

# Uncontrolled Randomness in Blockchains: Covert Bulletin Board for Illicit Activity

Nasser Alsalami and Bingsheng Zhang

Lancaster University, UK  
{n.alsalami, b.zhang2}@lancaster.ac.uk

**Abstract.** Public blockchains can be abused to covertly store and disseminate potentially harmful digital content. Consequently, this threat jeopardizes the future of such applications and poses a serious regulatory issue. In this work, we show the severity of the problem by demonstrating that blockchains can be exploited as a covert bulletin board to secretly store and distribute arbitrary content. More specifically, all major blockchain systems use randomized cryptographic primitives, such as digital signatures and non-interactive zero-knowledge proofs, and we illustrate how the uncontrolled randomness in such primitives can be maliciously manipulated to enable covert communication and hidden persistent storage. To clarify the potential risk, we design, implement and evaluate our technique against the widely-used ECDSA signature scheme, the CryptoNote's ring signature scheme, and Monero's ring confidential transactions. Importantly, the significance of the demonstrated attacks stems from their undetectability, their adverse effect on the future of decentralized blockchains, and their serious repercussions on users' privacy and crypto funds. Finally, besides presenting the attacks, we examine existing countermeasures and devise two new steganography-resistant blockchain architectures to practically thwart this threat in the context of blockchains.

**Keywords:** Blockchain, Steganography, Covert Broadcast Channels, Content Insertion, Wallet Subversion

# Table of Contents

1	Introduction	4
1.1	Paper Roadmap	8
2	Preliminaries	8
2.1	Notations	8
2.2	Blockchain	8
2.3	(Ring) Signature schemes	9
2.4	Brief description of CryptoNote	9
2.5	Brief description of Monero (Version 0.12.0.0)	10
2.6	Steganography	12
2.7	Kleptography/Algorithm-substitution attacks	13
2.8	ECDSA	14
3	Generic Steganographic Attack	14
3.1	Our generic steganographic attack on CryptoNote	15
3.2	Security	17
3.3	Robustness and Efficiency	18
4	Case studies: Bytecoin and Monero	19
4.1	Implementation in Bytecoin	19
4.2	Implementation for Monero (version 0.12.0.0)	22
5	Attack Scenarios	23
5.1	Attack Scenario 1: Covert Broadcast Channel	23
5.2	Attack Scenario 2: Covert Data Storage and Distribution	24
5.3	Attack Scenario 3: Wallet Subversion	24
6	Countermeasures	27
6.1	Existing Countermeasures	28
6.2	Stego-Resistant Blockchain Framework (SRBF)	30
6.3	Randomness-Externalized Architecture	31
7	Related Work	34
8	Conclusion	35
	References	36
	Appendices	40
A	Security Proofs	40
A.1	Proof of Theorem 1	40
A.2	Proof of Theorem 2	40
B	Bytecoin Psuedo-code and Demo Transaction	41
C	Detailed Implementation of Steganographic Attack in Monero	41
D	ECDSA-Signature Rejection-Sampling Experiment	43
E	Signature subversion	44
F	Ciphertext Stealing Technique	45
G	Summary of Existing Countermeasures	46
G.1	ASA-Resistant Techniques	46
G.2	Preventative Trust-Based Countermeasures	47

Uncontrolled Randomness in Blockchains	3
G.3 Blockchain-based techniques	49
H Non-interactive Zero-knowledge	49

## 1 Introduction

The blockchain technology has pioneered a new paradigm to realize large-scale immutable, persistent, and append-only distributed ledgers. Nowadays, blockchain-powered systems have become largely ubiquitous across various sectors including technology, academia, medicine, economics, finance, etc. While the blockchain technology is promising in a great number of application scenarios, it can also be abused to anonymously store and disseminate potentially harmful digital content. A recent study [1] has shown that 1.4% of all Bitcoin transactions contain non-financial data, some of which contain objectionable content, e.g. links to child pornography. Though the absence of a central censor makes blockchains appealing in some use cases, the increasing amount of illicit content posted to the blockchains poses a serious regulatory issue [2]. Subsequently, several techniques have been discussed to either filter unwanted content before it is added to the ledger [3] or remove content from the blockchain [4, 5].

However, all of the proposed countermeasures can only be effective if the malicious content attached to the transactions can be detected. The situation gets worse when the attackers hide data into normal transactions and use blockchain platforms for covert communications. Naively, one can encrypt the malicious content and attach its ciphertext to a transaction, but it is noticeable to the public that there is suspicious data attached. In 2018, Partala [6] showed a proof-of-concept steganography technique that allows an adversary to covertly embed one bit into a standard Bitcoin transaction’s recipient address without being distinguished from an innocuous transaction and without burning the funds.

In this work, we further advance this line of research by demonstrating an effective steganographic method that offers high throughput and can be launched against any blockchain platforms that use randomized, i.e. probabilistic, cryptographic primitives, such as digital signatures and non-interactive zero-knowledge proofs. The main observation is that all randomized cryptographic algorithms need to consume random coins somewhere along the execution, and these random coins are not audited or certified publicly. By intentionally manipulating the random coin supplied to a randomized algorithm, an attacker is able to embed arbitrary information into the outputs of the algorithm. The output that contains steganographic data is computationally indistinguishable from normal output.

Besides using the demonstrated attack for covert channels and persistent storage, the same attack is applicable in another scenario. The attacker(s) may try to subvert, or mis-implement, cryptocurrency wallets and re-distribute them to unsuspecting users. The subverted wallets can then surreptitiously leak the victim’s secret, such as the signing key, via standard transactions. Importantly, the transactions generated by the subverted wallets are computationally indistinguishable from normal transactions for any *black-box* observer.

Moreover, the current focus of research regarding blockchain subversion vulnerabilities is mainly on the trusted parameter setup process, such as common reference string (CRS) generation [7, 8], while software subversion vulnerabilities in blockchain cryptocurrency applications has not been extensively studied.

The plausibility of algorithm-substitution attacks against cryptocurrency can be attributed to the following three reasons. Firstly, cryptocurrencies have very complex cryptographic primitives and structures that makes them prone to unseen mis-implementations. An example of such mis-implementations is shown in [9] where the Tencent’s QQ browser is said to have used textbook RSA algorithm with no padding, which is well-known to be insecure as it is a deterministic encryption scheme. This is further demonstrated in [10] which notes that over 1/3 of the open-source smart contracts contain at least one bug, and some of them are maliciously embedded and can be triggered later by the attackers in a similar manner to the infamous Ethereum DAO hack [11] (\$ 55 million).

Secondly, although many cryptocurrencies are marketed as decentralized projects, studies have found that the development of many blockchain applications is highly centralized. For example, 30% of the source files in Bitcoin are written by a single author, and 7% of the code is written by the same author [12]. Similarly, 20% of the source code in Ethereum is attributed to the same author [12]. This high centralization may cause bias and introduce intentional and unintentional flaws. Thirdly, most end users lack the ability and the means to check the conformity of an executable wallet with its reference source code. In fact, in some platforms, such as iOS, users can not directly access the binary files without jailbreaking their devices, which paradoxically is not advisable and may render a device unsafe to run a cryptocurrency wallet. Besides, it is uncommon for users to compile the source code of any application by themselves; instead, they usually rely on downloading readily prepared executable applications. The difficulty to examine the implementation of a cryptocurrency wallet is even more pertinent to hardware wallets, such as the various *Swiss-Army-Knife* hardware wallets [13]. These hardware wallets are typically manufactured in an outsourced loosely-controlled environment, and it is practically impossible to audit the integrity of their implementation through the standard functionality ‘correctness’ test by observing input/output pairs in a black-box manner.

**Our contributions.** The primary objective of this work is to draw attention to the potential threat of abusing uncontrolled randomness in blockchain algorithms, and attempt to devise practical countermeasures. To the best of our knowledge, this work is the first in literature that discusses such a widely spread vulnerability in the blockchain context. More specifically, we summarize our contributions as follows:

- **Novel blockchain steganographic technique.** We propose a steganographic technique that greatly increases the throughput of the state-of-the-art blockchain steganographic attack that affects many cryptocurrencies. We present our general attack against the widely-used CryptoNote framework, and as a demonstration, we design, implement and evaluate the attack on Monero and Bytecoin currencies.
- **Covert broadcast channels.** As an immediate application, we show blockchain platforms can be exploited to act as covert broadcast channels. Once deployed, this would be the world’s first practical covert broadcast channel. The exis-

tence of such a channel will be untraceable, unlinkable, and even unobservable. Such broadcast channels could be disastrous if used by outlaws, e.g. terrorists.

- **Persistent storage.** With the proposed steganographic technique, anyone can use the blockchain as a cheap hidden persistent storage along with their daily transactions. For instance, this can be used for uncensorable cyberlockers. At the time of submission, persistently storing 1 GB of data on Bytecoin blockchain and using its P2P network as CDN costs less than \$ 3. In theory, data storage is just a communication channel between the current user and the user himself in the future. Nevertheless, there is a subtle difference between hidden storage and covert channels, that is how long the channel (data) would exist. Also some countermeasures are effective against persistent storage but not against covert channels.
- **Wallet subversion attacks.** For the first time, we point out that there is a troubling high risk of massive coin stealing among all of the current cryptocurrency wallets by demonstrating efficient and effective subversion attacks within the realm of *Kleptography* and *Algorithm Substitution Attacks*. This attack possesses the following properties:
  - *Passive attack.* After the victim user downloads and installs the subverted wallet, the attacker does not need to interact directly with the victim’s wallet. The communication channel between the subverted wallets and the attacker is simply through the transactions posted on the blockchain.
  - *Undetectability in black-box setting.* The transactions generated by compromised wallets are computationally indistinguishable from the honestly-generated transactions. Therefore, no online/offline *watchdog* can detect the subversion. It is important to note that *undetectability in this context does not mean the in-ability to detect source code discrepancies between genuine and subverted wallets, or the in-ability to reverse-engineer the subverted wallet*, but rather the computational indistinguishability between transactions generated by the two types of wallets.
  - *Interoperability.* The subverted wallets transact seamlessly with normal wallets; i.e. they can send to and receive from other wallets regardless whether other wallets are subverted or not.
  - *Subtlety.* In accordance with the definition of *kleptography and Algorithm-Substitution Attacks*, shown in Sec. 2, we consider our attack *exclusively* in the black-box setting. However, if optimized, the difference between a subverted wallet source code, e.g. Bytecoin wallet, and the original code is only about ten lines of code in two functions. This subtlety makes it difficult even for technology-savvy users to review and detect the subversion even if the subverted wallet is open source.

We have implemented our subversion attacks against the ECDSA signature scheme, and the ring signature used in the CryptoNote framework which is implemented by many cryptocurrencies, such as Bytecoin. Because ECDSA and ring signature are widely used among cryptocurrencies, this work has direct impact on 18 of the top 25 cryptocurrencies in terms of market capitalization [14] (as of the time of writing) as depicted in Table 1.

Cryptocurrencies' Signatures					
#	Cryptocurrency	ECDSA	EdDSA	Ring Signature	Note
1	Bitcoin	✓			
2	Ethereum	✓			
3	Ripple	✓	✓		
4	Bitcoin Cash	✓			
5	Litecoin	✓			
6	Cardanos		✓		
7	Stellar		✓		
8	Zcash		✓		
9	IOTA				Winternitz
10	Monero			✓	
11	Dash	✓			
12	NEM		✓		
13	Ethereum Classic	✓			
14	Komodo	✓			
15	Verge	✓			
16	Lisk		✓		
17	Dogecoin	✓			
18	Decred	✓	✓		
19	Nano		✓		
20	Wanchain	✓		✓	
21	Bytecoin			✓	
22	Siacoin		✓		
23	Bitcoin Diamond	✓			
24	BitShares	✓			
25	Waves		✓		

**Table 1.** Cryptocurrencies and Digital Signature Schemes (currencies checked with either the ECDSA signature or the Ring signature are potentially susceptible to the our wallet subversion attacks.)

- **Countermeasures.** We provide two practical and effective countermeasures to prevent the stego-use of the cryptographic components in a transaction. In particular, we propose a stego-resistant blockchain framework (SRBF) that can be readily applied to all of the off-the-shelf blockchain systems. It is a generic solution that is suitable to any blockchain with randomized signature schemes. More specifically, in this proposal, the miners are trusted. Upon receiving a (randomized) signature associated with a transaction, instead of directly including it into the transactions in the next block, the miner replaces the signature with a non-interactive zero-knowledge proof showing that the miner has seen a valid signature for the transaction. Subsequently, the miner drops the possibly stego-generated signature and puts the transaction together with the miner’s proof on the next block. As a long-term solution to malicious substitution attacks, we also propose randomness-externalized architecture as to enable implementation auditing.

## 1.1 Paper Roadmap

The rest of this document is organized as follows: Sec. 2 provides background and definitions, and explains some preliminary concepts. In Sec. 3, we illustrate our generic steganographic attack against CryptoNote-based cryptocurrencies, assess its effectiveness, and prove its security. After that, we demonstrate our implementation of the generic steganographic attack in Bytecoin and Monero in Sec. 4. In Sec. 5 we explore three different scenarios in which our generic attack could be applied. Besides, Sec. 5.3 presents two more subversion attacks on ECDSA-signature wallets. Also, we discuss the existing countermeasures and suggest two new techniques in Sec. 6. In addition, we present the related work in Sec. 7. Finally, Sec. 8 restates the main objectives and findings, concludes the document, and explains potential future work.

## 2 Preliminaries

Below we describe the necessary notations to used in this document, and provide description of some preliminary concepts that are related to this work.

### 2.1 Notations

We use the following notations throughout this paper. The notation  $[n]$  stands for the set  $\{1, 2, \dots, n\}$ . For a randomized algorithm  $A()$ , we write  $y = A(x; r)$  to denote the unique output of  $A$  on input  $x$  and randomness  $r$ , and write  $y \leftarrow A(x)$  to denote the process of picking randomness  $r$  uniformly at random and setting  $y = A(x; r)$ . We use  $s \stackrel{\$}{\leftarrow} S$  to denote sampling an element  $s$  uniformly at random from a set  $S$ . We use  $\lambda \in \mathbb{N}$  as the security parameter. Let  $\text{poly}(\cdot)$  denote a polynomially-bounded function and  $\text{negl}(\cdot)$  denote a negligible function. Unless specified in the context, we use  $\text{hash}_p : \{0, 1\}^* \mapsto \mathbb{Z}_p$  and  $\text{hash}_g : \{0, 1\}^* \mapsto \mathbb{G}$  as two collision resistant hash functions that map an arbitrary length string to a group element in  $\mathbb{Z}_p$  and  $\mathbb{G}$ , respectively.  $m_{[a:b]}$  stands for the truncation that contains from the  $a$ -th bit to the  $b$ -th bit of  $m$ .

### 2.2 Blockchain

The term *blockchain* encompasses a broader range of distributed ledger technologies initiated by Bitcoin [15]. There are two types of blockchains; permissioned (private) and permissionless (public). In this work, we mainly focus on permissionless blockchains. Typically, a permissionless blockchain uses a *Proof-of-X* mechanism, such as Proof-of-Work (PoW) and Proof-of-Stake (PoS), to randomly nominate a node which will propose the next block. The valid transactions contained in a block need to be signed by the owner(s) of the corresponding consumed coins. Most blockchain systems use randomized signature algorithms, which makes them vulnerable to our attacks.



