertexts that can either include the sender identity in plaintext or include messaging metadata such as counters used for in-order processing that would leak information useful for linking senders.

<b>SealedSender.Send(<math>m</math>)</b>	<b>SealedSender.Rcv(<math>pk_e, ct_{id}, ct_{ss}</math>)</b>
$ct_m \leftarrow \text{ratchet.Enc}(m)$	$\text{salt}_1 \leftarrow (pk_r, pk_e)$
$(pk_e, sk_e) \leftarrow \text{KeyGen}()$	$(e_{\text{chain}}, k_e) \leftarrow \text{HKDF}(\text{salt}_1, pk_e^{sk_r})$
$\text{salt}_1 \leftarrow (pk_r, pk_e)$	$pk_s \leftarrow \text{AE.Dec}(k_e, ct_{id})$
$(e_{\text{chain}}, k_e) \leftarrow \text{HKDF}(\text{salt}_1, pk_r^{sk_e})$	$\text{salt}_2 \leftarrow (e_{\text{chain}}, ct_{id})$
$ct_{id} \leftarrow \text{AE.Enc}(k_e, pk_s)$	$k \leftarrow \text{HKDF}(\text{salt}_2, pk_s^{sk_r})$
$\text{salt}_2 \leftarrow (e_{\text{chain}}, ct_{id})$	$\text{cert}    ct_m \leftarrow \text{AE.Dec}(k, ct_{ss})$
$k \leftarrow \text{HKDF}(\text{salt}_2, pk_r^{sk_s})$	$b \leftarrow \text{Verify}(pk_s, \text{cert})$
$ct_{ss} \leftarrow \text{AE.Enc}(k, \text{cert}    ct_m)$	If $b = 0$ then return $\perp$
Return $(pk_e, ct_{id}, ct_{ss}), ak$	$m \leftarrow \text{ratchet.Dec}(ct_m)$
	Return $m$

Figure 1: Pseudocode for Signal’s sealed sender feature.

$pk_s$  via a variant of hashed ElGamal [2] to produce the identity ciphertext  $ct_{id}$ . In particular, it generates ephemeral key pair  $(pk_e, sk_e)$  and makes use of a hash-based key derivation function HKDF and authenticated encryption scheme AE. The sender then encrypts the plaintext  $m$  using the original double ratchet algorithm  $\text{ratchet.Enc}(m)$ . It bundles the resulting ciphertext  $ct_m$  and sender certificate  $\text{cert}$  and encrypts this with a key derived from long-lived identity keys  $pk_s$  and  $pk_r$  to produce the content ciphertext  $ct_{ss}$ . The sender indicates the intended recipient and sends the triple  $(pk_e, ct_{id}, ct_{ss})$  along with the recipient’s access key  $ak$  to the platform.

Upon receipt of the sender’s message, the platform checks that the intended recipient’s registered access key matches  $ak$ . If this check passes, then the platform forwards the triple  $(pk_e, ct_{id}, ct_{ss})$  to the recipient. The recipient decrypts as shown in Figure 1. Once it recovers  $\text{cert}$  and  $ct_m$ , it verifies the sender as a valid account using the certificate and the recovered identity key  $pk_s$ . If the sender’s identity is authenticated, then  $ct_m$  is decrypted using the double ratchet algorithm.

## 2.2 Limitations of Sealed Sender

There are limitations to Signal’s sealed sender protocol for sender anonymity, which we raise here in the form of three different classes of attacks.

**Traffic analysis of sender-anonymous messages.** An inherent leakage of the sender-anonymous messaging setting (as opposed to the sender- and recipient-anonymous setting) is that the recipient of each message is inherently leaked to the platform. Martiny et al. [49] demonstrate a set of statistical disclosure attacks that use this leakage to infer communicating partners, for example, by searching for users with interleaving messages suggesting a back-and-forth conversation pattern. They provide a modification to Signal’s sealed sender that protects against traffic analysis of sender-anonymous messages, which they call “sender-anonymous conversations”. This mitigation approach, as well as another separate approach which instead relies on random message delays and/or noise messages [53], do not provide solutions for blocklisting. The techniques we introduce for supporting blocklists compose well with these traffic analysis mitigations. Given this prior

work, we do not explicitly address traffic analysis of sender-anonymous messages beyond considering the anonymity set, as we discuss next.

**Traffic analysis of non-sender-anonymous messages.** Recall that access keys are distributed through Signal’s original non-sender-anonymous channel. While this setup is still encrypted, the platform nevertheless observes with whom the user exchanged non-sender-anonymous messages. Thus, when a sender anonymously authenticates using  $ak$ , the set of users that could correspond to the sender (i.e., the *anonymity set* of the sender) is restricted and known to the platform. This means, for example, if a recipient only has a single contact with which they have communicated, there is no sender anonymity at all. Furthermore, if a user rotates their access key to revoke sending access, this resets their anonymity set of senders, as their new access key must be redistributed.

Martiny et al. [49] assume in their threat model that these access keys have already been exchanged between communicating parties. Their attack can therefore be improved by tracking the sender anonymity set of a recipient learned by the platform. Notably, our solution for blocklisting will prevent such improvements.

**Griefing attack by evading identification.** Sealed sender relies on the sender to self-identify to the recipient: the platform can not check for malformed messages. Instead, the recipient must decrypt and check validity of the sender identity key and certificate, dropping messages that do not verify. This allows for a straightforward griefing attack in which an attacker can spam the recipient with untraceable messages, causing the recipient’s device to suffer battery drain and to consume bandwidth, a type of user-mounted DoS attack.

We demonstrate through a proof-of-concept implementation that this griefing attack is effective. Our attack simply modifies  $pk_e$  in  $(pk_e, ct_{id}, ct_{ss})$  to a random value  $pk_f$ . To the platform this is indistinguishable from a legitimate sealed sender message, but the recipient’s decryption will fail when trying to decrypt  $ct_{id}$ . The recipient cannot recover any information about the sender. Running experiments on a Google Pixel phone running Android 9, we find that sending just 1 message every 10 seconds causes the battery to drain at an increased rate of  $9\times$  over baseline. We provide more extensive measurements of this attack in Appendix A.

Ultimately, there are no satisfying mitigation options available to victims (see last section of Appendix A). If the victim of the attack has opted in to accepting sealed sender messages from non-contacts, the attack can be mounted by anyone. Otherwise the attacker needs the recipient’s access key, meaning the attacker must be one of the victim’s contacts (or has found some other way to obtain the access key). While this limits who can mount the attack in the default case, it is still problematic: The victim can rotate their access key  $ak$  and attempt to redistribute a new  $ak'$  to their communicating partners. If the attacker is not able to get access to the new access key, the attack will be stopped by the platform and no messages will

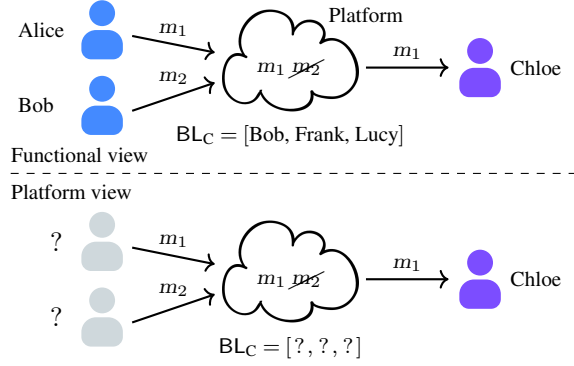


Figure 2: Privacy-preserving, outsourced blocklisting for sender-anonymous messaging. The platform is able to block messages from users on Chloe’s blocklist without learning their identity. The top view shows the functionality of outsourced blocklisting, while the bottom view shows what is revealed to the platform. Not shown, Chloe can also efficiently identify the sender of message  $m_1$  as Alice and update her blocklist  $BL_C$  if needed.

reach the victim’s client. But since the attack leaves no information about which of the victim’s communicating partners is responsible, the victim can only make a guess as to whom they should block.

Realistically to maintain usability of their mobile device, a user may limit Signal to only a few highly trusted contacts, or will push the user off Signal to a less private messenger. We consider both of these outcomes to be highly damaging to vulnerable users that would benefit from a metadata-private messenger. Looking forward, we will want a mechanism that provides the user more granular recourse against misbehaving senders.

### 3 Outsourced Blocklisting

We now turn to building a new sender-anonymous messaging protocol that avoids the current weaknesses of sealed sender. Our approach is to enable what we call privacy-preserving outsourced blocklisting (see Figure 2).

**Goals.** Such a system should enjoy the following features:

- *Sender anonymity:* Messages cryptographically hide the sender identity from the platform.
- *Sender attribution:* Recipients can cryptographically verify the sender of any ciphertexts delivered by the platform.
- *Blocklisting:* Recipients can register a blocklist with the platform and update it efficiently. The platform can use the blocklist to drop sender-anonymous messages from senders that the recipient has added to the blocklist.
- *Blocklist anonymity:* The blocklist should not reveal the identities of the senders blocked by the recipient.

Together these properties prevent the type of griefing attacks that affect sealed sender: a client receiving problematic messages can identify the sender and instruct the platform to drop

them on the client’s behalf.

We would also like the system to support:

- *Non-interactive initialization:* Users can begin sending sender-anonymous messages without previous interaction with the intended recipient.

This property obviates the use of non-sender-anonymous channels to initiate sender-anonymous communication. In particular, the platform should not be able to attribute messages to some smaller subset of users, as messages can have originated from any registered user of the system.

Orca is designed to accompany a sender-anonymous E2EE messaging protocol to provide the functionality of outsourced blocklisting while carrying over both the sender-anonymity and message confidentiality properties of the underlying protocol. As such, we assume the underlying E2EE protocol is sender-anonymous, and if it is not, can easily be made so using encapsulation techniques similar to sealed sender (see Figure 1). Our protocol will provide a registration process in which users interact with the platform to generate the required keys for the protocol; this will be done at the same time users register for the underlying E2EE protocol. To send a message, the sender first encrypts the message plaintext  $pt$  to the recipient as specified by the E2EE protocol. Then, Orca will concern itself with authenticating the delivery of the produced E2EE ciphertext; the authenticity of the underlying message plaintext needs to be provided by the E2EE protocol. We will refer to the E2EE ciphertext as the “message” from Orca’s perspective.

**Threat model.** We assume an active, persistent adversary that controls the messaging platform and an arbitrary number of users. We assume the clients of legitimate users are not compromised and that they correctly abide by the protocol.