### 2.3 (Ring) Signature schemes

For notation simplicity, we unify the syntax of signature schemes and ring signature schemes. A (ring) signature scheme consists of a tuple of algorithms  $\mathcal{S} = (\text{Setup}, \text{KeyGen}, \text{Sign}, \text{Verify})$  as follows:

- $\text{param} \leftarrow \text{Setup}(1^\lambda)$  is the setup algorithm that takes as input the security parameter  $1^\lambda$ , and it outputs a system parameter  $\text{param}$ . The rest of the algorithms implicitly take  $\text{param}$  as an input.
- $(\text{PK}, \text{SK}) \leftarrow \text{KeyGen}(\text{param})$  is the key generation algorithm that takes as input the setup parameter  $\text{param}$ , and outputs a pair of public and secret keys  $(\text{PK}, \text{SK})$ .
- $\sigma \leftarrow \text{Sign}(\mathcal{P}, \text{SK}, \ell, m)$  is the signing algorithm that takes as input a set of public keys  $\mathcal{P} := \{\text{PK}_1, \dots, \text{PK}_n\}$ , the secret key  $\text{SK}$ , the index  $\ell$  such that  $\text{SK}$  is the corresponding secret key of  $\text{PK}_\ell$ , and the message  $m$ , and it outputs the signature  $\sigma$ . (For a standard signature scheme, we have  $|\mathcal{P}| = 1$  and  $\ell = 1$ .)
- $b \leftarrow \text{Verify}(\mathcal{P}, m, \sigma)$  is the verification algorithm that takes as input a set of public keys  $\mathcal{P}$ , the message  $m$  and the signature  $\sigma$ , and it outputs  $b := 1$  if only if the signature is valid.

**Signature Unforgeability.** In a blockchain application, a signature scheme needs to achieve existential unforgeability under an adaptive chosen-message attack (EUF-CMA). While there are various unforgeability definitions for ring signature schemes; in this work, we adopt the most commonly used unforgeability against fixed-ring attacks and unify it with EUF-CMA. We refer interested readers to [16] for more ring signature security definition variants and their differences.

**Definition 1.** We say a (ring) signature scheme  $\mathcal{S} = (\text{Setup}, \text{KeyGen}, \text{Sign}, \text{Verify})$  is EUF-CMA if for any PPT adversary  $\mathcal{A}$ , any integer  $\lambda \in \mathbb{N}$ , any  $n = \text{poly}(\lambda)$ , any  $\text{param} \leftarrow \text{Setup}(1^\lambda)$ , any  $\{(\text{PK}_i, \text{SK}_i)\}_{i=1}^n$  output by  $\text{KeyGen}(\text{param})$ , we have:

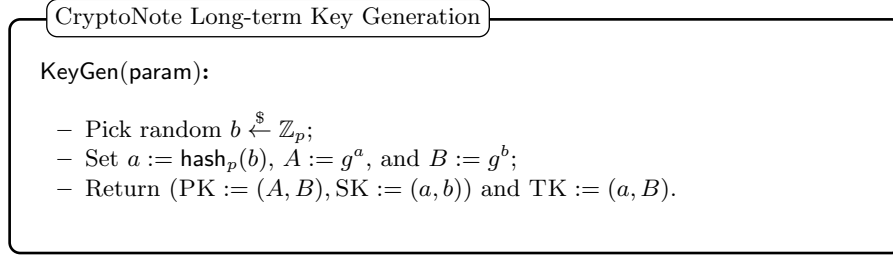
$$\Pr \left[ \begin{array}{l} (m^*, \sigma^*) \leftarrow \mathcal{A}^{\mathcal{O}(\cdot, \cdot)}(\{\text{PK}_i\}_{i=1}^n) : \\ \text{Verify}(\{\text{PK}_i\}_{i=1}^n, m^*, \sigma^*) = 1 \wedge m^* \notin \mathcal{Q} \end{array} \right] = \text{negl}(\lambda)$$

where  $\mathcal{O}(s, m) := \text{Sign}(\{\text{PK}_i\}_{i=1}^n, \text{SK}_s, s, m)$  be the signing oracle, and  $\mathcal{Q} := \{m_1, \dots, m_q\}$  is the set of queries to the signing oracle  $\mathcal{O}(\cdot, \cdot)$ .

### 2.4 Brief description of CryptoNote

CryptoNote is a protocol proposed by Nicolas van Saberhagen [17], and it has been implemented in many emerging cryptocurrencies, such as Bytecoin [18], CryptoNoteCoin [19], Fantomcoin [20], etc. Compared to Bitcoin-like cryptocurrencies, CryptoNote offers two main features: (i) *stealth address via non-interactive key exchange* and (ii) *set anonymity via (linkable) ring signatures*.

More specifically, the user's private key consists of  $a, b \in \mathbb{Z}_p$ , and the corresponding public key  $(A, B)$  consists of  $A := g^a$  and  $B := g^b$ . Note that in a



**Fig. 1.** CryptoNote Long-term Key Generation Algorithm.

standard CryptoNote implementation,  $a$  is usually defined as  $\text{hash}_p(b)$ ; therefore,  $b$  is the actual secret key. In CryptoNote, to transfer funds to a recipient, the payer needs to generate a transaction public key  $R := g^r$  and compute the corresponding one-time address  $T := (g^{\text{hash}_p(A^r)} \cdot B)$ . The recipient is then able to compute the corresponding one-time private key as  $t := (\text{hash}_p(R^a) + b)$ . By the property of Diffie-Hellman exchange, we have  $A^r = R^a$ . With regards to the one-time ring signature schemes, it is transformed from the OR-composition of Schnorr’s identification Sigma protocols. There exists a LNK algorithm that can link two signatures together if they are produced by the same signing key. By design, the one-time signature key can only be used once, and it can be detected if the same key is used to sign two transactions, which prevents double spending. More specifically, let  $T := g^t$  be the one-time public key, and define  $I := (\text{hash}_g(T))^t$  as a “key image” as part of the signature. The ring signatures signed by the same secret key would have identical key image; therefore, double spending can be defeated efficiently by simply checking if the key image has already been used.

Let  $\mathcal{P} := \{P_i\}_{i=1}^n$  be a set of public keys, and the signer knows the secret key  $t_\ell$  such that  $P_\ell = g^{t_\ell}$ ,  $\ell \in [n]$ . Denote  $I := \text{hash}_g(P_\ell)$  as the key image. The **param** is defined as the parameters of the ED25519 twisted Edwards curve. For completeness, we provide the key generation and signing algorithms in Fig. 1 and Fig. 2, respectively. Fig. 1 shows a third key called the *tracking key* TK that can be given to a third party to track all transactions destined to the owner of this key without revealing their secret key SK.

## 2.5 Brief description of Monero (Version 0.12.0.0)

Monero [21] is one of the most successful CryptoNote-based cryptocurrencies, and its source code is available on GitHub [22]. Although the original Monero was based on the CryptoNote protocol, its transaction signature has evolved beyond this protocol<sup>1</sup>. As mentioned in [24], CryptoNote suffers from a shortcoming where amounts in transactions are not hidden. To address this issue,

<sup>1</sup> The Monero project is very active and evolves rapidly. In fact they have two major releases each year. In Oct. 2018, Monero released version 0.13.0.0 “Beryllium Bullet”, which switched to Bulletproofs [23]. Since the technical specification of the latest

## CryptoNote Signing Algorithm

**Sign**( $\{P_i\}_{i=1}^n, t_\ell, \ell, m$ ):

- Set  $I := \text{hash}_g(P_\ell)$ ;
- For  $i \in [k]$ , pick  $q_i \xleftarrow{\$} \mathbb{Z}_p$ ;
- For  $i \in [k], i \neq \ell$ , pick  $w_i \xleftarrow{\$} \mathbb{Z}_p$ ;
- For  $i \in [k]$ :
  - Set  $L_i := g^{q_i}$  if  $i = \ell$ ;
  - Set  $L_i := g^{q_i} \cdot P_i^{w_i}$  if  $i \neq \ell$ ;
  - Set  $R_i := (\text{hash}_g(P_i))^{q_i}$  if  $i = \ell$ ;
  - Set  $R_i := (\text{hash}_g(P_i))^{q_i} \cdot I^{w_i}$  if  $i \neq \ell$ ;
- Set  $c := \text{hash}_p(m, L_1, \dots, L_k, R_1, \dots, R_k)$ ;
- For  $i \in [k]$ :
  - Set  $c_i := w_i$  if  $i \neq \ell$ ;
  - Set  $c_i := c - \sum_{j=1}^k c_j$  if  $i = \ell$ ;
  - Set  $r_i := q_i$  if  $i \neq \ell$ ;
  - Set  $r_i := q_\ell - c_\ell t_\ell$  if  $i = \ell$ ;
- Return  $\sigma := (I, c_1, \dots, c_k, r_1, \dots, r_k)$ .

**Fig. 2.** CryptoNote Signing Algorithm.

*Ring Confidential Transaction* (RingCT) [24] has been developed and deployed in Monero since January 2017. It combines (linkable) ring signature and Pedersen commitment schemes [25], and also adopts Multilayered Linkable Spontaneous Anonymous Group Signature (MLSAG).

In Monero, suppose a user wants to spend  $m$  coins from his wallet, denoted as  $A_s := \{(\text{PK}_s^{(i)}, \text{CN}_s^{(i)})\}_{i=1}^m$  where  $\text{PK}_s^{(i)}$  is the user's  $i$ -th account address and  $\text{CN}_s^{(i)}$  is the balance of the account. The user first chooses  $k$  output accounts  $\{(\text{PK}_r^{(j)}, \text{CN}_r^{(j)})\}_{j=1}^k$  such that the sum of balances of the input accounts equals the output accounts, and sets  $R := \{\text{PK}_r^{(j)}\}_{j=1}^k$  as the output addresses. In addition, the user selects  $n - 1$  groups of input accounts with each containing  $m$  different accounts to anonymously spend  $A_s$ , i.e. set anonymity. Whenever receiving this transaction from the P2P blockchain network, the miners check the validity of the transaction along with its public information. The commitments are used to hide account balance. There are several special properties required for the RingCT protocol. Public keys generated by the key generation algorithm of ring signature should be homomorphic. Commitments should be homomorphic w.r.t. the same operation as public keys. Commitments to zero are well-formed public keys, each corresponding secret key of which can be derived from the randomness of commitments.

---

version is not well documented yet, our work is for Monero version 0.12.0.0 and earlier versions.

## Borromean Signing Algorithm

$\text{Sign}(\mathcal{P}, \{t_i\}_{i=0}^{n-1}, \{j_i^*\}_{i=0}^{n-1}, m)$ :

- For  $i \in [0, n-1]$ :
  - Pick  $k_i \xleftarrow{\$} \mathbb{Z}_p$ ;
  - Set  $e_{i,j_i^*+1} := \text{hash}_p(m, g^{k_i}, i, j_i^*)$ ;
  - For  $j \in [j_i^*, m_i - 1]$ , pick  $s_{i,j} \xleftarrow{\$} \mathbb{Z}_p$  and compute  $e_{i,j+1} := \text{hash}_p(m, g^{s_{i,j}} \cdot P_{i,j}^{-e_{i,j}}, i, j)$ ;
- For  $i \in [0, n-1]$ , pick  $s_{i,m_j} \xleftarrow{\$} \mathbb{Z}_p$  and compute  $e_0 := \text{hash}_p(g^{s_{i,m_j}} \cdot P_{i,j}^{-e_{i,m_j}}, \dots, g^{s_{n,m_j}} \cdot P_{i,j}^{-e_{n,m_j}})$ ;
- For  $i \in [0, n-1]$ :
  - For  $j \in [0, j_i^* - 1]$ , pick  $s_{i,j} \xleftarrow{\$} \mathbb{Z}_p$  and compute  $e_{i,j+1} := \text{hash}_p(m, g^{s_{i,j}} \cdot P_{i,j}^{-e_{i,j}}, i, j)$ ;
  - Set  $s_{i,j_i^*} := k_i + t_i e_{i,j_i^*-1}$ ;
- Return  $\sigma := (e_0, \{s_{i,j}\}_{i \in [0,n], j \in [0,m_i]})$ .

Fig. 3. Borromean Signing Algorithm.

In particular, we will explore our subversion attack against the Borromean ring signature [26]. In a high-level abstraction, Borromean ring signature is a Fiat-Shamir transformation of an AND/OR composition Sigma protocol of the Schnorr's identity protocol. More specifically, let  $\mathcal{P} := \{P_{i,j}\}_{i \in [0,n-1], j \in [0,m_i-1]}$  be a set of public keys, and the signer knows the secret key  $t_i$  such that  $P_{i,j_i^*} = g^{t_i}$ ,  $i \in [0, n-1]$ , where  $j_i^*$  are fixed and unknown indices. Moreover, we provide pseudo-code to explain the Borromean signing algorithms in Fig. 3.

## 2.6 Steganography

Steganography refers to the techniques that allow a sender to send a message covertly over a *communication channel* so that the mere presence of the hidden message is not detectable by an adversary who monitors the channel [27, 28]. Modern steganography techniques can be applied to various media, such as images, audios, HTML files, etc. A stegosystem consists of three PPT algorithms  $\mathcal{ST} := (\text{KeyGen}, \text{Embed}, \text{Extract})$  as follows:

- $(\text{ek}, \text{dk}) \leftarrow \text{KeyGen}(1^\lambda)$  is the key generation algorithm that takes as input the security parameter  $1^\lambda$ , and it outputs an embedding key  $\text{ek}$  and an extraction key  $\text{dk}$ .
- $\text{st} \leftarrow \text{Embed}_{\mathcal{H}, \text{ek}}(m)$ . Given an embedding key  $\text{ek}$ , a hidden message  $m \in \{0, 1\}^*$  and channel history  $\mathcal{H} \in \{0, 1\}^*$ ,  $\text{Embed}$  generates a stegotext message  $\text{st} \in \{0, 1\}^*$  that is *indistinguishable* from the normal channel distribution  $\mathcal{D}$  of innocent cover text objects  $\text{ct}$ .

- $m \leftarrow \text{Extract}_{\text{dk}}(\text{st})$ ,  $\text{Extract}$  takes as input an extraction key  $\text{dk}$  and the stego-text  $\text{st} \in \{0, 1\}^*$  and outputs a hidden message  $m \in \{0, 1\}^*$ .

**Definition 2 (Correctness).** We say a stegosystem  $\mathcal{ST} := (\text{KeyGen}, \text{Embed}, \text{Extract})$  is correct if for all  $(\text{ek}, \text{dk}) \leftarrow \text{KeyGen}(1^\lambda)$  we have

$$\Pr[\text{Extract}_{\text{dk}}(\text{Embed}_{\mathcal{H}, \text{ek}}(m))] \geq 1 - \text{negl}(\lambda) .$$

**Stegosystem Security.** The stegosystem’s goal is to communicate a hidden message covertly by hiding the mere existence of the hidden communication. Therefore, a stegosystem is considered to be *secure* if an observer is not able to distinguish stegotext  $\text{st}$  from objects randomly picked from the channel distribution  $\mathcal{D}$ . More formally, this is defined as a *chosen hidden-text attacks* (CHA) game/experiment.

$\text{Expt}_{\mathcal{A}}^{\text{CHA}}(1^\lambda)$

1.  $\mathcal{A}(1^\lambda)$  outputs a message  $m$ ;
2.  $(\text{ek}, \text{dk}) \leftarrow \text{KeyGen}(1^\lambda)$ ;
3.  $b \leftarrow \{0, 1\}$ ;
4. **If**  $b = 0$ :  $c \leftarrow \text{Embed}_{\mathcal{H}, \text{ek}}(m)$ ;  
**Else**:  $c \leftarrow \mathcal{D}$ ;
5.  $\mathcal{A}(c)$  outputs a bit  $b'$ ;
6. **Return**  $b \stackrel{?}{=} b'$ ;

We say a stegosystem  $\mathcal{ST} := (\text{KeyGen}, \text{Embed}, \text{Extract})$  is CHA-secure if:

$$\text{Adv}_{\mathcal{A}, \mathcal{ST}}^{\text{CHA}}(1^\lambda) = \left| \Pr[\text{Expt}_{\mathcal{A}}^{\text{CHA}}(1^\lambda)] - \frac{1}{2} \right| = \text{negl}(\lambda) .$$

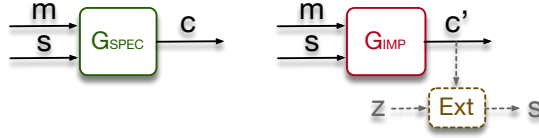
**Stegosystem Efficiency and Robustness.** Besides security, the following properties are also important to a stegosystem.

- Reliability/Efficiency. The probability that an *embedded* message is *extracted* when the stegosystem does not achieve not perfect correctness.
- Robustness. The inability of a challenger/warden to alter the sender’s communication transcript (that contain hidden message), and possibly prevent the receiver from recovering the hidden message.

## 2.7 Kleptography/Algorithm-substitution attacks

Our wallet subversion attacks can be classified as kleptographic attacks [29–31] and algorithm-substitution attacks (ASA) [32, 33]. As a high-level definition, in such attacks, the adversary maliciously tampers with the implementation of a cryptographic algorithm  $\mathbf{G}_{\text{IMP}}$  and changes it from its specification  $\mathbf{G}_{\text{SPEC}}$  algorithm, with the aim to subliminally and exclusively leak the user’s secret information to the adversary while evading detection in the black-box setting. The depiction in Fig. 4 illustrates how an adversarial implementation  $\mathbf{G}_{\text{IMP}}$  of the algorithm  $\mathbf{G}_{\text{SPEC}}$  can allow the adversary, given their secret key  $z$ , to detect the

subverted ciphertext  $c'$  and extract the user's secret  $s$ . Kleptographic attacks are significant due to their undetectability in the black-box setting and their severe consequences on the security of the users. See App. E for more details about signature subversion and *detectability*.



**Fig. 4.** Kleptography/ASA: specification  $G_{SPEC}$  takes input as the message  $m$  and the secret  $s$ , and outputs  $c$ ; whereas, the malicious implementation  $G_{IMP}$  outputs a subverted ciphertext  $c'$  which can leak the secret  $s$  exclusively to the attacker who knows  $z$ .

## 2.8 ECDSA

ECDSA is a randomized-signature scheme over the NIST elliptic curves that has been widely used in cryptocurrencies, such as Bitcoin, Ethereum, etc.

**Elliptic Curve Over  $\mathbb{F}_p$**  Let  $\text{param} := (p, a, b, g, q, \zeta)$  be the elliptic curve parameters over  $\mathbb{F}_p$ , consisting of a prime  $p$  specifying the finite field  $\mathbb{F}_p$ , two elements  $a, b \in \mathbb{F}_p$  specifying an elliptic curve  $E(\mathbb{F}_p)$  defined by  $E : y^2 \equiv x^3 + ax + b \pmod{p}$ , a base point  $g = (x_g, y_g)$  on  $E(\mathbb{F}_p)$ , a prime  $q$  which is the order of  $g$ , and an integer  $\zeta$  which is the cofactor  $\zeta = \#E(\mathbb{F}_p)/q$ . We denote the cyclic group generated by  $g$  as  $\mathbb{G}$ , and it is assumed that the DDH assumption holds over  $\mathbb{G}$ , that is for all PPT adversary  $\mathcal{A}$ :

$$\text{Adv}_{\mathbb{G}}^{\text{DDH}}(\mathcal{A}) = \left| \Pr \left[ \begin{array}{l} x, y \leftarrow \mathbb{Z}_q; b \leftarrow \{0, 1\}; h_0 = g^{xy}; \\ h_1 \leftarrow \mathbb{G} : \mathcal{A}(g, g^x, g^y, h_b) = b \end{array} \right] - \frac{1}{2} \right|$$

is negligible in  $\lambda$ .

**ECDSA description** The ECDSA signature scheme is depicted in Fig. 5.

## 3 Generic Steganographic Attack

Many cryptocurrencies use ring signatures to preserve users' privacy. For example, the CryptoNote framework [17], which is adopted by around 20 cryptocurrencies, uses ring signatures. As a demonstration, we describe how the *uncontrolled randomness* in CryptoNote's ring signature can be maliciously exploited. Namely, we show how the randomness within the ring signatures can be used to communicate covertly, store arbitrary information, and surreptitiously leak private keys. Note that the same principles are applicable to any other uncontrolled randomness in blockchain primitives.

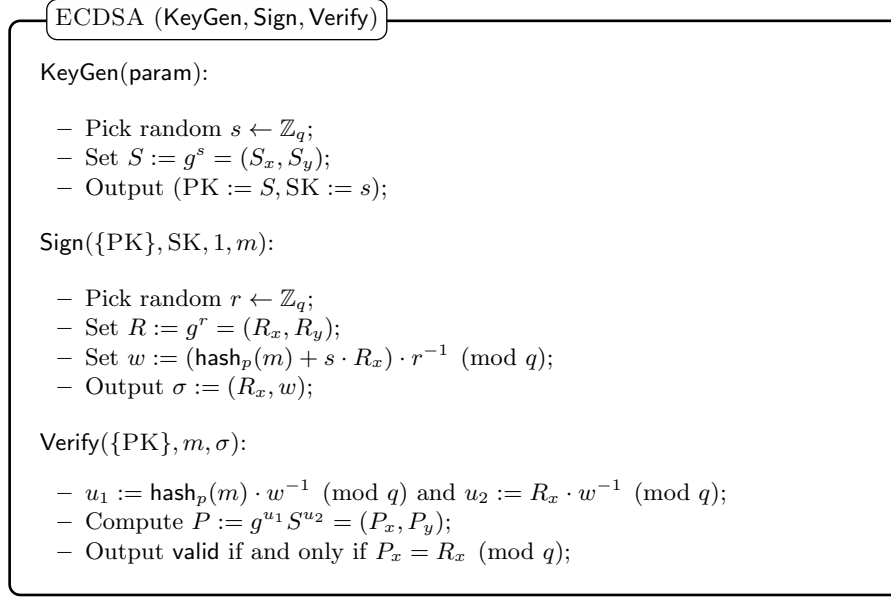


Fig. 5. ECDSA Signature Scheme.

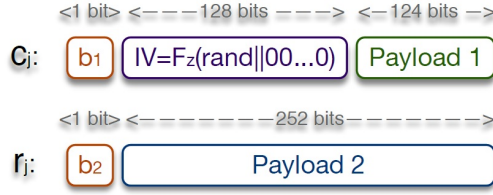
### 3.1 Our generic steganographic attack on CryptoNote

We now describe a generic steganographic attack against all CryptoNote-based cryptocurrencies and their variants. As mentioned in Sec. 2.4, the CryptoNote protocol uses the ED25519 twisted Edwards curve, and the group order is a 253-bit prime  $p$ . The long term secret key of a user consists of two group elements  $a, b \in \mathbb{Z}_p^*$ , but  $a := \text{hash}_p(b)$  is commonly used in practical implementation. Therefore, the long term secret key of a CryptoNote account is effectively 253 bits.

As part of the one-time (linkable) ring signature, a *one-out-of-many* non-interactive zero knowledge proof is included. More specifically, for a ring of size  $k$ , the format of the ring signature is  $\sigma = (I, c_1, \dots, c_k, r_1, \dots, r_k)$ . Suppose the sender's public key is  $\text{PK}_i$ ,  $i \in [k]$ . For all  $j \in [k]$  and  $j \neq i$ , the components  $c_j$  and  $r_j$  are uncontrolled random group elements in  $\mathbb{Z}_p$  and can be used for covert communication (cf. Fig. 2, above). Hence, our attack is premised on steganographically embedding arbitrary information on the ring signature's random numbers  $(c_j, r_j)$ .

In our attack example,  $\text{ek} = \text{dk}$ , which is a simple 128-bit random key  $z$ , is the common shared secret. The attack is explained as a three-step process carried out by two parties: a sender called Alice and a receiver called Bob.

*Step 1: embedding hidden messages (Embed).* As the most significant bit of a random  $\mathbb{Z}_p$  element does not have uniform distribution (which is more biased to 0), to ensure (computational) indistinguishability between stegotext  $\text{st}$  and



**Fig. 6.** Steganographic attack against CryptoNote: Format of one pair of random numbers  $(c_j, r_j)$  with 47-byte embedded stegotext.

innocuous random elements  $(c_j, r_j) \in \mathbb{Z}_p$ , Alice embeds her secret message  $m$  in the least significant 252 bits of  $c_j$  and  $r_j$ , whereas, the most significant bits  $b_1$  and  $b_2$  are sampled according to the real distribution of  $c_j$  and  $r_j$ . As depicted in Fig. 6, the rest of the bits consist of a 128-bit IV, 124-bit Payload 1, 252-bit Payload 2. Let  $F : \{0, 1\}^{128} \times \{0, 1\}^{128} \mapsto \{0, 1\}^{128}$  be a block cipher that takes as input a 128-bit plaintext and a 128-bit key, and outputs a 128-bit (pseudo-random) ciphertext. Moreover, Alice uses synthetic IV to allow Bob to efficiently identify which transactions on the blockchain contain stegotext  $st$ . In particular,  $IV := F_z(\text{rand}||00\dots0)$ , where  $\text{rand}$  is a 64-bits random string, and  $00\dots0$  is a 64-bit string of 0's. As a result, to check if a signature contains any  $st$ , Bob can simply try to decrypt a suspected IV, obtaining  $d := F_z^{-1}(IV)$ . If the lower half of  $d$  consists of 64 bits of 0's, then this signature contains stegotext  $st$ .

In our attack, Payload 1 and Payload 2 are jointly used to convey a 376-bit hidden message ( $m = \text{Payload 1}||\text{Payload 2}$ ). The payloads are encrypted via a semantically secure symmetric encryption under the secret key  $z$  and using IV. Also, to handle an arbitrary-length hidden message and ensure the resulting ciphertext has the same length as the message (besides the IV), Alice can use *Ciphertext Stealing* (CTS) as described in App. F.

*Step 2: identifying stegotext.* Unlike conventional P2P covert communication, before attempting to extract a hidden message from a transaction, Bob should first identify if the target transaction contains a stegotext  $st$ . As mentioned before, Bob can accomplish this by parsing IV from the first two  $c_j$ 's of the ring signature  $\sigma$  in a transaction, and checking whether the decryption of IV contains pattern 64 bits of 0's as shown in Fig. 6. Note that Embed embeds the hidden message  $m$  in one of the first two pairs of  $(c_j, r_j)$ . If  $c_1$  does not yield the IV, then Alice's secret index  $i$  must be 1, and Bob moves on to decrypt  $c_2$  which must contain the IV, otherwise, the signature is an innocent cover text  $ct$  that does not contain  $st$ .

*Step 3: extracting hidden messages (Extract).* Once a steganographic ring signature is successfully identified, Bob can use the Extract algorithm to extract the hidden message. More specifically, Bob collects Payload 1 and Payload 2 as depicted in Fig. 6. Bob then uses the extraction key  $dk := z$  to decrypt the payload, obtaining  $m := \text{CTS-Dec}_z(\text{IV}, \text{Payload 1}||\text{Payload 2})$ .



## A Generic CryptoNote Stegosystem

```

KeyGen( $1^{128}$ ):
  - Pick random  $z \leftarrow \{0, 1\}^{128}$ ;
  - Return  $ek := dk := z$ ;

Embed $_{\mathcal{H}, ek}(m)$ :
  - Pick random  $rand \leftarrow \{0, 1\}^{64}$ ;
  -  $IV := F_z(rand || 00 \dots 0)$ ;
  -  $\hat{m} := \text{CTS-Enc}_z(IV, m)$ ;
  - Payload 1 :=  $\hat{m}_{[0:123]}$ ;
  - Payload 2 :=  $\hat{m}_{[124:375]}$ ;
  - Sample random  $c \leftarrow \mathbb{Z}_p$ , and  $r \leftarrow \mathbb{Z}_p$ ;
  -  $c_{[1:128]} := IV$ ;
  -  $c_{[129:252]} := \text{Payload 1}$ ;
  -  $r_{[1:252]} := \text{Payload 2}$ ;
  - Return  $(c, r)$ ;

Extract $_{dk}(c, r)$ :
  -  $\alpha := F_z^{-1}(c_{[1:128]})$ ;
  - If  $\alpha_{[64:127]} \neq (00 \dots 0)$ :
    • Return  $\perp$ ;
  - Else:
    •  $IV := c_{[1:128]}$ ;
  - Payload 1 :=  $c_{[129:252]}$ ;
  - Payload 2 :=  $r_{[1:252]}$ ;
  -  $m := \text{CTS-Dec}_z(IV, \text{Payload 1} || \text{Payload 2})$ ;
  - Return  $m$ ;

```

**Fig. 7.** Pseudo-code for a generic stegosystem  $ST := (\text{KeyGen}, \text{Embed}, \text{Extract})$  to covertly communicate a 376-bit message  $m$  in *one* pair of innocuous-looking  $(c, r)$ , where  $\lambda = 128$ .

The pseudo-code in Fig. 7 further illustrates the generic steganographic attack on CryptoNote currencies. Note that, in practice, the **IV** and **Payload** can be encrypted under two different keys derived from a single master key  $z$ . However, for notation simplicity, we use the same key here.

### 3.2 Security

The security of the proposed generic stegosystem against all CryptoNote-based cryptocurrencies is examined for undetectability under the chosen hidden-text attacks (CHA) game/experiment. We remark that the content-insertion techniques that use non-standard Bitcoin scripts or exchange the public key with

an arbitrary string with *printable* characters, as mentioned in [1, 34], can be detected. However, our proposed steganographic attack on CryptoNote simply replaces random numbers with pseudo-random ciphers which, by definition of semantic security, are computationally indistinguishable from each other. Assuming the CTS-Enc algorithm described in App. F uses  $F$  as the internal PRF function, we have the following theorem.

**Theorem 1.** *If  $F$  is a secure pseudorandom function, the stegosystem  $ST := (\text{KeyGen}, \text{Embed}, \text{Extract})$  as shown in Fig. 7 is CHA secure.*

*Proof.* See App. A.

### 3.3 Robustness and Efficiency

In terms of robustness, it is easy to see that, unlike image steganography, the stegotext embedded in the signatures can never be removed while still preserving the functionality of the signatures. Therefore, there is no filter that can remove our stegotext.