Our primary concern is the *cryptographic anonymity* of the messaging protocol. The adversary, even with active deviations from the protocol, should not be able to learn sender identity information from the contents of protocol messages.

Even in the case that anonymity is achieved at the message protocol layer, identification information can leak through the network layer, e.g., by associating IP addresses or by making inferences based on timing. We consider preventing such leakage to be orthogonal to the goal of providing a blocklisting solution for the message protocol layer: existing solutions for mitigating network leakage will compose. Sender-anonymous channels resilient to linking attacks that exploit IP addresses can be constructed using services such as Tor [27]; linking attacks performed by stronger global network adversaries with the ability to observe and inject traffic along any network link can be mitigated using prior academic solutions for anonymous broadcasting [24, 41, 47, 53, 61]. Lastly, as discussed in Section 2.2, given a sender-anonymous channel, timing analysis of messages with designated recipients can be mitigated using existing techniques [49, 53].

It is trivial for an active adversary that controls the platform



to deny service to arbitrary users by not delivering messages. In future work, it may be valuable to provide a mechanism for honest users to provably expose such misbehavior, but in this work we leave platform-mounted denial-of-service (DoS) attacks out of scope. On the other hand, we do want to protect against user-mounted DoS attacks, in which a malicious user can interact with an honest platform to deny service to other users, as in the grieving attack.

**Overview.** We will now provide an overview of Orca’s design by stepping through a series of strawman constructions.

Sender-specific one-time use access tokens. Instead of having all senders authenticate by reusing the same shared access token, the recipient can deal unique access tokens to each sender. Reusing a sender-specific token allows linking by the platform, so these tokens will necessarily be one-time use only. We outline a version of this approach that is taken by the Pond messaging system [43, 44].

On registration, recipients register a key  $k$  to a pseudorandom function  $F$ , e.g. HMAC, with the platform. Recipients distribute one-time use tokens of the form  $(x, y = F(k, x))$  for random values  $x$  to senders. The platform verifies these tokens using  $k$  and the recipient can identify senders since they know to whom they dealt  $(x, y)$ . A sender’s tokens are refreshed in the normal exchange of messages. Now a recipient can block by reporting the unused tokens of a sender to the platform; the platform tracks these tokens along with previously spent tokens for a recipient in a strikelist and rejects incoming messages that authenticate with struck tokens. The platform’s strikelist grows unbounded as more messages are sent, but this cost can be managed by scheduled key rotations.

This blocklisting approach improves significantly over sealed sender as it effectively removes the grieving attack vector, however it does not address the concerns around leakage during initialization: the recipient still initially distributes the access tokens over non-sender-anonymous channels to senders, revealing to the platform a small set of possible senders for future messages. A different approach is needed to provide stronger sender anonymity with non-interactive initialization.

Group signatures. A promising starting point for sender-anonymous blocklisting with non-interactive initialization is *group signatures*, a well-studied cryptographic primitive [7, 8, 11, 17, 22]. Group signature schemes allow users to sign messages anonymously on behalf of a group whose membership is controlled by a *group manager*. Signatures appear anonymous to everyone except to a special *opening authority* who has the ability to deanonymize the signer and revoke their signing ability.

Our next strawman solution has the platform maintain a separate group signature scheme for each registered user, where the user is the opening authority and the platform is the group manager. A sender registers with the platform under the desired recipient’s group signature scheme. The sender sends

their message along with a signature on the message under the recipient’s group to the platform. The platform then verifies the anonymized signature. For blocklisting, we use a group signature scheme that supports *verifier-local revocation* [14]. This means that the recipient can revoke senders by communicating *only* with the platform (i.e., verifier).

This strawman provides effective sender attribution and blocklisting. It also allows senders to acquire group signature credentials without previous interaction with the recipient. However, messages to a recipient can be attributed by the platform to the set of users that registered under the recipient’s group signature scheme, so we do not achieve our stronger anonymity goal. Furthermore, existing group signatures that meet our requirements use expensive bilinear pairing operations, adding on to the efficiency concerns of managing a separate scheme for each registered user.

We resolve these issues by proposing a new type of group signature that introduces two novel features. The first is support for *multiple opening authorities*. This will dispense with the per-recipient group signature schemes and the need to register separately for each recipient that you wish to send to. The second feature is *keyed-verification*, in which we observe that the platform is also the only verifier. Removing public verifiability improves efficiency of client-side operations.

This new group signature, presented in Section 4, makes up the core of Orca. However even with our optimizations, e.g., keyed-verification, the group signature approach incurs significant computational cost, in particular for the platform, owing to the use of verifier-local revocation: verifying a signature incurs work linear in the size of the recipient’s blocklist.

Hybrid: Group signature with one-time tokens. This leads us to our final construction which combines the use of group signatures for non-interactive initialization with one-time use tokens for efficient authentication of subsequent messages. Here, the group signature is used to allow the sender to acquire its first batch of tokens from the platform. The main contribution of this approach is a new protocol for allowing the platform to dispense tokens on behalf of the recipient. This is challenging because the platform should not be able to link newly minted tokens to a sender, but it must provide a way for the recipient to learn to whom new tokens were dealt (for future sender attribution). We construct this protocol by adapting techniques from blinded issuance of anonymous credentials [20]. After this (relatively) expensive initialization procedure, users exchange new tokens in the normal flow of conversation and the system enjoys all the efficiency benefits of the token-based protocol. We describe Orca’s one-time token extension in Section 5.

## 4 Orca’s Group Signature

Our main construction is based on a novel group signature scheme. In this section, we will introduce our new group signature abstraction, describe how to use it to construct an

outsourced blocklisting protocol, and lastly provide an instantiation of such a group signature,

#### 4.1 Group Signature Syntax and Security

Group signatures [22] allow users to sign messages anonymously on behalf of a group. The basic setting is as follows. The membership of a group is coordinated by a *group manager*, with whom users register with in order to join the group. Additionally, anonymous group signatures can be opened (traced) to identify the signing user in the group by a designated *opening authority*.

We make use of three extensions to the basic group signature setting.

- (1) *Verifier local revocation*: A group signature supporting revocation allows the opening authority to additionally revoke the signing ability of group members. *Verifier local* revocation means that to revoke a member, the opening authority need only communicate a revocation message to verifying parties (as opposed to both verifying parties and group members); revocation does not affect the way group members sign messages.
- (2) *Multiple opening authorities*: An opening authority is created through registration with the group manager. Group members sign messages designated to one of many opening authorities, and only the opening authority that a signature is designated to is able to open the signature to the signer's identity. Revocation is handled separately per opening authority, meaning a group member may be able to sign messages designated for some opening authorities, but be revoked from signing messages to others.
- (3) *Keyed verification*: Verification of group signatures can only be completed by a secret key owned by the group manager and shared to verifying parties. This is particularly useful in cases where the group manager is the only party verifying signatures and allows for more efficient schemes than those that achieve public verifiability.

Verifier local revocation has been previously studied [14], but the other two extensions are novel to the best of our knowledge. The model and following security definitions for our new setting are derived from [8, 14].

**Syntax.** A multi-opener, keyed-verification group signature scheme  $GS$  is run between three types of participating parties: (1) users  $U$  that join the group and sign messages, (2) opening authorities  $OA$  that can trace signatures to signers, and (3) a group manager  $GM$  to coordinate registration and perform verification. It consists of the following algorithms:

- $pp \leftarrow GS.Setup(\lambda)$ : The setup algorithm defines the public parameters  $pp$ . We will assume  $pp$  is available to all algorithms, and all parties have assurance it was created correctly.
- $(gmpk, gmsk) \leftarrow GS.Kg_{GM}^{pp}()$ : The key generation al-

gorithm is run by the group manager to generate a public key  $gmpk$  and secret key  $gmsk$ .

- $GS.JoinU_U^{pp} \leftrightarrow GS.IssueU_{GM}^{pp}$ : Group registration is an interactive protocol implemented by  $GS.JoinU$  and  $GS.IssueU$  run between a user and the group manager, respectively. If execution is successful, the user will receive a public, secret key pair  $(upk, usk)$  and the group manager will receive  $upk$ , else both parties receive  $\perp$ . If the protocol accepts, the group manager will store  $upk$  in a global registration table and reject duplicate  $upk$  registrations.
- $GS.JoinOA_{OA}^{pp} \leftrightarrow GS.IssueOA_{GM}^{pp}$ : Opening authority registration is an interactive protocol run between a prospective opening authority and the group manager. If execution is successful, the opening authority will receive a public, secret key pair  $(oapk, oask)$  and the group manager will receive and store  $oapk$  in the registration table, else both parties receive  $\perp$ .
- $\sigma \leftarrow GS.Sign_U^{pp}(usk, gmpk, oapk, m)$ : The signing algorithm is run by a group member to produce a group signature  $\sigma$  on a message  $m$  designated for opening authority  $oapk$ .
- $upk \leftarrow GS.Open_{OA}^{pp}(oask, m, \sigma)$ : The opening algorithm is run by an opening authority to learn the identity of the signing user  $upk$ , and returns  $\perp$  upon failure.
- $\tau_R \leftarrow GS.Revoke_{OA}^{pp}(oask, upk)$ : The revocation algorithm is run by an opening authority to create a revocation token  $\tau_R$  for a user  $upk$ . The opening authority sends the revocation token to the group manager who includes it in a revocation list  $RL$  used for verification.
- $b \leftarrow GS.Ver_{GM}^{pp}(gmsk, oapk, RL, m, \sigma)$ : The verification algorithm is run by the group manager to determine if an input signature  $\sigma$  and  $m$  are valid under a designated opening authority  $oapk$  and revocation list  $RL$ .

As mentioned, we assume some global registration table that contains all user public keys  $upk$  and opening authority public keys  $oapk$  that succeed registration. In practice, such a table might be implemented with a public key infrastructure (PKI) supporting key transparency audits [50] allowing it be hosted by the untrusted platform. Additionally, for simplicity, we may drop the executing party from the subscript and the public parameters from the superscript if their use is clear from context.

**Correctness and security notions.** We extend the standard notions of correctness and security from [8, 14]. Here, we describe correctness and then the three security properties: anonymity, traceability, and non-frameability. The properties are formalized via security games involving an adversary in the full version [59].

The *correctness* property concerns signatures generated by honest group members. An honestly generated signature should pass verification under all honestly generated revoca-

tion lists that do not include a revocation token for the signing user created by the designated opening authority. An honestly generated signature should also be opened to the correct signing user by the designated opening authority.

The *anonymity* property captures that an adversary without access to the designated opening authority’s key should not be able to determine the signer of a signature among unrevoked group members. The adversary has the power of an actively malicious group manager and may adaptively compromise group members and opening authorities. More specifically, we target *CCA-selfless-anonymity* [11] meaning signatures are not anonymous to the signer (selfless) and the adversary has access to an opening oracle throughout the security game (CCA). We consider rogue key attacks, allowing the adversary to create public keys for corrupted parties, but require the adversary to prove knowledge of secret keys. We model this, for simplicity, by asking the adversary to produce the secret key for generated public keys following the knowledge of secret key model of [10], which can be instantiated with extractible proofs of knowledge. We also provide an extension of our anonymity game to capture anonymity of revocation tokens (in addition to signatures) that is, to our knowledge, the first definitional attempt at doing so.

*Traceability* ensures that every signature that passes verification can be opened by the designated opening authority to a registered user. Traceability necessarily considers an adversary that does not control the group manager since it is trivial for the group manager to craft signatures for unregistered public keys. However, traceability is accompanied by *non-frameability* which ensures that it is not possible to forge a signature that opens to an honest user; non-frameability considers a stronger adversary that controls the group manager as in anonymity.

**Bilinear pairing groups.** Our construction will make use of bilinear pairing groups for which we will use the following notation. (1) Groups  $\mathbb{G}_1, \mathbb{G}_2, \mathbb{G}_T$  are cyclic groups of prime order  $p$ . (2) Group element  $g_1$  is a generator of  $\mathbb{G}_1$ ,  $g_2$  is a generator of  $\mathbb{G}_2$ . (3) Pairing function  $e : \mathbb{G}_1 \times \mathbb{G}_2 \rightarrow \mathbb{G}_T$  is a computable map with the following properties: *Bilinearity*:  $\forall u \in \mathbb{G}_1, v \in \mathbb{G}_2$ , and  $a, b \in \mathbb{Z}$ ,  $e(u^a, v^b) = e(u, v)^{ab}$ , and *Non-degeneracy*:  $e(g_1, g_2) \neq 1$ . We assume an efficient setup algorithm that on input security parameter  $\lambda$ , generates a bilinear group,  $(p, \mathbb{G}_1, \mathbb{G}_2, \mathbb{G}_T, g_1, g_2, e) \leftarrow \mathcal{G}(\lambda)$ , where  $|p| = \lambda$ .

## 4.2 Outsourced Blocklisting from Group Signatures

Given a keyed-verification, multi-opener group signature with verifier-local revocation, we build our core protocol, detailed in Figure 3. The platform plays the role of the group manager. Users register with the platform as both a user of the group and as an opening authority, receiving keys  $(usk_i, oask_i)$ . For user  $i$  to send a message to user  $j$ , assume for now that user  $i$  has user  $j$ ’s public keys  $(upk_j, oapk_j)$ . We will describe how user  $i$  obtains these keys shortly.

User  $i$  signs their message with  $usk_i$  under the group signature scheme designating  $oapk_j$  as the opening authority. The platform verifies the anonymous group signature against user  $j$ ’s revocation list, and if it verifies, delivers the message and signature to user  $j$ , who can then identify the sender,  $upk_i$ , by opening the signature. Users can blocklist a sender  $upk_i$  to the platform by generating a revocation token under their opening authority key  $oask_j$  and sending it to the platform. Anonymity of the group signature and revocation tokens ensure that the platform does not learn sender identity information from messages or from the blocklist; and traceability and non-frameability ensure recipients will be able to properly attribute received messages to a sender.

To achieve our stronger sender anonymity goal, user  $i$  must be able to read the public key information of user  $j$  needed to start a conversation *without* revealing their own identity to the platform. Since public key information is not sensitive, the platform can provide unrestricted access to PKI lookups that do not require user authentication. Note that the platform can observe the *number* of lookups to a recipient’s public key, but learns no information on which users are making those lookups. We discuss how the platform can restrict access to resources and maintain anonymity in Section 8.

## 4.3 Construction of Group Signature

Our group signature follows closely the “certified signature” recipe that many group signatures take [34]. In this recipe, the group manager registers users by certifying their public key  $Y = g^y$ ; the user’s group key is made up of their secret identity key  $y$  along with the group manager’s certificate  $t$ . To sign a message under the group, the user encrypts their public key to the opening authority creating an *identity ciphertext* where  $Z$  is the opening authority’s encryption key.

$$ct_{id} \leftarrow (g_1^{\alpha_{ct}}, YZ^{\alpha_{ct}}) \quad \alpha_{ct} \leftarrow \mathbb{Z}_p$$

They then prove in zero knowledge that they have a certificate from the group manager on the same public key that is enclosed in the ciphertext *and* that they know the secret key associated to it. The signature is verified by verifying the zero knowledge proof and can be opened by the opening authority simply by decrypting the identity ciphertext.

This recipe naturally extends to support a scheme with multiple opening authorities. The identity ciphertext is encrypted using the public key of the designated opening authority.

**Supporting verifier-local revocation.** An opening authority registers with two keys: (1) an encryption key  $(z, Z = g_1^z)$ , and (2) a revocation key  $(w, W = g_1^w)$ , where  $oapk = (W, Z)$ . We have described how a user with identity key  $(y, Y = g_1^y)$  encrypts their public key  $Y$  to the opening authority. To revoke a user’s signing ability, the opening authority constructs a user-specific *revocation token* as the Diffie-Hellman value between the user’s public key and their own revocation key,  $\tau_R = Y^w$ . Intuitively, these revocation tokens are anonymous since a Diffie-Hellman value looks random to a verifier that



---

### Protocol 1: Orca Outsourced Blocklisting Protocol

---

#### Setup:

- (1) Public parameters for the group signature scheme are generated,  $pp \leftarrow \text{GS.Setup}(\lambda)$ .
- (2) The platform initializes its state as the group manager of the group signature scheme.
  - (a)  $(gmpk, gmsk) \leftarrow \text{GS.Kg}_{\text{GM}}^{pp}()$
  - (b)  $T_U \leftarrow []$ : Table tracking user public keys.
  - (c)  $T_R \leftarrow []$ : Table tracking user revocation tokens.

#### Registration:

- (1) User registers with platform to acquire group signature signing key with which to send messages,  $\text{GS.JoinU}_U^{pp} \leftrightarrow \text{GS.IssueU}_{\text{GM}}^{pp}$ . User stores  $usk$  and platform stores  $upk$ .
- (2) User registers as opening authority and generates keys with which to receive messages,  $\text{GS.JoinOA}_{\text{OA}}^{pp} \leftrightarrow \text{GS.IssueOA}_{\text{GM}}^{pp}$ . User stores  $oask$  and platform stores  $oapk$ .
- (3) Platform stores public keys in  $T_U[upk] \leftarrow oapk$ .
- (4) Platform initializes empty revocation token list for user,  $T_R[oapk] \leftarrow []$ .

#### Sending a message:

- (1) [Optional] Sender anonymously requests recipient public key ( $oapk$ ) and/or rate-limited pre-keys from platform (described in Section 8).
- (2) Sender signs message specifying the recipient as the opening authority (with recipient's  $oapk$ ),  $\sigma \leftarrow \text{GS.Sign}_U^{pp}(usk, gmpk, oapk, m)$ . Sender sends message, signature, and recipient to platform,  $(m, \sigma, oapk)$ .
- (3) Platform checks validity of signature against recipient's revocation list,  $b \leftarrow \text{GS.Ver}_{\text{GM}}^{pp}(gmsk, oapk, T_R[oapk], m, \sigma)$ . If  $b = 1$ , then platform delivers  $(m, \sigma)$  to recipient.

#### Blocklisting a user:

- (1) Recipient generates and sends anonymous revocation token to platform,
    - (a)  $upk \leftarrow \text{GS.Open}_{\text{OA}}^{pp}(oask, m, \sigma)$
    - (b)  $\tau_R \leftarrow \text{GS.Revoke}_{\text{OA}}^{pp}(oask, upk)$
  - (2) Platform adds revocation token to recipient's blocklist,  $T_R[oapk] \leftarrow T_R[oapk] \cup \{\tau_R\}$ .
  - (3) [Optional] Recipient stores identities of blocklisted senders and/or reports sender identity to platform (described in Section 8).
- 

Figure 3: Core protocol based on group signature.

does not know the secret keys  $y$  or  $w$ .

To allow a verifier in possession of a user's revocation token to identify signatures from a user, we need something more. In addition to the identity ciphertext, the user also constructs a *revocation ciphertext* enclosing their revocation token,  $\tau_R = W^y$ . This "ciphertext" is constructed to be undecryptable, but includes a backdoor for testing whether a plaintext  $pt$  is enclosed (following the approach of Boneh and Shacham [14]).

$$ct_R \leftarrow (M_1^{\alpha_T}, \tau_R N_1^{\alpha_T}) \quad \alpha_T \leftarrow \mathbb{Z}_p \quad M_1, N_1 \leftarrow \mathbb{G}_1$$

The backdoor of  $ct_R$  consists of the isomorphic  $\mathbb{G}_2$  elements  $M_2, N_2$ . The verifier can check whether  $\hat{\tau}_R$  is enclosed in  $ct_R$  via the following test using the pairing function  $e$ :

$$e(T_2/\hat{\tau}_R, M_2) \stackrel{?}{=} e(T_1, N_2) \quad (T_1, T_2) \leftarrow ct_R$$

The verifier performs this test for each revocation token in an

opening authority's revocation list and outputs 1 if no revocation token matches and the signature's proof verifies. The signature's proof now additionally proves the well-formedness of  $ct_R$  with respect to user public key  $Y$ .

**Improving efficiency with keyed-verification.** A central part of the group signature is that the user must prove they have a certificate on their public key from the group manager. Creating this proof, even for certificate signatures designed for this purpose [11, 19], is relatively expensive, with known constructions requiring multiple pairings to be evaluated. In our setting, the platform plays the role of both the group manager and the sole verifier; all messages pass through the platform. This setting allows us to bring in techniques from keyed-verification anonymous credentials [20]. Specifically, during user registration, instead of receiving a signature from the group manager, users receive a MAC  $t$  on their public key from an algebraic MAC scheme; our construction uses  $\text{MAC}_{\text{GGM}}$  from [20, 29]. Proving knowledge of a valid MAC is more efficient and, in particular, does not require pairing evaluations. The resulting proof can only be verified using the secret MAC key (held by the group manager), hence our introduction of the keyed-verification setting for group signatures (i.e., "group MACs"). This optimization limits the use of pairings in our group signature only to the revocation token tests made by the group manager during verification.