**Throughput.** The only similar attack in literature is the proof-of-concept attack in [6] which sends a hidden message bit-by-bit through the rejection-sampling of the transaction address. Besides sending one bit of the hidden message  $m$  per transaction, their attack also sends one transaction per block. As a result, with 10 minutes to add a new block in Bitcoin, a sender needs more than 24 hours to send a message of 20 bytes. On the other hand, our steganographic attack takes advantage of the randomness within each ring signature in CryptoNote transactions. In fact, a CryptoNote transaction contains a ring signature for each input. Therefore, if a transaction  $\text{tx}$  has  $y$  number of inputs, and  $n$  public keys in the ring of each signature, then the total number  $N$  of random numbers  $(c_j, r_j)$  in  $\text{tx}$  is  $N = y * (n - 1) * 2$ . Whereas, the available bandwidth  $B$  in bytes is  $B = 32N$ . Hence, the available bandwidth in one transaction of 10 inputs and 10 public keys is more than 5KB. In comparison, other techniques that replace segments of the transaction, e.g. replacing P2SH scripts in Bitcoin transactions as done in Tithonus [35], can at maximum transmit 1KB per transaction. This further proves the efficiency of exploiting cryptographic randomness as opposed to replacing segment of the transaction itself.

Note that many blockchains offer an API to retrieve certain transactions and blocks. Therefore, if the receiver know the heights, i.e. indices, of the blocks that contain the steganographically communicated data, he does not need to check the whole blockchain. Therefore, a sender can communicate this kind of information off-line to minimize the receiver’s computational effort.

**Robustness against DoS attacks.** Censors can discover censorship-resistant proxies, e.g. Tor bridges, and block them. On the other hand, censors can not distinguish steganographically-subverted blockchain transactions, hence, they can not launch any targeted DoS attack unless they blacklist the whole blockchain which might have other financial ramifications. Additionally, from an attacker

perspective, exploiting uncontrolled randomness is advantageous over other content-insertion approaches that simply replace segments of the transactions, as done in Tithonus [35] and Catena [36]. Namely, other techniques are susceptible to policy changes where certain transactions, or scripts, become conspicuous or are no longer accepted, forcing the adoption of alternative techniques.

**Cost.** Content-insertion through the use of OP\_RETURN transactions and the arbitrary replacement of transaction addresses [34] render the funds unspendable. Therefore, these techniques burn funds. On the contrary, our proposed steganographic attack does not incur any additional cost, except for minimal transaction fees, as the sender can always send transactions to his own addresses. Technically, we can choose arbitrarily large ring size in a transaction. In practice, however, we found that a value between 20 and 30 is the optimal ring size to get a transaction included quickly with minimum transaction fees. To further clarify the cost per Byte, a Bytecoin transaction tx with 2 inputs and 21 public keys can take about 2 KB of covert data and costs 0.01 BCN as the minimum transaction fee which, given the current price of Bytecoin is \$ 0.000619 [14], costs \$ 0.0000619. Therefore, the cost of transmitting 1 GB covertly is  $\approx$  \$ 2.4. On the contrary, As shown in [35], Bitcoin-based Tithonus can covertly transmit up to 1650 Bytes in one transaction by replacing segments of the P2SH script of a multisignature transaction. Assuming the minimum transaction fee of 1 Satoshi/Byte and \$ 3657 [14] per Bitcoin, the cost of transmitting 1 GB is more than \$ 36,000.

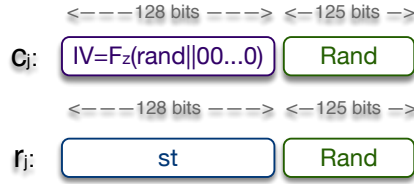
## 4 Case studies: Bytecoin and Monero

This section contains specific implementation of the proposed attack in Sec. 3 where we have implemented and evaluated the attack in two *real* cryptocurrencies; Bytecoin and Monero. Namely, we implemented the steganographic attack in the most recent release of Bytecoin (v 3.3.3) which has a market cap of around \$142 millions as of the time of writing [14]. Similarly, we implemented and tested the attack in Monero which is ranked 11 among currencies and has a market cap of around \$1 billion. It is important to note that as of October 2018, Monero (v 0.13.0.0) has replaced Borromean ring signatures, that is exploited by our attack, by a succinct zero-knowledge proof called Bulletproofs, which is not covered by this work. *Consequently, all of our discussion in relation with Monero is regarding previous versions of the source code mainly (v 0.12.0.0) and older.*

Although Monero is based on CryptoNote protocol, it uses Borromean ring signature which is different from the ring signature used in CryptoNote protocol as previously shown in Sec. 2. Nevertheless, our generic attack in Sec. 3 is still applicable to Monero. This emphasizes that the same attack can be extended to all public blockchain applications with randomized cryptographic primitives.

### 4.1 Implementation in Bytecoin

Bytecoin is an open-source cryptocurrency project [37] that implements the CryptoNote protocol described in Sec. 2.4. Accordingly, Bytecoin uses the ED25519



**Fig. 8.** Bytecoin: embedding a 16-byte  $st$  in one pair of  $(c_j, r_j)$  in transaction’s ring signature

twisted Edwards curve and CryptoNote (linkable) ring signature to sign its transactions. As previously mentioned in Sec. 2.4, this protocol has sufficiently many uncontrolled random numbers that could be exploited to covertly communicate arbitrary information. Since Bytecoin closely follows the specifications of the CryptoNote framework, it can directly be attacked using the generic steganographic attack described in Sec. 3. However, for code simplicity and clarity of demonstration, Ciphertext Stealing (CTS) is not used, and AES128 is used in the stegosystem because AES is already implemented in Bytecoin source code.

As a proof-of-concept experiment and due to ethical reasons, we only covertly transfer 16 bytes in the real-world Bytecoin without significantly abusing the blockchain system. Following the description of the generic attack in Sec. 3, we have implemented our steganographic attack on Bytecoin wallet in the following three steps.

*Step 1: embedding a hidden message  $m$  and generating a signature that contains  $st$ .* To embed a 16-byte hidden message  $m$  in a pair of random numbers  $(c_j, r_j)$ , Alice generates a synthetic  $IV := AES_z(rand||00\dots0)$  where  $rand$  is a 64-bit random string, and  $00\dots0$  is a 64-bit string of 0’s. Alice then places  $IV$  as the most significant 16 bytes of  $c_j$  and sets the rest of  $c_j$  randomly. She later uses this  $IV$  along with  $z$  to generate  $st$  that is embedded in the most significant 16 bytes of  $r_j$ . Namely,  $st := AES_z(m \oplus IV)$ . The format of  $(c_j, r_j)$  containing  $st$  is illustrated in Fig. 8.

Furthermore, to implement this step of the attack, the Bytecoin wallet’s source code is changed by mainly modifying one source file: `crypto.cpp`. The modified wallet simply alters the random numbers in the transaction’s ring signature(s) by producing *one* pair of  $(c_j, r_j)$  as aforementioned. Note that  $j \neq i$  where  $i$  is the signer’s *secret* index within the ring. Particularly, the changes introduced to `crypto.cpp` affect the following two functions within the source file:

- `generate_ring_signature()`: This function is slightly modified to pass a counter value to the `random_scalar` function.
- `random_scalar()`: This function is modified by including an additional parameter in its input to specify the counter. When this counter is 0 and 1, `random_scalar()` generates  $c_j$  and  $r_j$  respectively which are stegotexts that hide a 16-byte message as depicted in Fig. 8.

## Bytecoin covert communication pseudo code

```

KeyGen( $1^{128}$ ):
  - Pick random  $z \leftarrow \{0, 1\}^{128}$ ;
  - Return  $ek := dk := z$ ;

Embed $_{\mathcal{H},z}(m)$ :
generate_ring_signature():
  - If( $(j \neq i) \& (j == 0)$ ):
    •  $c_j := \text{random\_scalar}(0)$ ;
    •  $r_j := \text{random\_scalar}(1)$ ;
  - Else:
    • process as per normal;

random_scalar( $n$ ):
  -  $\text{rand} \leftarrow \mathbb{Z}_p$ ;
  - if( $n == 0$ ):
    •  $IV := \text{rand}_{[0:63]} || \text{zeros}$ ;
    •  $IV := \text{AES}_z(IV)$ ;
    •  $\text{rand}_{[0:127]} := IV$ ;
  - if( $n == 1$ ):
    •  $\text{rand}_{[0:127]} := \text{AES}_z(m \oplus IV)$ ;
  - Return  $\text{rand}$ ;

Extract $_z(c, r)$ :
  - for( $j = 0; j < 2; j++$ )
    •  $IV' := \text{AES}_z^{-1}(c_{j,[0:127]})$ ;
    • if( $IV'_{[64:127]} == \text{zeros}$ ):
      *  $m := \text{AES}_z^{-1}(r_{j,[0:127]}) \oplus c_{j,[0:127]}$ ;
      * Return  $m$ ;
  - Return 0;  % No hidden message

```

**Fig. 9.** Pseudo code for the implementation of covert communication in Bytecoin and similar currencies.

After generating the subverted signature that contains the stegotext, the transaction is sent as per normal over the blockchain. The sender does not need to modify other parts of the wallet source code.

*Step 2: identifying signature containing stegotext st.* To distinguish and identify signatures containing stegotext  $st$ , Bob checks every new transaction added to the ledger. To implement this, `BlockChainState.cpp` is slightly modified to check each signature by decrypting each pair of  $(c_j, r_j)$  numbers. In particular, Bob uses his key  $z$  to decrypt the most significant 16 bytes of  $c_j$  to check if it contains 64 bits of zeros as in Fig. 8. If he detects such a pattern, Bob identifies the existence of

a stegtext and sets IV as the most significant 16 bytes of  $c_j$ . If, however, no such pattern is detected, then the signature does not contain any hidden message.

*Step 3: extracting hidden message  $m$ .* After identifying a stegotext  $\mathbf{st}$ , Bob decrypts the most significant 16 bytes of  $r_j$  to extract  $m$ , that is  $m := \text{AES}_z^{-1}(r_{j,[0:127]}) \oplus (c_{j,[0:127]})$ . This process is further clarified by the pseudo code in Fig. 9.

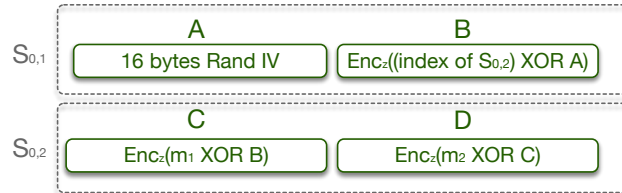
To further demonstrate the attack over the real Bytecoin blockchain, App. B provides a demo transaction included in the block at height 1671177 that contains a 16-byte hidden message “steganography”, and a tool to extract the steganographically embedded message.

## 4.2 Implementation for Monero (version 0.12.0.0)

Monero has a very complex cryptographic structure and ring signature scheme in particular. The core of Monero’s wallet involves Multilayered Linkable Spontaneous Anonymous Group Signature (MLSAG) and Borromean ring signature [26]. MLSAG is similar to the 1-out-of- $n$  ring signature that is used as part of the CryptoNote protocol; however, rather than using a ring signature on a set of  $n$  keys, MLSAG uses a ring signature on a set of  $n$ -key vectors. Using MLSAG, the signer proves to know all the private keys corresponding to one column in the public keys’ matrix. Despite the massive one-time secret key, the long-term secret key is still a single group element in  $\mathbb{Z}_p$ .

Borromean ring signature [26], which is a generalization and based on the 1-out-of- $n$  signature [38], is used to mask the transferred amount while enabling the receiver to know how much they have received by revealing the mask [39].

In our experiment, we chose to exploit the Borromean ring signature as it offers higher throughput. However, though with lower throughput, different primitives could also be exploited to mount steganographic attacks. Our attack on Monero is based on embedding a 32-byte hidden message  $m$  in the randomly generated  $s_{i,j}$  numbers as part of the Borromean ring signature [26]. Specifically, two vectors of  $s_{i,j}$  numbers are generated by the `genBorromean()` function:  $s_{0,j}$  and  $s_{1,j}$ . In addition,  $s_{0,j}$ ’s are randomly generated when the  $j^{\text{th}}$  bit commitment is 1. Two of these randomly generated  $s_{0,j}$ ’s are used to embed  $m$  as shown in Fig. 10. In a similar manner to our attack on Bytecoin, we use AES because it is already available in the source code. More details about the implementation of the steganographic attack on Monero can be found in App. C.



**Fig. 10.** Monero: embedding a 32-byte hidden message ( $m_1 || m_2$ ) in two random numbers ( $s_{0,1}, s_{0,2}$ ) in the Borromean ring signature

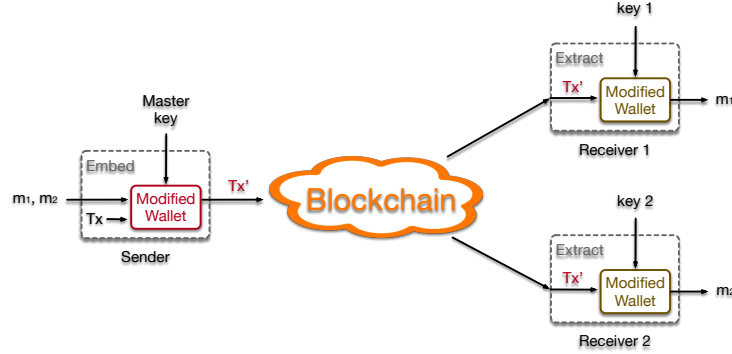


Fig. 11. Attack scenario 1: Covert broadcast communication.

## 5 Attack Scenarios

This section describes the following three attack scenarios: (i) covert broadcast communication, (ii) covert persistent storage, and (iii) wallet subversion attacks.

The first two scenarios represent direct applications of the steganographic attack in Sec. 3. Moreover, it is noticeable that both attack scenarios do not only facilitate objectionable behaviour, but can also hinder the very future of public blockchains. In particular, if a public blockchain is known to the authorities to be abused for covert communication or storage of malicious content, then authorities in any given country may criminalize the mere participation in such blockchains. Even if participation is not criminalized, users may choose not to store the full ledger, which defeats the purpose of decentralized blockchains, and leads to a more centralized setting, where few users participate in the consensus protocol. Unlike the first two scenarios, the third attack scenario is considered an *Algorithm-Substitution Attack* (ASA) and represents a scenario where the user is an oblivious victim of the attack.

### 5.1 Attack Scenario 1: Covert Broadcast Channel

Conventional steganographic techniques typically assume that the covert communication is between two parties – a sender and a receiver. In fact, our steganographic attack can be used as a covert broadcast channel, i.e. one sender and multiple receivers. As analyzed in Sec. 3, to steganographically send a hidden message of 1 KB, Alice can easily craft a transaction with 4 inputs and 5 public keys. Also, Fig. 11 shows that it is easy to use our steganographic technique in conjunction with some broadcast encryption scheme, e.g. [40], to enable a practical broadcast channel. The feasibility of this attack and the high throughput demonstrate the severity of this attack scenario, especially if abused by outlaws to use public blockchains as covert broadcast networks for their illicit communication.

## 5.2 Attack Scenario 2: Covert Data Storage and Distribution

Data storage can be viewed as a communication channel between the user and the user himself in the future. Unlike covert communication, covert persistent storage requires the uploaded content to be permanently stored and available on the blockchain. As discussed in Sec. 3, the cost of covertly storing 1 GB in Bytecoin’s blockchain is about \$ 2.4. Consequently, an adversary can use Bytecoin as a *cyberlocker* and abuse the P2P network of Bytecoin as a persistent content-distribution network (CDN). For example, it could be used to store pirated movies, wikileaks documents, etc.

An example special case of this scenario that shows the threat of such scenario is *blackmailing*. An adversary, Alice, can covertly store private information about a victim, Bob. Alice may even demonstrate this to Bob by sharing the key and the extraction tool with him. Alice can then threaten Bob that she can make the information publicly available by revealing the key to everyone.