**Summary.** In total, our group signature is composed of three components, (1) the identity ciphertext  $ct_{id}$  enclosing the signer's public key to the opening authority, (2) the revocation ciphertext  $ct_R$  enclosing the revocation token, and (3) a zero knowledge proof  $\pi$  that (1) and (2) were constructed properly with knowledge of a key pair  $(y, Y)$  and a MAC  $t$  on  $Y$ . The full details of the construction along with security proofs with respect to the formal definitions of anonymity, traceability, and non-frameability are presented in the full version [59].

As stated, every time a user sends a message, they create a group signature and the platform verifies the group signature. Even with our optimizations, this involves the platform running a verification algorithm that is linear in the size of the recipient's revocation list. We improve in the next section, extending Orca with one-time use sender tokens to make the need for a group signature a rare event.

## 5 Extending Orca with One-time Use Tokens

In this section, we describe how to reduce Orca's reliance on its core group signature protocol. Instead of creating and verifying a group signature for every message sent, the group signature will only be used periodically to mint new batches of one-time use sender tokens from the platform. Messages can be sent, with very little cost, by including a valid token for a recipient. Furthermore, once communication with a recipient has been established, a recipient can replenish a sender's tokens directly in a return message, avoiding the need to mint more token batches from the platform. The protocol is



detailed in Figure 4.

**Blinded MACs as one-time use tokens.** We want that a sender can anonymously mint a batch of tokens for a recipient from the platform. The platform should not be able to link the tokens (when they are spent) to the time of minting. To realize this, we again turn to algebraic MACs used by keyed-verification anonymous credentials [20]; we use  $\text{MAC}_{\text{GGM}}$ . Each user generates a MAC secret key  $sk \leftarrow (x_0, x_1) \in \mathbb{Z}_p^2$  and sends it to the platform. A valid MAC on input  $\nu \in \mathbb{Z}_p$  is of the form,

$$t \leftarrow (u_0, u_1 = u_0^{x_0+x_1\nu}) \quad u_0 \leftarrow \mathbb{G}_1.$$

To blindly evaluate a MAC on input  $\nu$ , a user generates a random ElGamal key pair  $(\gamma, D = g_1^\gamma)$  and encrypts  $g_1^\nu$  to  $D$ ,

$$ct = (ct_1 = g_1^r, ct_2 = g_1^\nu D^r) \quad r \leftarrow \mathbb{Z}_p.$$

The user blinds a batch of inputs  $[\nu]_i$  in this manner, creates a group signature  $\sigma$  over  $[ct]_i$  designating the recipient as the opening authority, and then sends  $(\sigma, [ct]_i, D)$  to the platform. The platform verifies the group signature under the recipient's revocation list, and if verification succeeds, proceeds with the blind evaluation using the recipient's MAC secret key. By the homomorphic properties of ElGamal, the platform can maul  $ct$  to form  $ct'$  as an encryption of a valid MAC on  $\nu$  without ever learning anything about  $\nu$ ,

$$ct' = (ct_1^{x_1 \cdot b} g_1^{r'}, ct_2^{x_1 \cdot b} u_0^{x_0} D^{r'}) \quad u_0 \leftarrow g_1^b \quad b, r' \leftarrow \mathbb{Z}_p.$$

The full details of the blind MAC evaluation is given in the full version [59]. The user decrypts  $ct'$  to learn  $u_1$  and stores token  $\tau \leftarrow (\nu, t = (u_0, u_1))$  as the input, tag pair.

To send a message, the user sends the message to the platform along with an unused token  $\tau$  for the recipient. The platform checks that the token  $(\nu, t) \leftarrow \tau$  is unused, i.e.,  $\nu$  is not in the strikelist of used tokens for a recipient, and that the token is valid, i.e., the MAC  $t$  is valid for  $\nu$  under the recipient's MAC key. If those checks pass, the platform delivers the message along with the token  $\tau$  to the recipient and adds  $\nu$  to the recipient's strikelist.

However, the recipient has no way identifying the sender from the token  $\tau$ . The generation of  $\tau$  was (necessarily) blinded to prevent linking by the platform, but that also prevents linking by the recipient.

**Allowing a recipient to link tokens to senders.** Senders must communicate to the recipient the unblinded inputs  $\nu$  for which they are minting tokens. They do this by additionally encrypting the input  $\nu$  to the recipient under the recipient's public key  $Z$ ,

$$\hat{ct} = (\hat{ct}_1 = g_1^{\hat{r}}, \hat{ct}_2 = g_1^\nu Z^{\hat{r}}) \quad \hat{r} \leftarrow \mathbb{Z}_p,$$

and proving in zero knowledge that the input  $\nu$  enclosed in the blinded ciphertext  $ct$  is the same as that enclosed in the ciphertext  $\hat{ct}$  to the recipient. The sender signs the batch of recipient ciphertexts  $[\hat{ct}]_i$  under the group signature with the recipient as the designated opening authority. As before, if the signature  $\sigma$  verifies under the recipient's revocation list,

the platform proceeds with blind evaluation, *but also* sends  $(\sigma, [\hat{ct}]_i)$  to the recipient.

The recipient opens  $\sigma$  to the sender's identity  $upk$ , then decrypts and stores the token identifiers  $[g_1^\nu]_i$ . Later when a recipient receives a message and token  $(\nu, t) \leftarrow \tau$  from the platform, they can link the token to the sender by looking up  $g_1^\nu$ . To block a sender, the recipient generates and sends the revocation token for the sender's  $upk$  to the platform so the sender cannot mint new tokens, as well as sends the sender's remaining unused tokens  $[g_1^\nu]_i$  to add to the strikelist.

**Replenishing tokens directly from the recipient.** The motivation for one-time use tokens was to avoid the cost of the more expensive group signature for every message. However, in some sense, the gain from not running the group signature for every message is offset by the upfront cost of generating a proof to mint each token. While there are optimizations that can be made when batching proofs in this manner [37], this is still an unsatisfying result.

The real efficiency gain from one-time use tokens is when senders can replenish their tokens directly from the recipient, without going through the blind minting process with the platform. Once two users have established sender-anonymous communication, they can use their own secret MAC keys to generate and exchange tokens directly at very little cost.

**Summary.** In this protocol, the core group signature is used only to initiate conversations and mint the first batch of tokens. Once conversation has been established, messages can be exchanged and tokens can be replenished at almost no cost, beyond storage. With regards to storage, users must maintain lists of unused tokens in order to send messages and identify senders of received messages. The platform also needs to maintain an ever-growing strikelist for each user; in practice, users will need to periodically rotate their keys to refresh the platform strikelist, but can ensure that they have distributed tokens for the new key prior to doing so.

Using tokens does leak some information about user communication patterns in a nuanced way. An example might be that if senders need to often mint tokens from the platform for a particular user, the platform can infer that user is not active in responding and replenishing sender tokens.

A second nuance is that in both our scheme and the token strawman [43, 44] presented in Section 3, the message ciphertext of a sender is not bound to the token. The platform can forward the sender's token to the recipient, but swap out the ciphertext, so the recipient will incorrectly attribute it to the sender. In Section 6, we discuss why the impact of such an attack is not large if the underlying E2EE protocol provides message authentication. Nevertheless, we provide a proposal for modifying our token showing protocol to bind the sender's message ciphertext using a BLS signature [13] in the full version [59].

Despite these nuances, we feel Orca with one-time use tokens represents an attractive design choice.

---

## Protocol 2: Orca with One-time Use Tokens

---

### Setup:

- (1) Public parameters for the group signature scheme, algebraic MAC scheme, and public key encryption scheme are generated,  $pp \leftarrow \$ \text{GS.Setup}(\lambda)$ ,  $pp_M \leftarrow \$ \text{MAC.Setup}(\lambda)$ ,  $pp_{PKE} \leftarrow \$ \text{PKE.Setup}(\lambda)$ .
- (2) The platform initializes its state as the group manager of the group signature scheme.
  - (a)  $(gmpk, gmsk) \leftarrow \$ \text{GS.Kg}_{GM}^{pp}()$
  - (b)  $T_U \leftarrow []$ : Table storing user public keys.
  - (c)  $T_R \leftarrow []$ : Table storing user revocation tokens.
  - (d)  $T_k \leftarrow []$ : Table storing user token MAC key and encryption key.
  - (e)  $T_\tau \leftarrow []$ : Table storing strikelist of previously-used tokens for user.

### Registration:

- (1) User generates keys for protocol and initializes recipient state:
  - (a) User registers with platform to acquire group signature signing key with which to send messages,  $\text{GS.JoinU}_U^{pp} \leftrightarrow \text{GS.IssueU}_{GM}^{pp}$ . User stores  $usk$  and platform stores  $upk$ .
  - (b) User registers as opening authority and generates keys with which to block senders,  $\text{GS.JoinOA}_{OA}^{pp} \leftrightarrow \text{GS.IssueOA}_{GM}^{pp}$ .
  - (c) User generates algebraic MAC key used for creating sender tokens,  $(tsk, tpk) \leftarrow \$ \text{MAC.Kg}^{pp_M}()$ , and sends both  $tsk$  and  $tpk$  to platform.
  - (d) User generates keys for public key encryption scheme,  $(ek, dk) \leftarrow \$ \text{PKE.Kg}()$ , stores  $dk$  and sends  $ek$  to platform.
  - (e) User initializes two tables,  $T_x$  and  $T_x^{-1}$ , to identify (and blacklist) senders and their associated sender tokens.
- (2) Platform stores keys and initializes table entries for user:
 
$$T_U[upk] \leftarrow (oapk); T_k[oapk] \leftarrow (tsk, tpk, ek)$$

$$T_R[oapk] \leftarrow []; T_\tau[oapk] \leftarrow []$$

### Sending a message:

- (1) Sender selects unused sender token for recipient and sends message, token, and recipient,  $(m, \tau, oapk)$ , to platform.
- (2) Platform checks if token  $(x, t) \leftarrow \tau$  is valid under recipient's MAC key  $(tsk, tpk, ek) \leftarrow T_k[oapk]$  and if token was not already used (i.e., is not on strikelist).
 