## 5.3 Attack Scenario 3: Wallet Subversion

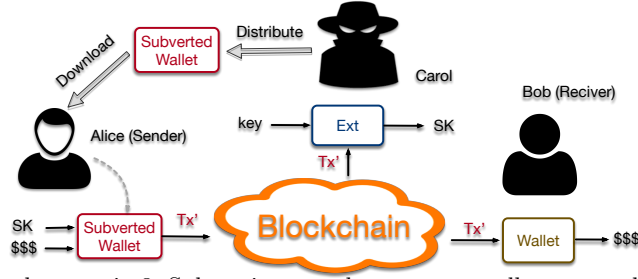
In the aforementioned attack scenarios, the sender, Alice, is complicit in the malicious attacks. This section presents another scenario where the sender is oblivious and is in fact a victim of the attack. Although this scenario may be applicable to *open-source* blockchain applications due to their complexity, it is more applicable to *close-source* and hardware-based applications, e.g. hardware wallets. The significance of this attack scenario stems from its undetectability in the black-box setting, where secrets are leaked via normal transactions posted on the blockchain, and its serious repercussions on the victim’s privacy and funds.

As depicted in Fig. 12, in this scenario, Alice is an innocent user who has downloaded, or bought, a wallet that is produced by a third party Carol who has maliciously implemented the wallet. In particular, Carol used a *subversion attack* to modify a wallet and redistribute it so to leak the signer’s private key, while evading detection in black-box settings. The way in which Carol modifies the wallet depends on the used cryptographic primitives and signature algorithms.

Below we present three subversion attacks that realize the scenario in Fig. 12. The first is a direct application of the generic steganographic attack described in Sec. 3 and its demo implementation in Bytecoin and Monero. Additionally, we present two more wallet subversion attacks targeting ECDSA-signature cryptocurrencies. Preliminary description of ECDSA can be found in Sec. 2.8. Namely, the first attack on ECDSA-signature crypto wallets uses synthetic ephemeral key to covertly leak the entire signer’s secret key over two signatures. However, it requires that the wallet is stateful in the sense that the wallet needs to store some variables from the previous signing execution. The second attack on ECDSA-signature crypto wallets is stateless and has lower throughput compared to the stateful attack. Note, in both ECDSA attacks, it is assumed that the attacker can identify the transactions generated by the victim user.

**Subverting Ring-Signature Crypto Wallets.** In the following we describe how the generic steganographic attack described in Sec. 3 is used by a third





**Fig. 12.** Attach scenario 3: Subversion attack on crypto wallets to steal users' private keys

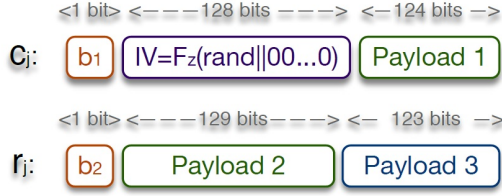
party, Carol, to subvert a ring-signature currency wallet, e.g. Bytecoin, to steal private keys. Similar attack is also applicable to Monero's Borromean signature.

Carol modifies the wallet by adding an embedding algorithm  $\text{Embed}_z(b)$ , where  $b$  is Alice's 253-bit private key  $b \in \mathbb{Z}_p$ . The subverted wallet secretly executes  $\text{Embed}_z(b)$  to generate one pair of  $(c_j, r_j)$  as shown in Fig. 13. In this scenario, the subverted  $c_j$  contains 128-bit IV that consists of an encryption of 64 random bits  $\text{rand}$  and 64 bits of zeros, i.e.  $\text{IV} := F_z(\text{rand}||00\dots0)$ .  $c_j$  also contains Payload 1 which is 124 bits of  $b$ .  $r_j$  contains Payload 2 which is 129 bits of  $b$  and Payload 3 which is the least significant 123 bits of Alice's public key  $B$ . The payloads are encrypted via a symmetric encryption under the same secret key  $z$  using IV.

Carol checks every added transaction for any exfiltrated private keys by decrypting the first 16 bytes of  $c_j$ 's from each signature, and checking if the decrypted text contains 64 bits of 0's as in Fig. 13. Note, Carol only needs to check the first two pairs of  $(c_j, r_j)$  to identify any subverted signature.

After successfully identifying a subverted signature, Carol parses and collects IV, Payload 1, Payload 2, and Payload 3. Carol then uses her secret key  $z$  to decrypt the payloads, obtaining  $b \in \mathbb{Z}_p$  and  $\text{LSB}_{123}(B)$ . After that she computes  $a := \text{hash}_p(b)$  and retrieves the corresponding public key  $(A, B)$  from the blockchain. After checking that  $A = g^a$  and  $B = g^b$ , Carol returns the secret key  $(a, b) \in (\mathbb{Z}_p)^2$ . Carol can now recover all the one-time addresses and transactions and even impersonate the compromised signer, Alice, to spend her money.

**Subverting ECDSA: Synthetic Randomness.** Our first proposed subversion attack on ECDSA is a simplified version of the attack proposed in [41]. The subverted algorithm is depicted in Fig. 14. Let  $z \in \mathbb{Z}_p$  be the adversary's secret key, and set  $Z := g^z$ . Let  $R \leftarrow \text{map}(R_x)$  be a mapping function that takes as input the x-coordinate and outputs the corresponding point on the curve. The subverted wallet needs to use algorithms  $\text{Sign}^{(1)}$  and  $\text{Sign}^{(2)}$  in turn to leak the signing key  $s$ . For the first time,  $\text{Sign}^{(1)}$  is identical to the original signature algorithm; however, the subtle difference is that  $\text{Sign}^{(1)}$  stores the ephemeral key  $r_1$  in a long-term memory, which can be accessed during the next signature invocation.  $\text{Sign}^{(2)}$  is also similar to the original signature algorithm except that it deterministically generates  $r_2 := \text{hash}_p(Z^{r_1})$ , where  $Z$  is hardcoded in the wallet. Once the adversary obtains two signatures  $\sigma_1, \sigma_2$ , he can use his secret



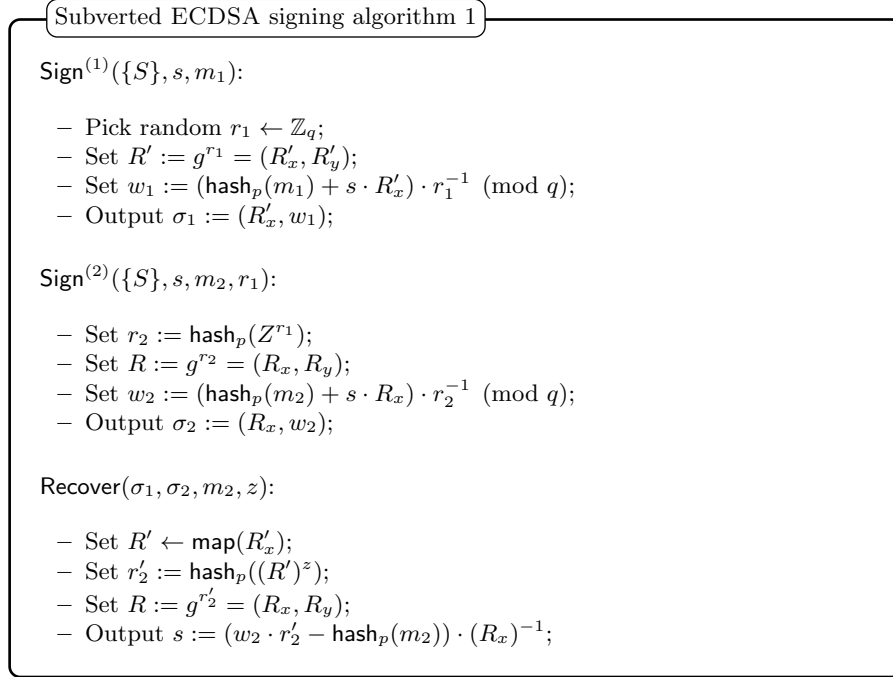
**Fig. 13.** Attack scenario 3: Covertly leaking the signer’s private key in *one* pair of  $(c_j, r_j)$ .

key  $z$  to recover the victim’s signing key  $s$ . First, he parses  $\sigma_1, \sigma_2$  as  $(R'_x, w_1)$  and  $(R_x, w_2)$ . The attacker then finds the point on the curve that corresponds to  $R'_x$ , using  $R' \leftarrow \text{map}(R'_x)$ . After that, the attacker computes  $r'_2 := \text{hash}_p((R')^z)$ . Note that if  $r'_2$  is equal to  $r_2$  then everything is correct. Let  $R := g^{r'_2} = (R_x, R_y)$ . The secret key can be extracted as  $s := (w_2 \cdot r'_2 - \text{hash}_p(m_2)) \cdot (R_x)^{-1}$ . This attack illustrates how the entire long term signing key  $s$  can be leaked exclusively to the adversary over two subverted signatures.

**Subverting ECDSA: Rejection Sampling.** While our first ECDSA subversion attack has a very high throughput, it has few drawbacks. First of all, it is a stateful algorithm, so it is not suitable for all scenarios, especially for software wallets. Furthermore, the first attack can only leak the signing key by the nature of its design, and not any other confidential information. Note that most cryptocurrency wallets are able to avoid the re-use of the address and signing key. As a result, the leaked signing key in our first attack, may never be used again even if the signing algorithms are executed twice with the same signing key. Nevertheless, for most wallets, there is a master key that is used to deterministically derive all the one-time signing keys.

As a result, our second subversion attack on ECDSA is stateless and is designed to leak arbitrary confidential information. As depicted in Fig. 15, the subverted signing algorithm takes as input the signing key  $s$ , the message  $m_i$ , and the secret  $x \in \{0, 1\}^n$  to be leaked. The signing algorithm leaks a random bit of  $x$  per signature. Let  $\text{PRF} : \{0, 1\}^* \times \{0, 1\}^\lambda \mapsto \{0, 1\}^{\log n} \times \{0, 1\}$  be a pseudo-random function that takes as input an arbitrary length message and the  $\lambda$ -bit PRF key, and it outputs a random number of  $(\log n + 1)$  bits. The first  $\log n$  bits is interpreted as an index  $j$ , and the last 1 bit is viewed as  $b$ . The subverted signing algorithm performs a rejection-sampling to find a random  $R = (R_x, R_y)$  such that  $(j, b) \leftarrow \text{PRF}_z(R_x)$  and  $x[j] = b$ . The rest signing process is identical to the original signature algorithm. Note that the rejection-sampling is efficient, and the expected repetition per signature is 1.5 times.

To recover the secret, the adversary needs to obtain a collection of the signatures generated by the subverted algorithm. We emphasize that when the secret is a master key that can be tested for its correctness, it is *not* necessary to leak the entire key in practice. Assuming the master key is 256 bits, to ob-

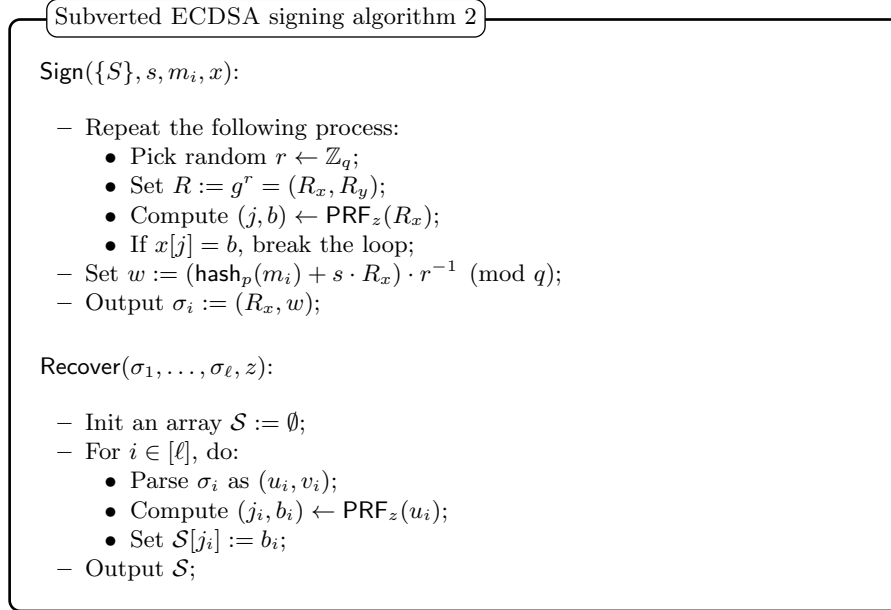


**Fig. 14.** The subverted ECDSA signing algorithm 1.

tain 50% distinct key bits, the expected number of signatures is bounded by approximately 256 signatures. Asymptotically, to obtain  $n$  secret bits, we need  $\theta(n \log n)$  signatures. To demonstrate this, we performed an experiment using our rejection sampling technique to empirically test the needed number of signatures to 32, 64, 96, 128, 160, 192 and 224 bits out of the total 256 key bits. This experiment was run 20 times to record the number of needed signatures to leak some bits of the secret key. As shown in Table 3 in App. D, the average number of signatures that should be intercepted by an attacker to retrieve 50% of the key, i.e. 128 bits, is about 179 signatures.

## 6 Countermeasures

In this section, we first examine the state-of-the-art of known countermeasures against subversion attacks and some techniques that can prevent exploiting uncontrolled randomness for arbitrary content insertion in blockchains. We then propose two countermeasures that are tailor-made for the blockchain scenarios. Note that our countermeasures aim to eliminate any steganographic messages hidden inside the cryptographic components (such as signatures and non-interactive zero-knowledge proofs) attached in a blockchain transaction; whereas,



**Fig. 15.** The subverted ECDSA signing algorithm 2

solutions to general content-insertion, e.g., inserting arbitrary content in unspendable OP\_RETURN Bitcoin transactions, is beyond the scope of this work.

## 6.1 Existing Countermeasures

In the literature, there are several known techniques that were proposed against generic steganography and substitution attacks. Here we systematically assess those techniques in the context of blockchain. In particular, they are divided into the following three categories: (I) *ASA-resistance techniques* which have been proposed to immunize cryptographic primitives against malicious implementation attacks, (II) *proactive trust-based countermeasures* where a trusted entity or a trusted initial state is used to prevent possible subversion, and (III) *blockchain-based techniques* that may and may not be primarily proposed to counter malicious content but can still be used to deter adverse covert communication and/or persistent storage. More details on the existing countermeasures can be found in App. G.