$$b_1 \leftarrow \text{MAC.Ver}^{pp_M}(tsk, x, t)$$

$$b_2 \leftarrow (x \notin T_\tau[oapk])$$
 If  $b_1 = 0$  or  $b_2 = 0$ , platform aborts.
- (3) Platform adds token to strikelist,  $T_\tau[oapk] \leftarrow T_\tau[oapk] \cup \{x\}$ .
- (4) Platform forwards message and token value,  $(m, x)$ , to recipient.
- (5) Recipient removes token from list of valid tokens for sender,
 
$$T_x[T_x^{-1}[x]] \leftarrow T_x[T_x^{-1}[x]] \setminus \{x\}; T_x^{-1}[x] \leftarrow \perp.$$

### Acquiring sender tokens (from platform):

- (1) [Optional] Sender anonymously requests public key information,  $(oapk, tpk, ek)$ , for desired recipient from platform.
- (2) Sender authenticates to platform as a non-blocklisted sender for the recipient using a group signature.
  - (a) Sender signs set of recipient ciphertexts  $[\hat{ct}]_i$  (constructed in (3)) with recipient as opening authority, and sends  $(\sigma, oapk)$  to platform,  $\sigma \leftarrow \$ \text{GS.Sign}_U^{pp}(usk, gmpk, oapk, [\hat{ct}]_i)$ .
  - (b) Platform checks validity of signature against recipient's revocation list,  $b \leftarrow \text{GS.Ver}_{GM}^{pp}(gmsk, oapk, T_R[oapk], [\hat{ct}]_i, \sigma)$ . If  $b = 0$ , then platform aborts.
- (3) Sender engages in token generation protocol with platform.
  - (a) Sender samples  $m$  inputs,  $[x]_i^m \leftarrow \$ \text{MAC.In}(\lambda)^m$ .
  - (b) Sender encrypts inputs to recipient,  $\hat{ct}_i \leftarrow \$ \text{PKE.Enc}(ek, x_i)$ .
  - (c) Sender and platform engage in MAC blind evaluation for each token,  $\text{MAC.BlindInp}^{pp_M}(tpk, x_i) \leftrightarrow \text{MAC.BlindEv}^{pp_M}(tsk)$ , for recipient keys  $(tsk, tpk, ek) \leftarrow T_k[oapk]$ . Sender also sends proof that the input used in the MAC protocol is properly well-encrypted in the ciphertext to the recipient:
 
$$\pi_i \leftarrow \$ \text{NiZK}\{x_i : \text{MAC.BlindInp}^{pp_M}(tpk, x_i) \wedge ct_i = \text{PKE.Enc}^{pp_{PKE}}(ek, x_i)\}$$
 If  $\pi_i$  does not verify, platform aborts the blind MAC protocol.
  - (d) If blind MAC protocol succeeds, sender receives MAC  $t_i$  as output and stores token,  $\tau_i \leftarrow (x_i, t_i)$ .
- (4) Platform sends  $(\sigma, [\hat{ct}]_i^m)$  to recipient.
- (5) Recipient stores tokens to later identify sender.
  - (a) Recipient traces sender,  $upk \leftarrow \text{GS.Open}_{OA}^{pp}(oask, [\hat{ct}]_i, \sigma)$ .
  - (b) Recipient decrypts token ciphertexts and stores tokens.
 
$$x_i \leftarrow \text{PKE.Dec}^{pp_{PKE}}(dk, \hat{ct}_i)$$

$$T_x[upk] \leftarrow T_x[upk] \cup \{x_1, \dots, x_m\}; T_x^{-1}[x_i] \leftarrow upk$$

### Acquiring sender tokens (from recipient):

- (1) Recipient samples  $m$  inputs,  $(x_1, \dots, x_m) \leftarrow \$ \text{MAC.In}(\lambda)^m$ , and MACs them,  $t_i \leftarrow \text{MAC.Ev}^{pp_M}(tsk, x_i)$ .
- (2) Recipient sends tokens  $\tau_i \leftarrow (x_i, t_i)$  to sender associated with  $upk$  out-of-band or via secure channel.
- (3) Recipient stores tokens to later identify sender.
 
$$T_x[upk] \leftarrow T_x[upk] \cup \{x_1, \dots, x_m\}; T_x^{-1}[x_i] \leftarrow upk$$

### Blocklisting a user:

- (1) Recipient looks up sender identity associated with token,  $upk \leftarrow T_x^{-1}[x]$ , and generates revocation token,  $\tau_R \leftarrow \$ \text{GS.Revoke}_{OA}^{pp}(oask, upk)$ . Recipient sends revocation token along with list of remaining sender tokens for sender to platform,  $(x_1, \dots, x_m) \leftarrow T_x[upk]$ .
  - (2) Platform updates blocklist state by adding revocation token to blocklist and remaining tokens to strikelist.
 
$$T_R[oapk] \leftarrow T_R[oapk] \cup \{\tau_R\}$$

$$T_\tau[oapk] \leftarrow T_\tau[oapk] \cup \{x_1, \dots, x_m\}$$
- 

Figure 4: Hybrid protocol based on group signature and tokens.

## 6 Composition with an E2EE Protocol

The main security properties of an E2EE messaging protocol are *message confidentiality* and *message authentication*. Modern forms of message confidentiality include forward secrecy and post-compromise security which ensure that, even in the event of key compromise, previous message content and future message content (after recovery) are not leaked, respectively [23]. Message authentication ensures that messages accepted by the recipient were those encrypted by the sender. A third property, *repudiability*, requires that the authentication mechanism cannot help non-conversation participants verify message authorship, even if secrets from a conversation participant are leaked [15]. For our setting, we will also require the E2EE messaging protocol to be *sender-anonymous*, meaning ciphertexts do not leak any information about the sender, which can be achieved using encapsulation as in sealed sender.

Orca composes with an E2EE messaging protocol to further provide anonymous, outsourced blocklisting (see Section 3). Public keys for Orca may be distributed using the same mechanism used to distribute public keys for the E2EE messaging protocol. Similar to E2EE messaging, to prevent ghost key attacks by a malicious PKI, in which a user's key is replaced by one owned by the adversary, users are expected to perform manual verification of key fingerprints out-of-band or perform periodic auditing of the PKI [50]. Without this assurance, ghost key attacks against Orca result in a break in anonymity, as the adversary can open group signatures using the ghost key. Of course, using Orca does not increase the damage of such attacks: such an adversary can read encrypted messages and break anonymity by subverting the E2EE.

In basic Orca (Figure 3), E2EE ciphertexts are sent along with a group signature over the ciphertext, and when extended with one-time sender tokens (Figure 4), E2EE ciphertexts are sent along with a token produced from a token minting protocol authenticated with a group signature. The composition preserves the message confidentiality and authentication properties of the underlying E2EE protocol: Orca composes generically with the E2EE ciphertexts and does not make further use of the message plaintext. However, Orca necessarily weakens sender-anonymity and repudiability to support blocklisting by a third-party (the platform).

With regards to anonymity, a necessary leakage of the outsourced blocklisting setting is that a ciphertext leaks (to the platform) whether or not the sender is present on the designated recipient's blocklist. Basic Orca meets this minimum leakage, following directly from the anonymity and revocation anonymity security properties of the group signature. Orca extended with one-time tokens leaks more: platform-assisted token minting leaks how many tokens for a recipient are minted, and blocking reveals how many valid tokens remain for the blocked sender. In addition to the anonymity properties of the group signature, achieving only this level of leakage relies on (1) randomly chosen MAC inputs, (2)

security of blind MAC evaluation, (3) confidentiality of the recipient ElGamal ciphertexts, and (4) zero knowledge of the well-formedness proof. We believe it is unlikely the additional leakage of token counts leads to damaging inference attacks, especially considering token counts are further obscured by tokens replenished directly by the recipient.

The minimum weakening to repudiability for the outsourced blocklisting setting is that the platform can at most verify authorship to *some* registered member of the platform, even with compromised secrets. However, our group signature construction does not meet this weakened notion; the platform and the recipient can together provide proof of authorship of a message for a sender. Future work may adapt techniques from deniable signatures (c.f., [58]) to recover repudiability.

Lastly, outsourced blocklisting requires sender attribution: messages delivered to recipients can be correctly attributed to a sender. Basic Orca achieves sender attribution following directly from the traceability security property of the group signature. The extension with one-time tokens achieves sender attribution additionally relying on the soundness of the well-formedness proof of recipient token-tracing ciphertexts.

We also note an optional non-frameability property: a malicious platform should not be able to frame a user as being a sender for a ciphertext they did not create. We do not see this property as security-critical for outsourced blocklisting. A break in non-frameability allows a platform to deliver ciphertexts that are misattributed, however, due to the message authentication property of the underlying E2EE protocol, these ciphertexts will not be accepted by the recipient. The recipient may choose to block the misattributed sender, mistakenly thinking they are spamming malformed ciphertexts. We view this as a special (slightly more damaging) case of a platform-mounted DoS attack, which is not a goal of Orca to defend against. Nevertheless, basic Orca does prevent this attack due to the non-frameability security property of the group signature. Orca with one-time tokens can be extended with token-binding (see full version [59]) to achieve non-frameability relying on the soundness of the blind MAC evaluation proof and the unforgeability of the token-binding signature.

**Formal analyses.** As mentioned, we provide in the full version [59] formal definitions and security analyses for our group signature, the core underlying component of Orca. These analyses do not cover the one-time token extension, nor the security of the composition informally discussed above. Developing formal models suitable for analysis of these higher level primitives remains an open problem. Our initial attempts suggest that this will be challenging, as it seems to require extending existing (already complex) confidentiality and authenticity models for messaging (e.g., [5, 9, 23, 38, 54]) to model sender anonymity, token distribution, blocklist maintenance, etc. An ideal functionality based approach may provide an alternative tack, though any resulting functionality will also be complex (possibly as complicated as our protocols).

## 7 Implementation and Evaluation

This section aims to evaluate the feasibility of deploying Orca at scale. Specifically, we answer the following questions:

- *Client costs*: What are the processing and storage costs that Orca incurs on user clients?
- *Platform costs*: What are the processing and storage costs incurred on the platform? What throughput (user activity) can be reasonably supported given these costs?
- *Bandwidth costs*: How large are Orca protocol messages? What additional networking costs does Orca introduce?