Table. 2 summarizes our assessment of the existing countermeasures w.r.t. their effectiveness against the attack scenarios in Sec. 5. The first and most intuitive countermeasure is to stop using any randomized primitives as noted in [32, 33]. As seen in Table 2, this countermeasure is theoretically effective, but its usage is very limited in the context of blockchains. Similarly, signatures with synthetic randomness, as proposed in [42], do not address covert channels

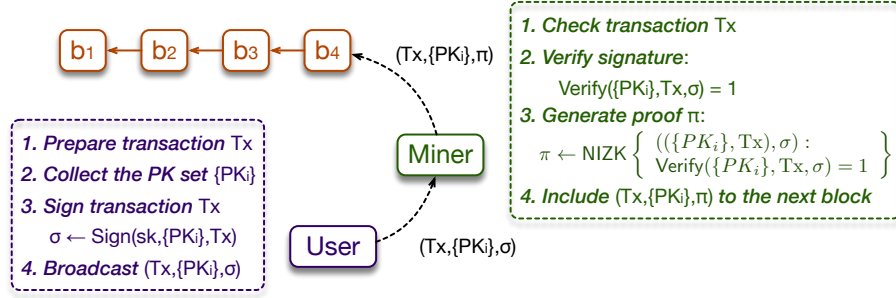
Technique	Attack Scenarios			Input Trigger	Practical in blockchains	References
	Attack 1 Covert Channels	Attack 2 Persistent Storage	Attack 3 Wallet Subversion			
<b>New Countermeasure</b>						
SRBF	/	✓	/	✓	✓	this work
Randomness-Externalized Arch.	x	x	✓	✓	✓	this work
<b>ASA-Resistant Techniques</b>						
Deterministic crypto primitives	✓	✓	✓	✓	x	[32, 33]
Signature with synthetic randomness	✓	✓	✓	x	x	[42]
Verifiable Random Function	x	x	x	✓	x	[43]
Split model	x	x	✓	x	x	[44–46]
<b>Active Trust-Based Countermeasures</b>						
Reverse firewalls	✓	✓	✓	✓	x	[47]
Self-guarding protocols	✓	✓	✓	✓	x	[48]
<b>Blockchain-Based Techniques</b>						
Light chains	x	✓	x	✓	✓	[49, 50]
Redactable chains	x	✓	x	✓	✓	[4, 5]
Content Filters	x	x	x	x	x	[3]
Increasing Transaction Fees	x	x	x	x	x	[3]
Self-verifying Addresses	x	x	x	✓	x	[3]

**Table 2.** Effectiveness of new and current countermeasures against the three attack scenarios mentioned in Sec. 5. Current countermeasures are categorized into three categories: *ASA-resistant techniques*, *Active trust-based countermeasures*, and *blockchain-based techniques*. (✓) denotes that the relevant countermeasure is resistant against the corresponding attack scenario, while (x) means the countermeasure is vulnerable to the attack scenario. *Input-trigger* states if a countermeasure is resistant to ‘time bombs’ or input-triggered malicious behaviour. Finally, *Practical in blockchains* examines if the corresponding countermeasure is applicable in the context of blockchains, if it is not likely to produce other security ramifications, and its robustness against malicious users. (/) is used to indicate that SRBF does not deter covert channels nor wallet subversion in *one case*: when the attackers’ nodes are alive at the broadcast stage of the malicious transactions, while it practically defends against all 3 scenarios in other cases.

and storage, and they are also susceptible to input-trigger attacks; thus, they are not practical in this context. Another countermeasure that can be used to sanitize randomness is the use of verifiable random functions (VRF) [43]. Though they may be able to minimize the bandwidth available for steganography per transaction, VRFs are vulnerable to rejection-sampling attacks. Another subversion-resistant technique was proposed by Russell et al. [44–46] is to double-split a given randomness generation function  $RG$  into two separate components  $RG_0$  and  $RG_1$ , execute the two components independently, and compose their outputs using a deterministic function  $\phi$  to generate the final subversion-free random number. However, this countermeasure setting is not practical.

Another category of countermeasures is the use trusted entities to eliminate steganography. The first countermeasure of this kind is to use reverse firewalls (RF) as proposed by Ateniese et al. [47] which requires a trusted entity to immunize re-randomizable signature schemes, and the second countermeasure is the use of *self-guarding* protocols proposed by Fischlin and Mazaheri [48] which require a trusted initial state of the algorithm. Both countermeasures are effective but not generic, as we need to use special signature schemes.

The last group of countermeasures are related to blockchains techniques. The first two techniques, light chains [49, 50] and redactable chains [4, 5], were proposed to address sustainability; nevertheless, they can defend against mali-



**Fig. 16.** Stego-resistant blockchain framework (SRBF)

cious persistent storage. Finally, several countermeasures were proposed in [3] to address content insertion in Bitcoin, but they are not practical to deter the exploitation of uncontrolled randomness.

## 6.2 Stego-Resistant Blockchain Framework (SRBF)

A typical blockchain transaction contains one or more cryptographic component(s), such as the signatures and non-interactive zero-knowledge (NIZK) proofs. As mentioned before, the existing reverse-firewall-based solutions utilize special re-randomizable signature schemes, which is not generic. Here we propose a universal stego-resistant blockchain framework (SRBF) that can be readily deployed to any off-the-shelf blockchain system. Without loss of generality, we explain our technique in terms of signature schemes, and it can be applied to NIZK proofs analogously.

As depicted in Fig. 16, in our setting, the miners are assumed to be trustworthy, and the user’s client might be maliciously modified to exploit the cryptographic components, e.g., signature, attached with the transactions to broadcast the steganographic information over the blockchain. Conventionally, upon receiving a transaction  $tx$ , the miners would check the validity of its associated signature  $\sigma$ , using the signature verification algorithm  $Verify(PK, tx, \sigma) = 1$ . The miners then include the transaction together with its signature as it is to the next block, which will be eventually uploaded to the blockchain. Therefore, other miners and users can verify the validity of the transaction as well. Can we drop the signature from the transaction while still preserving public verifiability? In our solution, the miner, instead of showing the signature to every other blockchain users/miners, it replaces the signature with a NIZK proof (see App. H for NIZK definition). Informally, the proof says that “I have seen a valid signature such that the signature verification algorithm  $Verify(PK, tx, \sigma) = 1$ ”. More precisely, we have the following NP relation:

$$\mathcal{R}_{\text{sig}} = \{((\{PK_i\}_{i=1}^n, tx), \sigma) \mid \text{Verify}(\{PK_i\}_{i=1}^n, tx, \sigma) = 1\}$$

In our stego-resistant blockchain framework, only  $(\text{tx}, \{\text{PK}_i\}_{i=1}^n, \pi)$  will be posted on the blockchain, where

$$\pi \leftarrow \text{NIZK} \left\{ \left( (\{\text{PK}_i\}_{i=1}^n, \text{tx}), \sigma \right) : \text{Verify}(\{\text{PK}_i\}_{i=1}^n, \text{tx}, \sigma) = 1 \right\}$$

is generated by trusted miners. The security guarantee of our solution is obvious, as signature  $\sigma$  is the witness of the corresponding NIZK proof. By NIZK definition,  $\pi$  does not leak any information about  $\sigma$ . Therefore, all the steganographic information hidden in the signatures are filtered out from the blockchain. In practice, we can use Bulletproofs [23] to implement SRBF efficiently.

*Remark.* Although this solution prevents the permanent storage of steganographic information, the subverted signatures can still propagate through the P2P network in the broadcast stage. Therefore, we only consider it fully effective against persistent storage in Table 2.

### 6.3 Randomness-Externalized Architecture

As mentioned above, our first countermeasure cannot prevent steganographic information from being propagated via the P2P network in the first place. We now discuss a randomness-externalized architecture to eliminate this drawback. It is designed to prevent kleptographic attacks while still enabling randomized signatures.

Before illustrating the proposed architectural modification, let us first examine the existing running environment of a cryptocurrency wallet. It is safe to assume that the majority of users download wallets' executable binaries directly from the corresponding cryptocurrency website or a third-party software distribution platform, e.g. Apple AppStore. During the running time, depending on the functionality, the wallet may consume randomness collected by the underlying operating system (OS). For instance, Linux kernel gathers entropy from keyboard timings, mouse movements and IDE timings, and the randomness pool can be accessed via `/dev/random` and `/dev/urandom`. Although the executable binary files can be potentially subverted, unlike conventional kleptographic settings, the randomness generator is usually a part of the underlying operating system and can be trusted if we can ensure sufficient entropy. Later in this section, we also address the case of untrusted OS-level randomness generator. Nevertheless, as shown in our attacks, trusted randomness source alone does not guarantee subversion-immunity, because the actual use of randomness in implementation may deviate from the corresponding software specifications. In fact, the implementation may simply bypass the specified software or hardware randomness generator. For example, in our attack, the subverted wallet uses  $r' := \text{Enc}_z(\text{sk}; r)$  as the randomness instead of the given randomness  $r$ , where  $z$  is the adversarial key and  $\text{sk}$  is the victim's secret key to be leaked. Besides, if the randomness consumption is not restricted, the wallet can also perform rejective sampling to leak information. To control randomness usage, we propose the following modifications.

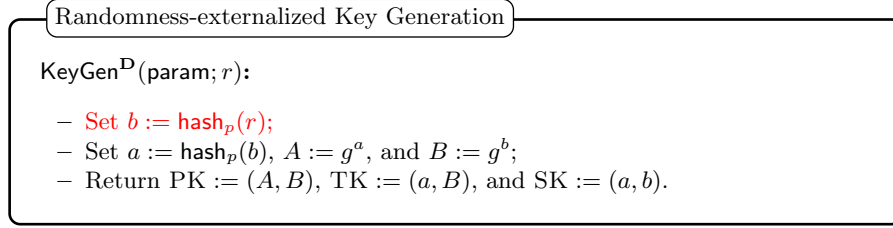


Fig. 17. Randomness-externalized key generation specification.

**Randomness-externalized Wallet.** To ensure correct usage of randomness, we need to make the cryptocurrency wallet deterministic by externalizing the software randomness draw. At the specification level, all the algorithms of the wallet are forbidden to have any internal randomness draw component. The algorithm specifications  $A_{\text{SPEC}}$  take a  $\lambda$ -bit randomness as an explicit input parameter, where  $\lambda \in \mathbb{N}$  is the security parameter, e.g. 256 in practice. For uniformity, all algorithms  $A_{\text{SPEC}}$  take the same amount of randomness. If more randomness is needed, they are *deterministically* derived from the input randomness  $r$  by setting  $r_i := \text{hash}(r, i)$ , where  $i \in \mathbb{N}$  is an index and  $\text{hash}$  is a cryptographically secure hash function, e.g. SHA3-256.

The main functionality of a cryptocurrency wallet is a signature scheme. Without loss of generality, a signature scheme consists of KeyGen, Sign, and Verify. (Of course, the linkable ring signature scheme used in CryptoNote also has a LNK algorithm to detect double spending, but it is deterministic.)

In the following, we modify the specification of KeyGen and Sign of CryptoNote wallet as examples. A typical key generation algorithm KeyGen takes as input the group parameter **param**, which defines the underlying group (or elliptic curve), e.g. Bitcoin uses SECP256K1 NIST curves and CryptoNote uses the ED25519 twisted Edwards curve. Of course, choosing a “nothing up the sleeve” group parameter is essential for subversion resilient cryptographic primitives; however, it is outside the scope of our work. We assume the commonly used blockchain group parameters are carefully examined, and are widely believed to be stego-free. The modified signature specification takes an explicit randomness  $r \in \{0, 1\}^\lambda$ . Moreover, Fig. 17 presents pseudo-code for our randomness-externalized architecture; whereas, Fig. 18 illustrates the architecture’s signing algorithm. It is easy to see that only a few number of lines need to be modified to make a signature scheme use external randomness. The difference between the modified version and the original version is marked in red.

**Hidden Trigger Elimination.** As mentioned before, in practice, a subverted signature algorithm may behave maliciously when a specific input message is given. Such a specific input message has high entropy, so an offline watchdog can only detect it with negligible probability. To remove hidden triggers, we need to randomize the input message. The proposed new signature scheme uses identical KeyGen, and the modified Sign\* and Verify\* algorithms are described in



## Randomness-externalized Signing Algorithm

**Sign**(param,  $\{P_i\}_{i \in [k]}, t_\ell, \ell, m; r$ ):

- Set  $I := \text{hash}_g(P_\ell)$ ;
- Set  $ctr := 0$ ;
- For  $i \in [k]$ , set  $q_i := \text{hash}_p(r, ctr)$  and  $ctr := ctr + 1$ ;
- For  $i \in [k], i \neq \ell$ , set  $w_i := \text{hash}_p(r, ctr)$  and  $ctr := ctr + 1$ ;
- For  $i \in [k]$ :
  - Set  $L_i := g^{q_i}$  if  $i = \ell$ ;
  - Set  $L_i := g^{q_i} \cdot P_i^{w_i}$  if  $i \neq \ell$ ;
  - Set  $R_i := (\text{hash}_g(P_i))^{q_i}$  if  $i = \ell$ ;
  - Set  $R_i := (\text{hash}_g(P_i))^{q_i} \cdot I^{w_i}$  if  $i \neq \ell$ ;
- Set  $c := \text{hash}_p(m, L_1, \dots, L_k, R_1, \dots, R_k)$ ;
- For  $i \in [k]$ :
  - Set  $c_i := w_i$  if  $i \neq \ell$ ;
  - Set  $c_i := c - \sum_{j=0}^k c_j$  if  $i = \ell$ ;
  - Set  $r_i := q_i$  if  $i \neq \ell$ ;
  - Set  $r_i := q_\ell - c_\ell t_\ell$  if  $i = \ell$ ;
- Return  $\sigma := (I, c_1, \dots, c_k, r_1, \dots, r_k)$ .

Fig. 18. Randomness-externalized signing algorithm specification.

## Signature scheme without hidden triggers

**Sign\***( $\{\mathbf{PK}\}, \mathbf{SK}, m; r$ ):

- Set  $s := \text{hash}(r, \text{"msg"})$  and  $m^* := \text{hash}(m, s)$ ;
- Compute  $\sigma \leftarrow \text{Sign}(\{\mathbf{PK}\}, \mathbf{SK}, m^*; r)$ ;
- Return  $\sigma^* := (\sigma, s)$ .

**Verify\***( $\{\mathbf{PK}\}, m, \sigma^*$ ):

- Parse  $\sigma^*$  as  $(\sigma, s)$ ;
- Compute  $m^* := \text{hash}(m, s)$ ;
- Return  $\text{Verify}(\{\mathbf{PK}\}, m^*, \sigma)$ .

Fig. 19. Signature with hidden trigger elimination.

Fig. 19. We now show that the new proposed signature scheme achieves strong existential unforgeability if the original signature scheme is strong existential unforgeable under adaptive chosen-message attack.

**Theorem 2.** *Let  $\mathcal{S} := (\text{Setup}, \text{KeyGen}, \text{Sign}, \text{Verify})$  be a signature scheme that achieves strong existential unforgeability under adaptive chosen-message attack. Let  $\text{hash} : \{0, 1\}^* \mapsto \{0, 1\}^\lambda$  be a cryptographic hash function. If  $\text{hash}$  securely*

realises a random oracle, then  $\mathcal{S}^* := (\text{Setup}, \text{KeyGen}, \text{Sign}^*, \text{Verify}^*)$  is also strong existential unforgeable under adaptive chosen-message attack.

*Proof.* See App. A.

**Extension.** In the extreme scenario when the underlying OS is also compromised, e.g. by malware, the randomness provided by the OS can no longer be trusted. We propose a “synthetic randomness” approach in the blockchain environment. Namely, in a cryptocurrency wallet, we can use the blockchain as a (weak) beacon for public randomness  $r_{\text{PUB}}$ , and use local randomness pool as private  $r_{\text{PRIV}}$ . The synthetic randomness is  $r := \text{hash}(r_{\text{PUB}} \| r_{\text{PRIV}})$ . The public randomness  $r_{\text{PUB}}$  can be computed by hashing the last  $\lambda - \kappa$  blocks,  $r_{\text{PUB}} := \text{hash}(B_{n-\lambda}, B_{n-\lambda+1}, \dots, B_{n-\lambda+\kappa})$ ; namely,  $\lambda - \kappa$  blocks of the common prefix [51]. The public randomness is used to ensure the min-entropy  $H_\infty(r) = \omega(\log \lambda)$  and thus eliminate any potential backdoors, whereas the local randomness is used to guarantee privacy. We assume it is difficult for an adversary to control both  $r_{\text{PUB}}$  and  $r_{\text{PRIV}}$  in practice. If this assumption is not satisfied, for paranoid users, they can use hash of stock market data or earthquake statistical data as a more trustworthy beacon source.