To answer these questions, we provide a prototype library in Rust of our group signature and token-based scheme. Our implementation is over the BLS12-381 pairing-friendly elliptic curve and uses the `zexe/algebra` Rust pairing library [16]. We instantiate the proofs of knowledge using standard Sigma protocols of discrete logarithm relations [17] made non-interactive using the Fiat-Shamir transform [31]. Our security proofs (see full version [59]) rely on a simulation-extractability property of the zero knowledge proofs which has been shown to hold in the algebraic group model [32] for the knowledge of discrete logarithm relation [3, 33]; we believe these techniques can be readily extended to the discrete logarithm relations used in this work. Our implementation consists of less than 1400 lines of code and is available open source<sup>2</sup>.

The experiments, including the microbenchmarks given in Figure 5, were performed using a `c5.12xlarge` Amazon EC2 virtual machine with 24 cores and 96 GB of memory running Ubuntu Server 20.04 LTS as the platform and desktop client (single-core) and on a Google Pixel device running Android 9 as the mobile client. The platform is implemented using an in-memory Redis database for storing revocation blocklists and token strikelists.

When evaluating Orca, recall that users can replenish their token supply directly from the recipient provided there is back and forth communication. Thus, we make the distinction between “initialization costs” of minting an initial token batch from the platform and the “steady-state costs” that occur when tokens are replenished directly from the communicating partner. We expect the majority of user communication to be in steady-state where costs are low.

**Client costs.** Clients must store, for each of their communicating partners, two lists of unused tokens, one for sending messages and one for identifying received messages. These tokens are not large (240B) and the lists can remain small as they can be replenished on next communication. Say a user has 200 communication partners and stores 20 tokens per list. This setup would incur  $\sim 1\text{MB}$  for the client.

The bulk of the processing costs incurred by Orca are concentrated at initialization when a client mints an initial batch

of tokens to start a conversation. On a mobile client, minting an initial batch of tokens takes  $\sim 150\text{ ms}$  for the group signature and an additional  $\sim 100\text{ ms}$  for each token in the batch (see Figure 5). This means it takes around 1 second for a sender to mint 10 tokens. While these costs are significant, we stress that a user only needs to mint enough tokens to initiate a conversation and await a response. If a response from a recipient is delayed, more tokens can be minted as needed. Once a conversation with back-and-forth communication is established, the amortized steady-state cost of sending a message is in creating a new token to replenish the recipient, which is done at very little cost ( $\sim 10\text{ ms}$ ) — approximately the same as sealed sender.

**Platform costs.** The platform stores per-recipient revocation blocklists and token strikelists. The revocation lists are on the order of 100B / revoked user; e.g., a recipient that has blocked 100 users would require a revocation list of size 10KB to be stored. We do not anticipate revocation lists to grow too large, since the platform has other mechanisms to ban users globally (see Section 8). In any case, a platform can impose limits on the size of revocation lists if necessary.

The per-recipient strikelists would grow in size with every message a user sends (32B / spent token). One can use Bloom filters or other data structures to compress the size of the strikelist as well as enforce periodic key rotations to reset its size. If each user sends  $\sim 100$  messages per day and token keys are rotated every two weeks, the platform can store a strikelist of  $\sim 5\text{KB}$  per user with a false positive rate of  $10^{-6}$ . Note the false positive rate can be traded off with storage size; messages that get rejected due to false positives will result in an error returned to the anonymous sender, who may resend with a different token.

The processing costs of the platform are similarly dominated by the token mint requests for initializing conversation as opposed to send requests during steady-state conversation. A request to mint a batch of 10 tokens given a recipient blocklist size of 100 takes  $\sim 200\text{ ms}$  to complete whereas a send request is just a simple algebraic MAC verification and strikelist lookup taking  $< 1\text{ ms}$  (see Figure 5).

Figure 6 demonstrates these workloads are easily parallelizable to achieve high levels of throughput. In this experiment, we run the platform with one million users, each with a blocklist of size 100 and a strikelist of size 1400 (100 messages/day/two weeks), and measure the rate at which the platform can process requests for different levels of hardware parallelism. We do not implement the Bloom filter optimization, so the Redis database stores  $\sim 50\text{KB}$  per user (50GB total), which can still easily fit in memory. The computationally expensive mint requests parallelize with essentially no loss, reaching a rate of 80 requests (for 10 token batches) per second on 24 cores. The inexpensive send requests also parallelize but top out at around 30000 requests per second on 12 cores, which is bottlenecked by the operation throughput of a single Redis database and can be unblocked via a different

<sup>2</sup><https://github.com/nirvantyagi/orca>



Operation		Platform		User (Desktop client)				User (Mobile client)			
				Sender		Recipient		Sender		Recipient	
Sealed sender		–		0.50	(0.02)	0.50	(0.02)	6.6	(0.2)	6.6	(0.2)
Orca	mint tokens with group signature	11.2	(0.2)	10.8	(0.1)	9.7	(0.2)	131.7	(0.8)	117	(2)
	+ cost per token minted	7.60	(0.09)	8.50	(0.08)	0.30	(0.01)	105.2	(0.9)	3.3	(0.1)
	+ cost per blocked user	1.70	(0.04)	–	–	–	–	–	–	–	–
	send message with token*	0.30	(0.01)	0.80	(0.02)	–	–	10.0	(0.2)	–	–

\*Steady-state cost of sending a message with a token that includes cost of replenishing one token

Figure 5: Processing time (ms) microbenchmarks of user and platform operations for Orca compared to sealed sender. Mean time is given with standard deviations shown in parentheses. Dashes indicate an operation that has negligible cost (e.g., a table lookup).

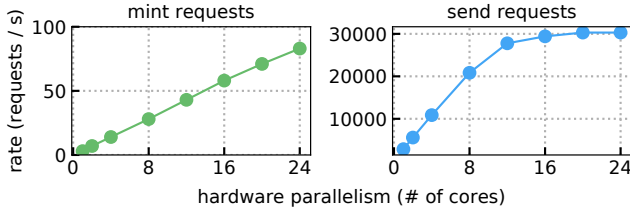


Figure 6: Platform request throughput for different levels of hardware parallelism over a one million user deployment with blocklists of size 100 and strikelists of size 1400. Each mint request corresponds to a request to mint a batch of 10 tokens.

database setup if needed (e.g. through sharding). The achieved bottlenecked throughput already demonstrates feasibility.

**Bandwidth costs.** Minting a token requires sending the group signature (1.6KB) and exchanging proofs for each token to be minted (0.7KB / token). These costs extend to the recipient who receives the signature and also a ciphertext for each token minted (0.2KB / token). Apart from these initialization costs, the steady-state bandwidth costs of sending a message, once again, compare quite favorably with sealed sender. In the steady state, the amortized bandwidth overhead of sending a message would be two tokens (240B / token) — the token being spent *and* the token being created to replenish the recipient. Thus we can achieve amortized per-message overheads of only 30B compared to sealed sender (450B / message).

## 8 Further Extensions

**Backwards unlinkability for revocation tokens.** A drawback of verifier-local revocation is that whenever a new revocation token is provided, the platform can replay the history of messages to link which ones were sent by the newly blocked sender. To prevent such leakage one can take the approach of [51] to rotate revocation keys in set epochs. Naively, this requires recipients to resupply their entire list of revocation tokens; future work may try to incorporate techniques from updatable encryption [12] to provide more efficient epoch transitions.

**Credential expiry and global banning.** Per-recipient blocklists are not a substitute for platform-wide banning of abusive users. The platform must maintain some mechanism for banning accounts in the case of identified user abuse, e.g., through

user reports [28, 35, 58] or account compromise. This can be done by enforcing periodic credential expiration, by for example, rotating the platform’s group manager key. Users must retrieve a new MAC on their public key, at which point, the platform can choose to deny their request.

**Sybil resistance and account recovery.** Outsourced blocklisting works by blocking a public key, not an identity. If malicious users are able to easily send messages under new public keys, either by registering with many accounts or continually rotating an account key after they are blocklisted, then our blocklisting protocol will be of little use. Signal ties accounts to phone numbers to mitigate the ability to easily register new accounts. On the other hand, rotating an account key is a legitimate operation that may need to be taken after account compromise or device loss. Blocking accounts with suspicious key rotation behavior or rate-limiting account recovery are possible mitigations.

**Rate-limited resources.** In Signal, in addition to needing the recipient’s long-lived identity public key, senders also need to pull a one-time use recipient “pre-key” which is used in the initial key agreement protocol to provide forward secrecy properties. Recipients store some number of pre-keys with the platform and replenish them as needed. If a recipient’s pre-keys run out, then conversations are initiated without the pre-key leading to weaker forward secrecy. To prevent malicious users from exhausting a recipient’s pre-key supply, these resources can be protected while preserving anonymous authentication using anonymous rate-limiting techniques [18].

## 9 Related Work

**Anonymous credentials.** Anonymous credentials [19] allow a user to present a cryptographic token proving some specific statement about their identity (e.g., their authorization to send messages to a particular recipient), without revealing anything else about their identity. A problem with anonymous credentials in our setting is that they are — by design — not attributable. While the server processing messages can verify the sender is authorized, the recipient cannot identify the sender. This means there is no way for the server to block the sender in the future, even if some revocation mechanism for the credentials did exist.

A notable design contrast to general-purpose anonymous

credential schemes is Privacy Pass [26], which offers single use credentials that encode only one bit — “I am authorized.” Privacy Pass mints tokens using a verifiable oblivious pseudorandom function [39, 40], which is more efficient than our approach of blind MACs [20], but does not provide the algebraic structure needed to prove relations on the input. We need this property to encrypt the input to the recipient to allow linking of tokens. Blind MACs have been previously suggested for use as one-time tokens [46] and have also been recently proposed as part of Signal’s new proposal for private group messaging [21].

**Anonymous blacklisting.** Anonymous blacklisting [36, 55, 56] systems cover a variety of cryptographic techniques. In general, these systems allow a user to authenticate anonymously to third parties in such a way that the third party can block them from subsequent authentications if they misbehave. In some systems, this blocking ability takes the form of an additional trusted third party that can de-anonymize users much like a group signature. In others, every time a user authenticates they provide a fresh anonymous cryptographic token derived from their identity and a proof that the current blacklist contains no tokens generated by their own keys. Such systems are cryptographically expensive, requiring work linear in the blacklist (for the sender). Moreover, much of the overhead across both settings comes from providing anonymity from the third party. Our setting differs in that the sender need not be anonymous (and in fact, should be identifiable) to the party *adding* to the blacklist (i.e., the recipient), but only be anonymous to the party *filtering* on the blacklist (i.e., the platform).

**Abuse reporting in E2EE messaging.** A complementary line of work [28, 30, 35, 58] considers reporting abusive content sent over an encrypted channel. These systems allow the recipient to verifiably reveal the content of a message to the platform to enable content moderation. They allow attribution of message content to a sender for a known sender identity. They do not allow the attribution of a malformed message with unknown sender as in the griefing attack we describe.

**Metadata-private messaging.** A number of messaging systems have been proposed that provide strong metadata-privacy even against strong network adversaries [4, 6, 24, 25, 42, 45, 47, 53, 57, 60, 61]. These systems incur significant costs on their users, e.g. to send and receive messages at frequent intervals. These costs may dwarf the costs of the types of abuse that Orca aims to prevent. Despite this, a subclass of these systems that could still make use of Orca for blacklisting are based on anonymous broadcasting [24, 41, 47, 53, 61]. Anonymous broadcasts can be converted to a sender-anonymous messaging service by having a messaging service collect, filter, and deliver the broadcast messages with designated recipients.

## 10 Conclusion

This paper explores the tensions between abuse mitigation and sender-anonymity in E2EE messaging. We highlighted several issues with Signal’s sealed sender feature, including weak anonymity set guarantees and vulnerability to griefing attacks.

Our solution, Orca, allows recipients to register privacy-preserving blocklists with the platform. Without learning the sender’s identity, the platform can check that the sender is not on the blocklist and that the recipient will be able to verify their identity. We introduced a new type of group signature tailored to Orca’s needs and propose a hybrid scheme that uses tokens to amortize the bandwidth and computational costs of group signatures.

## References

- [1] Battery Historian. <https://github.com/google/battery-historian>, 2017.
- [2] Michel Abdalla, Mihir Bellare, and Phillip Rogaway. The oracle diffie-hellman assumptions and an analysis of DHIES. In *CT-RSA*, 2001.
- [3] Michel Abdalla, Fabrice Benhamouda, and Philip MacKenzie. Security of the J-PAKE password-authenticated key exchange protocol. In *IEEE S&P*, 2015.
- [4] Nikolaos Alexopoulos, Aggelos Kiayias, Riivo Talviste, and Thomas Zacharias. Mcmix: Anonymous messaging via secure multiparty computation. In *USENIX Security*, 2017.
- [5] Joël Alwen, Sandro Coretti, and Yevgeniy Dodis. The double ratchet: Security notions, proofs, and modularization for the signal protocol. In *EUROCRYPT*, 2019.
- [6] Sebastian Angel and Srinath T. V. Setty. Unobservable communication over fully untrusted infrastructure. In *OSDI*, 2016.
- [7] Mihir Bellare, Daniele Micciancio, and Bogdan Warinschi. Foundations of group signatures: Formal definitions, simplified requirements, and a construction based on general assumptions. In *EUROCRYPT*, 2003.
- [8] Mihir Bellare, Haixia Shi, and Chong Zhang. Foundations of group signatures: The case of dynamic groups. In *CT-RSA*, 2005.
- [9] Mihir Bellare, Asha Camper Singh, Joseph Jaeger, Maya Nyayapati, and Igors Stepanovs. Ratcheted encryption and key exchange: The security of messaging. In *CRYPTO*, 2017.
- [10] Alexandra Boldyreva. Threshold signatures, multisignatures and blind signatures based on the gap-diffie-hellman-group signature scheme. In *PKC*, 2003.
- [11] Dan Boneh, Xavier Boyen, and Hovav Shacham. Short group signatures. In *CRYPTO*, 2004.
- [12] Dan Boneh, Kevin Lewi, Hart William Montgomery, and Ananth Raghunathan. Key homomorphic PRFs and their applications. In *CRYPTO*, 2013.
- [13] Dan Boneh, Ben Lynn, and Hovav Shacham. Short signatures from the weil pairing. In *ASIACRYPT*, 2001.

- [14] Dan Boneh and Hovav Shacham. Group signatures with verifier-local revocation. In *CCS*, 2004.
- [15] Nikita Borisov, Ian Goldberg, and Eric A. Brewer. Off-the-record communication, or, why not to use PGP. In *WPES*, 2004.
- [16] Sean Bowe, Alessandro Chiesa, Matthew Green, Ian Miers, Pratyush Mishra, and Howard Wu. Zexe: Enabling decentralized private computation. In *IEEE S&P*, 2020.
- [17] Jan Camenisch. *Group signature schemes and payment systems based on the discrete logarithm problem*. PhD thesis, ETH Zurich, Zürich, Switzerland, 1998.
- [18] Jan Camenisch, Susan Hohenberger, Markulf Kohlweiss, Anna Lysyanskaya, and Mira Meyerovich. How to win the clonewars: efficient periodic n-times anonymous authentication. In *CCS*, 2006.
- [19] Jan Camenisch and Anna Lysyanskaya. Signature schemes and anonymous credentials from bilinear maps. In *CRYPTO*, 2004.
- [20] Melissa Chase, Sarah Meiklejohn, and Greg Zaverucha. Algebraic MACs and keyed-verification anonymous credentials. In *CCS*, 2014.
- [21] Melissa Chase, Trevor Perrin, and Greg Zaverucha. The signal private group system and anonymous credentials supporting efficient verifiable encryption. In *CCS*, 2020.
- [22] David Chaum and Eugène van Heyst. Group signatures. In *EUROCRYPT*, 1991.
- [23] Katriel Cohn-Gordon, Cas J. F. Cremers, Benjamin Dowling, Luke Garratt, and Douglas Stebila. A formal security analysis of the signal messaging protocol. In *IEEE EuroS&P*, 2017.
- [24] Henry Corrigan-Gibbs, Dan Boneh, and David Mazières. Ripposte: An anonymous messaging system handling millions of users. In *IEEE S&P*, 2015.
- [25] Henry Corrigan-Gibbs and Bryan Ford. Dissent: accountable anonymous group messaging. In *CCS*, 2010.
- [26] Alex Davidson, Ian Goldberg, Nick Sullivan, George Tankersley, and Filippo Valsorda. Privacy pass: Bypassing internet challenges anonymously. *PoPETs*, 2018.
- [27] Roger Dingledine, Nick Mathewson, and Paul F. Syverson. Tor: The second-generation onion router. In *USENIX Security*, 2004.
- [28] Yevgeniy Dodis, Paul Grubbs, Thomas Ristenpart, and Joanne Woodage. Fast message franking: From invisible salamanders to encryption. In *CRYPTO*, 2018.
- [29] Yevgeniy Dodis, Eike Kiltz, Krzysztof Pietrzak, and Daniel Wichs. Message authentication, revisited. In *EUROCRYPT*, 2012.
- [30] Facebook. Messenger Secret Conversations technical whitepaper. <https://fbnewsroomus.files.wordpress.com/2016/07/messenger-secret-conversations-technical-whitepaper.pdf>, 2017.
- [31] Amos Fiat and Adi Shamir. How to prove yourself: Practical solutions to identification and signature problems. In *CRYPTO*, 1986.
- [32] Georg Fuchsbauer, Eike Kiltz, and Julian Loss. The algebraic group model and its applications. In *CRYPTO*, 2018.
- [33] Georg Fuchsbauer, Antoine Plouviez, and Yannick Seurin. Blind schnorr signatures and signed elgamal encryption in the algebraic group model. In *EUROCRYPT*, 2020.
- [34] Jens Groth. Fully anonymous group signatures without random oracles. In *ASIACRYPT*, 2007.
- [35] Paul Grubbs, Jiahui Lu, and Thomas Ristenpart. Message franking via committing authenticated encryption. In *CRYPTO*, 2017.
- [36] Ryan Henry and Ian Goldberg. Formalizing anonymous black-listing systems. In *IEEE S&P*, 2011.
- [37] Ryan Henry and Ian Goldberg. Batch proofs of partial knowledge. In *ACNS*, 2013.
- [38] Joseph Jaeger and Igors Stepanovs. Optimal channel security against fine-grained state compromise: The safety of messaging. In *CRYPTO*, 2018.
- [39] Stanislaw Jarecki, Aggelos Kiayias, and Hugo Krawczyk. Round-optimal password-protected secret sharing and T-PAKE in the password-only model. In *ASIACRYPT*, 2014.
- [40] Stanislaw Jarecki, Aggelos Kiayias, Hugo Krawczyk, and Jia-yu Xu. Highly-efficient and composable password-protected secret sharing (or: How to protect your bitcoin wallet online). In *EuroS&P*, 2016.
- [41] Albert Kwon, Henry Corrigan-Gibbs, Srinivas Devadas, and Bryan Ford. Atom: Horizontally scaling strong anonymity. In *SOSP*, 2017.
- [42] Albert Kwon, David Lu, and Srinivas Devadas. XRD: scalable messaging system with cryptographic privacy. In *NSDI*, 2020.
- [43] Adam Langley. Pond. <https://github.com/agl/pond>, 2016.
- [44] Adam Langley and Trevor Perrin. Replacing group signatures with HMAC in Pond. <https://moderncrypto.org/mail-archive/messaging/2014/000409.html>, 2016.
- [45] David Lazar, Yossi Gilad, and Nickolai Zeldovich. Karaoke: Distributed private messaging immune to passive traffic analysis. In *OSDI*, 2018.
- [46] Isis Agora Lovecruft and Henry de Valence. HYPHAE: Social secret sharing. <https://patternsinthevoid.net/hyphae/hyphae.pdf>, 2017.
- [47] Donghang Lu, Thomas Yurek, Samarth Kulshreshtha, Rahul Govind, Aniket Kate, and Andrew K. Miller. Honeybadgermpc and asynchromix: Practical asynchronous MPC and its application to anonymous communication. In *CCS*, 2019.
- [48] Joshua Lund. Technology preview: Sealed sender for Signal. <https://signal.org/blog/sealed-sender/>, 2017.
- [49] Ian Martiny, Gabriel Kaptchuk, Adam Aviv, Dan Roche, and Eric Wustrow. Improving Signal’s sealed sender. In *NDSS*, 2021.
- [50] Marcela S. Melara, Aaron Blankstein, Joseph Bonneau, Edward W. Felten, and Michael J. Freedman. CONIKS: bringing key transparency to end users. In *USENIX Security Symposium*, 2015.

- [51] Toru Nakanishi and Nobuo Funabiki. Verifier-local revocation group signature schemes with backward unlinkability from bilinear maps. In *ASIACRYPT*, 2005.
- [52] Trevor Perrin and Moxie Marlinspike. The double ratchet algorithm. <https://signal.org/docs/specifications/doubleratchet/>, 2016.
- [53] Ania M. Piotrowska, Jamie Hayes, Tariq Elahi, Sebastian Meiser, and George Danezis. The loopix anonymity system. In *USENIX Security*, 2017.
- [54] Bertram Poettering and Paul Rösler. Towards bidirectional ratcheted key exchange. In *CRYPTO*, 2018.
- [55] Patrick P. Tsang, Man Ho Au, Apu Kapadia, and Sean W. Smith. Blacklistable anonymous credentials: blocking misbehaving users without ttps. In *CCS*, 2007.
- [56] Patrick P. Tsang, Apu Kapadia, Cory Cornelius, and Sean W. Smith. Nymble: Blocking misbehaving users in anonymizing networks. *IEEE Trans. Dependable Sec. Comput.*, 2011.
- [57] Nirvan Tyagi, Yossi Gilad, Derek Leung, Matei Zaharia, and Nickolai Zeldovich. Stadium: A distributed metadata-private messaging system. In *SOSP*, 2017.
- [58] Nirvan Tyagi, Paul Grubbs, Julia Len, Ian Miers, and Thomas Ristenpart. Asymmetric message franking: Content moderation for metadata-private end-to-end encryption. In *CRYPTO*, 2019.
- [59] Nirvan Tyagi, Julia Len, Ian Miers, and Thomas Ristenpart. Orca: Blocklisting in sender-anonymous messaging. *IACR Cryptology ePrint Archive*, 2021.
- [60] Jelle van den Hooff, David Lazar, Matei Zaharia, and Nickolai Zeldovich. Vuvuzela: scalable private messaging resistant to traffic analysis. In *SOSP*, 2015.
- [61] David Isaac Wolinsky, Henry Corrigan-Gibbs, Bryan Ford, and Aaron Johnson. Dissent in numbers: Making strong anonymity scale. In *OSDI*, 2012.



## A Griefing Attack on Sealed Sender

We identify and implement a griefing attack against Signal’s sealed sender protocol. An attacker in possession of a recipient’s access key can spam the recipient with untraceable messages, causing the recipient’s system to suffer battery drain and to consume bandwidth.

**Attack vector.** The attack takes advantage of the fact that the platform cannot check for malformed sealed messages. Our proof-of-concept attack simply modifies the Signal client to modify the triple  $(pk_e, ct_{id}, ct_{ss})$  by replacing  $pk_e$  with a new, random value  $pk_f$ . To the platform this is indistinguishable from a legitimate sealed sender message, but the recipient’s decryption will fail when trying to decrypt  $ct_{id}$  and cannot recover any information about the sender. Our modification causes the recipient’s decryption to fail early. While technically one could force the recipient to perform more cryptographic steps, this would have small impact on the efficacy of the attack.

This approach required changing only two lines of code in the Signal Desktop client. We also wrote a small script to automate sending messages via the client.

**Attack efficacy.** We performed some measurements to assess whether the griefing attack can be used, particularly, to drain a target’s battery. In our experiments, we used as attacker our modified Signal Desktop application on a MacBook Pro 2017 machine running macOS Mojave using a 2.5 GHz Intel Core i7. We used as a stand-in for victim recipient an unmodified Signal Android application (version 4.54.3) on a Google Pixel phone running Android version 9. We used the Android Battery Historian tool [1] to inspect the effect of our attack on battery drainage. It reports the battery level rounded to the nearest percent.

In our experiments we only interacted with the Signal platform and with researcher devices. We purposefully experimented only with very low volume attacks in order to ensure we did not burden the Signal platform, and confirmed ahead of time with members of the Signal team that our experiments would not be problematic. In summary, the platform and its users were not negatively affected by our experiments.

We measured the rate of change in battery level per hour when sending one malformed sealed message every 1, 2, 5, or 10 seconds. As a baseline comparison, we also measured the rate of battery drainage when no messages were sent. Each of the four sending rates were measured over a period of 2 hours, while the baseline was measured over a period of 11 hours; the phone discharges slowly at rest so an extended measurement period was needed for the baseline. Before each experiment, the recipient phone was rebooted and charged to full capacity. During each experiment, the phone used its mobile data for network connectivity and was otherwise idle.

In the baseline case, where the phone received no malicious messages, the battery level dropped by only 0.45 levels per hour (dropping the battery by only 7% in 11 hours). In com-

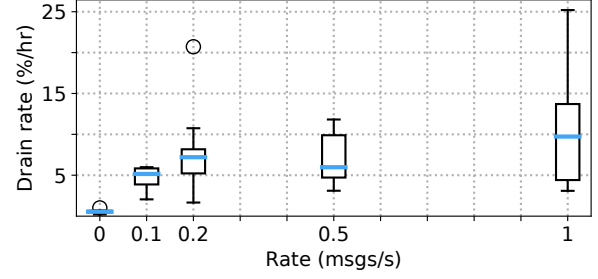


Figure 7: Battery drain rate of griefing attack for various rates of sending,  $x \in \{0, 0.1, 0.2, 0.5, 1\}$  / second. The box plot shows the variability of drain rates over trials, with the range, quartiles and median denoted by the whiskers, box, and line, respectively (outliers marked separately).

parison, the drop rates were 4.11, 5.37, 5.84, and 6.88 levels per hour when sending a message once every 10s, 5s, 2s, and 1s. Thus even the slowest attack rate speeds up battery drain by 9x; for one message a second it is 15x. We show a boxplot of these measurements in Figure 7.

The attack also consumes recipient bandwidth (which could be costly if they pay for data service per byte): at one message per second, the Signal Android application received 1.13 MB/hour, while as a baseline it receives 0.94 KB/hour.

A real attacker can of course trivially increase attack volume up to any general rate limiting enforced by the platform. While it is not public if Signal rate limits clients (and we did not want to stress test it), we believe even modest increases to the volume will allow draining batteries quickly. While battery drain rates will vary significantly based on target handset and other factors, we believe our proof-of-concept evidences sufficient impact on a victim to be a concern.

**Mitigation options for victims.** The receiver’s Signal client gives no obvious visual indication that messages are being received and filtered. To learn of message filtering, a user would have to inspect the client’s debug logs, making the attack essentially invisible for the majority of users. Even if detected, there are no particularly good ways to prevent the griefing attack.

The victim can rotate their access key  $ak$  and attempt to redistribute a new  $ak'$  to their communicating partners. If the attacker is not able to get access to the new access key, the attack will be stopped by the platform and no messages will reach the victim’s client. But since the attack leaves no information about which of the victim’s communicating partners is responsible, the victim can only make a guess as to whom they should block.

This issue might lead people to only add a few, highly trusted contacts. But this degrades anonymity significantly, since as discussed in Section 2.2, the platform knows that a sealed sender must be one of the recipient’s contacts.

## B Artifact Appendix

### B.1 Abstract

Our artifact contains source files of the Orca blocklisting protocol as a library in Rust. The cryptographic protocol is built on top of the open-source `arkworks` library for pairing-based cryptography. The implementation consists of three major parts: (1) an implementation of the Chase et al. algebraic MAC protocol, (2) an implementation of the Orca group signature, and (3) an implementation of the Orca one-time-use token protocol. The artifact also includes two benchmarks for reproducing the performance numbers reported on. These benchmarks can be easily run on any machine that can compile Rust from source, though we report performance numbers from running on high-memory AWS machines (for the server) and mobile devices (for the client). The artifact does not include source files for the griefing attack and battery-drain experiments against Signal, as they are potentially harmful and are not core to our work’s claimed contribution.

### B.2 Artifact check-list (meta-information)

- **Algorithm:** The Orca blocklisting protocol including group signature and one-time-use tokens.
- **Compilation:** Benchmarks are built from source using the Rust compiler.
- **Run-time environment:** Our artifact was run on a `c5.12xlarge` AWS EC2 virtual machine with 24 cores and 96 GB of memory running Ubuntu Server 20.04 LTS, as well as on a mobile device running Android 9.
- **Hardware:** The mobile microbenchmarks were run on a Google Pixel 2 device. The server throughput benchmark requires at least 64 GB, though comparable results can be reproduced with less memory.
- **Execution:** The microbenchmarks run in less than 5 minutes. The server throughput benchmark runs in under 2 hours on our test AWS machine.
- **Security, privacy, and ethical concerns:** We do not provide the source files for the griefing attack and battery-draining experiments.
- **Output:** The benchmarks produce summarized performance outputs printed to the terminal.
- **Experiments:** There are two benchmarks: (1) microbenchmarks for measuring the performance of the cryptographic primitives used in Orca, and (2) macrobenchmark for measuring server throughput of requests.
- **How much time is needed to prepare workflow (approximately)?:** The benchmark binaries are built from source in under 5 minutes. Setting up the AWS machine and/or the mobile device may take additional time.
- **Publicly available?:** The latest version of the library is available at <https://github.com/nirvantyagi/orca>.

### B.3 Installation

The setup consists of installing Rust and compiling the benchmark binaries from source. Compiling and running the microbenchmarks on a mobile device requires additional installation of the Android Native Development Kit (NDK) and related Rust toolchains. The macrobenchmark for server throughput additionally requires installing and running a Redis server locally. Detailed installation instructions are given on the README available at <https://github.com/nirvantyagi/orca>.

### B.4 Evaluation and Expected Results

There are two benchmark binaries that we report results on. The first is the microbenchmarks binary that is used to populate Figure 5. The platform and desktop client user columns are given from running the microbenchmark binary on a single core of the specified AWS machine. The mobile client user column is given from running the microbenchmark on the specified mobile device.

The second benchmark binary measures server throughput and is used to populate Figure 6. The reported numbers are based on experiments setting benchmark parameters of 200 requests for a blocklist size of 100, a strikelist size of 1400, and one million users, while varying the number of cores. This setup requires 64 GB of memory, however, the number of users can be reduced (e.g., to 200) to reproduce similar results without large memory requirements.

Detailed evaluation instructions are given on the README available at <https://github.com/nirvantyagi/orca>.

### B.5 Experiment Customization

The benchmark source code is available and can be customized beyond the existing parameterization.

### B.6 Notes

The cryptographic code has not been reviewed; it serves as a research prototype and is not suitable for deployment. If any bugs are discovered, please raise an issue on Github or send an email to the authors.