## 7 Related Work

This work is closely related to the topics of *malicious content insertion in blockchains*, *steganography*, *covert channels in blockchains*, and *kleptography/ASA* and their countermeasures.

**Content insertion in blockchains.** The authors of [34] provided insight regarding the various ways that could be exploited to store, possibly illegal, content onto the Bitcoin blockchain. Furthermore, using some heuristics to analyze the *plaintext* of 146 million transactions, the authors of [34] reported that 0.8% transactions store content on the blockchain or use non-standard scripts. Recently, the authors of [1] attempted to systematically analyze the non-financial content in Bitcoin’s blockchain. Specifically, they surveyed the methods and services that are used to store non-financial content, and provided a general categorization of objectionable content that could be found on Bitcoin’s blockchain. They found that 1.4% of all Bitcoin transactions contain non-financial data, and retrieved over 1600 files, some of which contain objectionable content. Nonetheless, there are some non-malicious applications that rely on Bitcoin-based content insertion. For example, Tithonus [35] offers a Bitcoin-based censorship-resistant system, and Catena [36] is an application that uses Bitcoin OP\_RETURN transactions to establish consensus among users on an application-specific log.

**Steganography.** The concept of *steganography* was introduced by Simmons’ *prisoner’s problem* [52]. Anderson et al. listed some of the limits of steganography and discussed the difficulty associated with formalizing a general proof of security for steganography [53, 54]. A number of works, e.g. [55–57], provided information-theoretic treatment of steganography security and robustness. More recently, Hopper et al. provided a cryptographic formalization of steganographic

security and robustness [27]. They presented a definition for the security of a steganographic system in terms of the *computational* indistinguishability of stegotext from cover text.

**Covert channels in blockchains.** While there is a relatively significant body of research about *content insertion* in Bitcoin’s blockchain [1, 3, 34, 58], the authors of [6] were the first to discuss the use of *steganography* to *covertly* communicate on Bitcoin’s blockchain. However, due to its limitation, the authors of [6] consider their attack to be a proof of concept rather than a practical attack.

**Kleptography and ASA.** Our wallet subversion attack scenario falls within the realm of *Algorithm-Substitution Attacks (ASA)* [32, 33], also called *Kleptography* [29, 30] and *Subversion Attacks (SA)* [47]. The notion of Kleptography was introduced by Young and Yung in 1996 [29, 30]. Subsequent work demonstrated the possible use of ASA in mass surveillance, and the susceptibility of all randomized symmetric encryption schemes to such attacks [33, 59]. Another demonstration of ASA attacks is found in the work of Goh et al. [60] which presented practical hidden key-recovery attacks against the SSL/TLS and SSH2 protocols by modifying the implementation of the OpenSSL library. In the context of signature schemes, Young and Yung [31] showed that DSA signature schemes can be subverted to leak secret information. Another kleptographic attack was proposed by the work of Teşeleanu [61] which describes a threshold kleptographic attack on the generalized ElGamal signature that can be extended to similar DL-based signatures. In addition, as a countermeasure against subversion, Russell et al. [44] modeled and proved a full domain hash-based signature scheme achieves subversion resilience. Recently, Russell et al. [45, 46] proposed the use of a splitting-randomness technique to secure a randomizable IND-CPA public-key encryption, but it is unknown how to apply their technique in the blockchain context with reasonable assumptions. Another countermeasure was proposed by Ateniese et al. [47] who proposed the use of trusted reverse firewalls to re-randomize the output of possibly subverted signature algorithms. Finally, Fischlin and Mazaheri [48] proposed a novel technique that proactively defends against ASA’s assuming *temporary initial trust* of the possibly subverted algorithm.

## 8 Conclusion

The main aim of this work is to highlight the potential threat of maliciously abusing uncontrolled randomness in randomized cryptographic primitives in blockchain applications, and design efficient and practical countermeasures. To illustrate the idea, we designed, implemented, and evaluated our attacks against the widely-used ECDSA signature scheme, the ring signature used in the CryptoNote framework, and the Ring Confidential Transaction used in Monero (up to version 0.12.0.0). The demonstrated attacks can be used in three malicious scenarios: covert communication, persistent storage of objectionable data, and wallet subversion attacks. Finally, we emphasize that this line of research is far from being

completed, and we hope that our work motivates the design of stego-resistant blockchains.

## Acknowledge

We would like to sincerely thank Kenneth G. Paterson for his constructive comments on an early version of this work, and other anonymous reviewers for their feedback.

## References

1. R. Matzutt, J. Hiller, M. Henze, J. H. Ziegeldorf, D. Mullman, O. Hohlfeld, and K. Wehrle, "A quantitative analysis of the impact of arbitrary blockchain content on bitcoin," in *FC 2018*, 2018.
2. J. Smith, J. Tennison, P. Wells, J. Fawcett, and S. Harrison, "Applying blockchain technology in global data infrastructure," tech. rep., Open Data Institute, June 2016. ODI-TR-2016-001.
3. R. Matzutt, M. Henze, J. H. Ziegeldorf, J. Hiller, and K. Wehrle, "Thwarting unwanted blockchain content insertion," in *IC2E 2018*, pp. 364–370, April 2018.
4. I. Puddu and A. Dmitrienko, "μchain: How to forget without hard forks," *IACR Cryptology ePrint Archive 2017/106*, 2017. <https://eprint.iacr.org/2017/106>.
5. G. Ateniese, B. Magri, D. Venturi, and E. Andrade, "Redactable blockchain or rewriting history in bitcoin and friends," in *Euro S&P 2017*, pp. 111–126, April 2017.
6. J. Partala, "Provably secure covert communication on blockchain," *Cryptography*, vol. 2, no. 3, 2018.
7. G. Fuchsbaauer, "Subversion-zero-knowledge snarks," in *PKC 2018*, pp. 315–347, 2018.
8. B. Abdolmaleki, K. Bagheri, H. Lipmaa, and M. Zajac, "A subversion-resistant snark," in *ASIACRYPT 2017*, pp. 3–33, 2017.
9. J. Knockel, T. Ristenpart, and J. R. Crandall, "When textbook RSA is used to protect the privacy of hundreds of millions of users," *CoRR*, vol. abs/1802.03367, 2018.
10. L. Luu, D.-H. Chu, H. Olickel, P. Saxena, and A. Hobor, "Making smart contracts smarter," in *CCS 2016*, pp. 254–269, 2016.
11. D. Siegel, "Understanding the dao attack," 2016. Available Online: <https://www.coindesk.com/understanding-dao-hack-journalists> (Last accessed 7-Feb-2018).
12. S. Azouvi, M. Maller, and S. Meiklejohn, "Egalitarian society or benevolent dictatorship: The state of cryptocurrency governance," in *5th Workshop on Bitcoin and Blockchain Research*, 2018.
13. Giza Device Ltd, "Giza wallet," 2017. Available Online: <https://www.gizadevice.com/> (Last accessed 7-Feb-2018).
14. CoinMarketCap, "Cryptocurrency market capitalizations," 2018. Available Online: <https://coinmarketcap.com/> (Last accessed 26-Nov-2018).
15. S. Nakamoto, "A Peer-to-Peer Electronic Cash System," 2008. Available Online: <https://bitcoin.org/bitcoin.pdf> (Last accessed 05-Nov-2018).

16. A. Bender, J. Katz, and R. Morselli, “Ring signatures: Stronger definitions, and constructions without random oracles,” *J. Cryptol.*, vol. 22, pp. 114–138, Dec. 2008.
17. N. V. Saberhagen, “Cryptonote v 2.0,” 2013. whitepaper, Available online: <https://cryptonote.org/whitepaper.pdf>, (Last accessed 23-Nov-2018).
18. Bytecoin Org., “Bytecoin (bcn),” 2018. Available Online: <https://bytecoin.org/> (Last accessed 23-Nov-2018).
19. CryptoNote Org., “Cryptonotecoin,” 2018. Available Online: <http://cryptonote-coin.org/> (Last accessed 23-Nov-2018).
20. Fantomcoin, “Fantomcoin,” 2014. Available Online: <http://fantomcoin.org/> (Last accessed 23-Nov-2018).
21. Monero, “Monero,” 2018. Available Online: <https://getmonero.org/> (Last accessed 07-Feb-2018).
22. R. Spagni, “Monero project github repository,” 2018. Available Online: <https://github.com/monero-project/monero> (Last accessed 07-Feb-2017).
23. B. Bünz, J. Bootle, D. Boneh, A. Poelstra, P. Wuille, and G. Maxwell, “Bulletproofs: Short proofs for confidential transactions and more,” in *SECP 2018*, vol. 00, pp. 319–338, 2018.
24. S. Noether, “Ring signature confidential transactions for monero.” Cryptology ePrint Archive, Report 2015/1098, 2015.
25. T. P. Pedersen, “Non-interactive and information-theoretic secure verifiable secret sharing,” in *CRYPTO '91*, pp. 129–140, 1992.
26. G. Maxwell and A. Poelstra, “Borromean Ring Signatures,” 2015. Available Online: <http://diyhpl.us/~bryan/papers2/bitcoin/Borromean%20ring%20signatures.pdf> (Last accessed 07-Feb-2018).
27. N. J. Hopper, J. Langford, and L. von Ahn, “Provably secure steganography,” in *CRYPTO 2002*, 2002.
28. N. Dedić, G. Itkis, L. Reyzin, and S. Russell, “Upper and lower bounds on black-box steganography,” *Journal of Cryptology*, vol. 22, pp. 365–394, Jul 2009.
29. A. Young and M. Yung, “The dark side of “black-box” cryptography or: Should we trust capstone?,” in *CRYPTO '96*, 1996.
30. A. Young and M. Yung, “Kleptography: Using cryptography against cryptography,” in *EUROCRYPT '97*, 1997.
31. A. Young and M. Yung, “The prevalence of kleptographic attacks on discrete-log based cryptosystems,” in *CRYPTO '97*, 1997.
32. M. Bellare, K. G. Paterson, and P. Rogaway, “Security of symmetric encryption against mass surveillance,” in *CRYPTO 2014*, (Berlin, Heidelberg), pp. 1–19, Springer Berlin Heidelberg, 2014.
33. M. Bellare and J. Jaeger, “Mass-surveillance without the State : Strongly Undetectable Algorithm-Substitution Attacks,” in *CCS*, 2015.
34. R. Matzutt, O. Hohlfeld, M. Henze, R. Rawiel, J. H. Ziegeldorf, and K. Wehrle, “Poster: I don’t want that content! on the risks of exploiting bitcoin’s blockchain as a content store,” in *CCS '16*, 2016.
35. R. Recabarren and B. Carbunar, “Tithonus: A bitcoin based censorship resilient system,” *Proceedings on Privacy Enhancing Technologies*, vol. 2019, no. 1, pp. 68 – 86, 2019.
36. A. Tomescu and S. Devadas, “Catena: Efficient non-equivocation via bitcoin,” in *2017 IEEE Symposium on Security and Privacy (SP)*, vol. 00, pp. 393–409, May 2017.
37. B. D. Team, “Bytecoin project github repository,” 2018. Available Online: <https://github.com/bcndev> (Last accessed 26-Nov-2018).

38. M. Abe, M. Ohkubo, and K. Suzuki, “1-out-of-n signatures from a variety of keys,” in *ASIACRYPT 2002*, 2002.
39. G. Maxwell, “Confidential Transactions,” 2018. Available Online: [https://people.xiph.org/~greg/confidential\\_values.txt](https://people.xiph.org/~greg/confidential_values.txt) (Last accessed 07-Feb-2018).
40. A. Fiat and M. Naor, “Broadcast encryption,” in *CRYPTO’ 93* (D. R. Stinson, ed.), (Berlin, Heidelberg), pp. 480–491, Springer Berlin Heidelberg, 1994.
41. E. Mohamed and H. Elkamchouchi, “Kleptographic Attacks on Elliptic Curve Cryptosystems,” *Journal of Computer Science*, vol. 10, no. 6, pp. 213–215, 2010.
42. T. Pornin, “Deterministic DSA and ECDSA,” RFC 6979, RFC Editor, August 2013. Available online: <https://tools.ietf.org/html/rfc6979> (Last accessed 13-Feb-2019).
43. S. Micali, M. Rabin, and S. Vadhan, “Verifiable random functions,” in *40th Annual Symposium on Foundations of Computer Science*, pp. 120–130, Oct 1999.
44. A. Russell, Q. Tang, M. Yung, and H.-S. Zhou, “Cliptography: Clipping the power of kleptographic attacks,” in *ASIACRYPT*, 2016.
45. A. Russell, Q. Tang, M. Yung, and H.-S. Zhou, “Destroying steganography via amalgamation: Kleptographically CPA secure public key encryption.” *Cryptology ePrint Archive*, Report 2016/530, 2016.
46. A. Russell, Q. Tang, M. Yung, and H.-S. Zhou, “Generic semantic security against a kleptographic adversary,” *CCS ’17*, (New York, NY, USA), pp. 907–922, ACM, 2017.
47. G. Ateniese, B. Magri, and D. Venturi, “Subversion-resilient signature schemes,” in *CCS ’15*, pp. 364–375, 2015.
48. M. Fischlin and S. Mazaheri, “Self-guarding cryptographic protocols against algorithm substitution attacks,” in *2018 IEEE 31st Computer Security Foundations Symposium (CSF)*, pp. 76–90, July 2018.
49. A. Molina and H. Schoenfeld, “PascalCoin Whitepaper v2,” 2017. Available Online: <https://www.pascalcoin.org/PascalCoinWhitePaperV2.pdf> (Last accessed 08-November-2018).
50. J.D. Bruce, “The Mini-Blockchain Scheme - Rev. 3,” 2017. Available Online: <http://cryptonite.info/files/mbc-scheme-rev3.pdf> (Last accessed 08-November-2018).
51. J. Garay, A. Kiayias, and N. Leonardos, “The bitcoin backbone protocol: Analysis and applications,” in *EUROCRYPT 2015* (E. Oswald and M. Fischlin, eds.), (Berlin, Heidelberg), pp. 281–310, Springer Berlin Heidelberg, 2015.
52. G. J. Simmons, *The Prisoners’ Problem and the Subliminal Channel*, pp. 51–67. Boston, MA: Springer US, 1984.
53. R. Anderson, “Stretching the limits of steganography,” in *Information Hiding*, pp. 39–48, Springer Berlin Heidelberg, 1996.
54. R. J. Anderson and F. A. P. Petitcolas, “On the limits of steganography,” *IEEE Journal on Selected Areas in Communications*, vol. 16, pp. 474–481, May 1998.
55. J. A. O’Sullivan, P. Moulin, and J. M. Ettinger, “Information theoretic analysis of steganography,” in *Proceedings. 1998 IEEE International Symposium on Information Theory*, 1998.
56. T. Mittelholzer, “An information-theoretic approach to steganography and watermarking,” in *Information Hiding*, 2000.
57. C. Cachin, “An information-theoretic model for steganography,” in *Information Hiding*, 1998.
58. K. Shirriff, “Hidden surprises in the bitcoin blockchain and how they are stored: Nelson mandela, wikileaks, photos, and python software,” 2014. Available Online:

- <http://www.righto.com/2014/02/ascii-bernanke-wikileaks-photographs.html> (Last accessed 01-November-2018).
59. M. Bellare and V. T. Hoang, “Resisting randomness subversion: Fast deterministic and hedged public-key encryption in the standard model,” in *EUROCRYPT 2015*, 2015.
  60. E.-J. Goh, D. Boneh, B. Pinkas, and P. Golle, “The design and implementation of protocol-based hidden key recovery,” in *Information Security*, 2003.
  61. G. Tesseanu, “Threshold kleptographic attacks on discrete logarithm based signatures.” Cryptology ePrint Archive, Report 2017/953, 2017. <https://eprint.iacr.org/2017/953>.
  62. J. Katz and V. Vaikuntanathan, “Signature schemes with bounded leakage resilience,” in *ASIACRYPT 2009*, 2009.
  63. E. Boyle, G. Segev, and D. Wichs, “Fully leakage-resilient signatures,” in *EUROCRYPT 2011*, 2011.
  64. M. J. Dworkin, “Recommendation for block cipher modes of operation: Three variants of ciphertext stealing for cbc mode.” NIST Pubs, Report Number 800-38A Addendum, 2010. Available Online: <https://www.gpo.gov/fdsys/pkg/GOV PUB-C13-c0b0bae5f66880bf051f6d4ac2d8f07d/pdf/GOV PUB-C13-c0b0bae5f66880bf051f6d4ac2d8f07d.pdf> (Last accessed 08-Feb-2018).
  65. M. Cardinal, *Executable Specifications with Scrum*. Upper Saddle River, NJ: Addison-Wesley, 2014.
  66. L. Hanzlik, K. Kluczniak, and M. Kutyłowski, “Controlled randomness – a defense against backdoors in cryptographic devices,” in *Paradigms in Cryptology – Mycrypt 2016. Malicious and Exploratory Cryptology* (R. C.-W. Phan and M. Yung, eds.), (Cham), pp. 215–232, Springer International Publishing, 2017.
  67. Docker-Inc., “Docker,” 2018. Available Online: <https://www.docker.com> (Last accessed 7-Feb-2018).
  68. Mini-blockchain Project, “Cryptonite Cryptocurrency,” 2018. Available Online: <http://cryptonite.info/> (Last accessed 01-November-2018).
  69. J. Camenisch, D. Derler, S. Krenn, H. C. Pöhls, K. Samelin, and D. Slamanig, “Chameleon-hashes with ephemeral trapdoors,” in *PKC 2017*, 2017.

## Appendices

### A Security Proofs

#### A.1 Proof of Theorem 1

In this section, we provide a full proof of Theorem 1.

*Proof.* We prove this theorem by reduction. Assume there exists a PPT adversary  $\mathcal{A}$  who can break  $\mathcal{ST}$  with a non-negligible  $\text{Adv}_{\mathcal{A}, \mathcal{ST}}^{\text{CHA}}(1^\lambda)$  advantage w.r.t. the CHA game. We need to construct a PPT adversary  $\mathcal{B}$  who can break the PRF game for  $F$ . During the reduction game,  $\mathcal{B}$  plays as a challenger for  $\mathcal{A}$  in the CHA game. Upon receiving  $m$  from  $\mathcal{A}$ ,  $\mathcal{B}$  picks random  $\text{rand} \leftarrow \{0, 1\}^{64}$  and sets  $x := \text{rand} || 00 \dots 0$ .  $\mathcal{B}$  then queries  $x$  to the PRF game challenger and obtains  $\text{IV}$ . Subsequently,  $\mathcal{B}$  queries  $\text{IV}, \text{IV} + 1, \text{IV} + 2$  to the PRF game challenger, and obtains  $k_1, k_2, k_3$ .  $\mathcal{B}$  then compute  $c_1, c_2, c_3$  according to the description shown in Fig. 22. It then computes  $(c, r)$  as described in Fig. 7.  $\mathcal{B}$  flips a coin  $b \leftarrow \{0, 1\}$ . If  $b = 0$ ,  $\mathcal{B}$  computes a ring signature using  $(c, r)$ ; otherwise,  $\mathcal{B}$  computes a ring signature normally.  $\mathcal{B}$  then sends the resulting signature to  $\mathcal{A}$ . Finally,  $\mathcal{A}$  outputs a guess  $b'$ . Assume the challenge bit in the PRF game is  $\beta$ , i.e.  $\beta = 0$  is in the PRF mode;  $\beta = 1$  is in the random function mode. If  $b = b'$ ,  $\mathcal{B}$  outputs  $\beta^* = 0$ ; otherwise,  $\mathcal{B}$  outputs  $\beta^* = 1$ .

$$\begin{aligned} \Pr[\mathcal{B} \text{ win}] &= \Pr[\beta^* = 0 | \beta = 0] \cdot \Pr[\beta = 0] + \\ &\quad + \Pr[\beta^* = 1 | \beta = 1] \cdot \Pr[\beta = 1] \\ &= \Pr[\mathbf{Expt}_{\mathcal{A}}^{\text{CHA}}(1^\lambda)] \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} \\ &= (\text{Adv}_{\mathcal{A}, \mathcal{ST}}^{\text{CHA}}(1^\lambda) + \frac{1}{2}) \cdot \frac{1}{2} + \frac{1}{4} \\ &= \frac{1}{2} \cdot \text{Adv}_{\mathcal{A}, \mathcal{ST}}^{\text{CHA}}(1^\lambda) + \frac{1}{2} \end{aligned}$$

Hence, the advantage of  $\mathcal{B}$  w.r.t to the PRF game is

$$\text{Adv}_{\mathcal{B}, F}^{\text{PRF}} = \left| \Pr[\mathcal{B}] \text{ win} - \frac{1}{2} \right| = \frac{1}{2} \cdot \text{Adv}_{\mathcal{A}, \mathcal{ST}}^{\text{CHA}}(1^\lambda) .$$

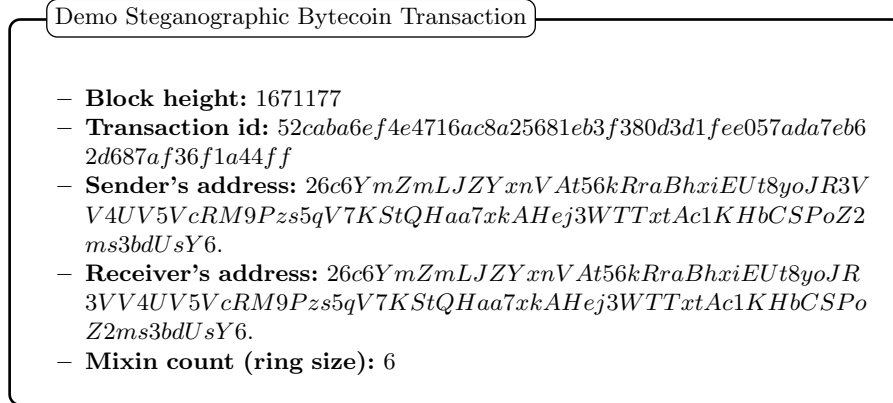
Since  $\text{Adv}_{\mathcal{A}, \mathcal{ST}}^{\text{CHA}}(1^\lambda)$  is non-negligible, we have  $\text{Adv}_{\mathcal{B}, F}^{\text{PRF}}$  is also non-negligible, which concludes the proof.

#### A.2 Proof of Theorem 2

In this section, we provide a full proof of Theorem 2.

*Proof.* We proof this theorem by reduction. Assume there exists a PPT adversary  $\mathcal{A}$  who can break  $\mathcal{S}^* := (\text{Setup}, \text{KeyGen}, \text{Sign}^*, \text{Verify}^*)$  with at most  $q$  signature queries and at least  $\varepsilon$  advantage. We need to construct a PPT adversary





**Fig. 20.** An example of the steganographically-generated bitcoin transaction

$\mathcal{B}$  who can break the  $\mathcal{S} := (\text{Setup}, \text{KeyGen}, \text{Sign}, \text{Verify})$ . During the reduction game, upon receiving query  $m_i$  from  $\mathcal{A}$ ,  $\mathcal{B}$  computes  $m_i^* := \text{hash}(m_i, s_i)$  with a randomly sampled  $s_i \in \{0, 1\}^\lambda$ , and then query the signing oracle with  $m_i^*$ . After getting  $\sigma_i$  back from the signing oracle,  $\mathcal{B}$  return  $\sigma_i^* := (\sigma_i, s_i)$  to  $\mathcal{A}$ . Whenever  $\mathcal{A}$  can provide a valid pair  $(\hat{m}, \hat{\sigma})$  as forgery ( $\hat{\sigma} = (\sigma', s')$ ), since the probability that  $\mathcal{A}$  can find a collision on  $\text{hash}$  is negligible,  $\mathcal{B}$  can compute  $m' := \text{hash}(\hat{m}, s')$ . By definition,  $(\hat{\sigma}, m')$  is a forgery against  $\mathcal{S}$ .

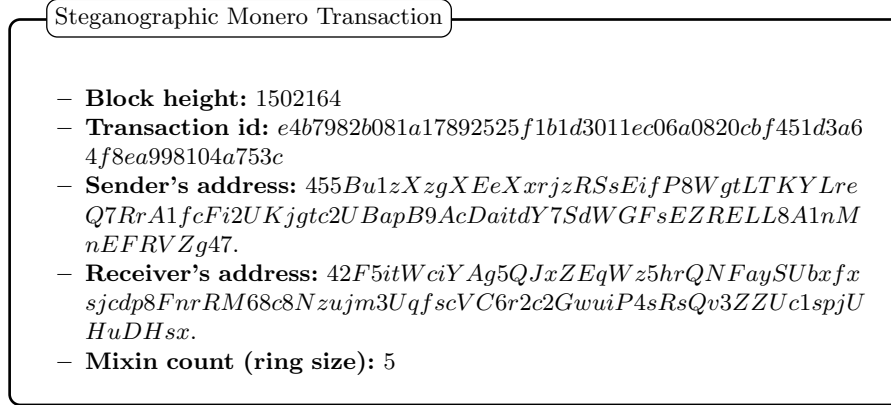
## B Bitcoin Pseudo-code and Demo Transaction

The attack described in Sec. 4 follows the pseudo-code in Fig. 9. This attack has been executed, and the transaction with the attributes shown in Fig. 20 was generated. To demonstrate how the hidden message is embedded and extracted we provide an extraction tool that can be downloaded and tested from: [https://github.com/NaLancaster/hash\\_and\\_extract.git](https://github.com/NaLancaster/hash_and_extract.git). The repository also contains the actual transaction binary in tx.txt and includes a pair  $(c, r)$  of random numbers containing a 16-byte hidden message, where  $(c, r)$  is found in cr.txt. The transaction hash in Fig. 20 can also be seen in any Bitcoin explorer, like <https://minergate.com/blockchain/bcn/blocks>, and the provided transaction binary should hash to the same hash value.

## C Detailed Implementation of Steganographic Attack in Monero

The following is the details of our implementation of the generic steganographic attack in Monero (v0.12.0.0) and older).

*Step 1: embedding a hidden message  $m$  and generating a signature that contains  $st$ .* The sender's wallet is modified to surreptitiously embed a 32-byte message  $m$



**Fig. 21.** An example of the steganographically-generated Monero transaction

in the randomly generated  $s_{i,j}$  numbers as part of the Borromean ring signature. Specifically, two vectors of  $s_{i,j}$  numbers are generated by the `genBorromean()` function:  $s_{0,j}$  and  $s_{1,j}$ . In addition,  $s_{0,j}$ 's are randomly generated when the  $j^{\text{th}}$  bit commitment is 1, and two of these randomly generated  $s_{0,j}$ 's are used in our attack to embed the stegotext  $st$ . For simplicity, we use  $s_{0,1}$  and  $s_{0,2}$  to denote the first two randomly generated numbers in  $s_{0,j}$  vector, although they might not necessarily correspond to  $j = 1$  and  $j = 2$  respectively.

Fig. 10 shows the two subverted numbers, in which  $s_{0,1}$  includes 16 bytes of random IV concatenated with 1 byte representing the index of  $s_{0,2}$  and 15 bytes of zeroes, where the last 16 bytes are sent encrypted using AES-CBC. The second subverted random number,  $s_{0,2}$ , contains hidden message  $m$  encrypted using AES-CBC under the key  $z$ .

This step of the attack is achieved by slightly modifying two functions: `genBorromean()` and `skGen()` in two files: `rctSig.cpp` and `rctOps.cpp`. `genBorromean()` is modified to pass two extra parameters to `skGen()`. The first parameter is the counter that indicates which of the two random numbers is to be generated, while the second parameter represents the index of  $j^{\text{th}}$  bit that corresponds to the second number  $s_{0,2}$  within the  $s_{0,j}$  vector. When the value of the counter is 0 or 1, `skGen()` generates random numbers according to Fig. 10, otherwise executes as normal.

*Step 2: identifying signature containing stegotext  $st$ , and extracting hidden message  $m$ .* To identify transactions containing  $st$ , the source `blockchain.cpp` file is modified to check the randomness within each new transaction and identify signatures containing stegotext. The receiver tests each number in the  $s_{0,j}$  vector by looking for a random IV that decrypts the second half of the tested number to a similar pattern as  $s_{0,1}$  in Fig. 10. Once this pattern is detected, the receiver concludes that this signature contains  $st$  and retrieves the index of  $s_{0,2}$  from the  $16^{\text{th}}$  byte of  $s_{0,1}$ .

**Table 3.** Number of signatures needed to leak bits of the long-term private key in our rejection-sampling ECDSA subversion attack

Exp. #	Number of Signatures to Leak Key Bits						
	32	64	96	128	160	192	224
1	35	71	117	173	227	314	520
2	33	73	119	189	253	344	485
3	32	74	120	170	230	339	471
4	34	70	114	179	238	335	491
5	32	69	119	189	280	385	552
6	32	76	122	170	233	333	526
7	36	76	120	180	262	400	576
8	32	71	127	197	259	368	566
9	33	72	115	162	242	348	528
10	32	71	121	177	243	363	498
11	31	70	121	180	246	345	524
12	35	76	120	181	260	386	563
13	33	69	124	190	251	352	506
14	32	72	121	180	255	353	518
15	33	78	124	178	246	355	539
16	34	72	113	168	232	331	522
17	35	72	111	162	228	340	512
18	31	79	125	201	281	378	544
19	37	76	122	174	243	329	475
20	34	75	121	175	260	369	570
<b>Average</b>	33.3	73.1	119.8	178.8	248.4	353.4	524.3
<b>Std. dev.</b>	1.6	2.9	3.9	10.2	15.2	21.6	30.5

When a malicious signature is detected, the receiver retrieves the index of  $s_{0,2}$  as above. The receiver then extracts the hidden message by decrypting  $s_{0,2}$  using his key  $z$  with AES-CBC. Fig. 21 shows a Monero transaction that has been steganographically subverted by our attack, and has been successfully posted to the Monero blockchain

## D ECDSA-Signature Rejection-Sampling Experiment

Table 3 illustrates the experimental results on how many signatures needed to obtain 32, 64, 96, 128, 160, 192 and 224 bits out of the total 256 key bits. This experiment was run 20 times to record the number of needed signatures to leak some bits of the secret key. As seen in Table 3, the average number of signatures that should be intercepted by the attacker to retrieve 50% of the key, i.e. 128 bits, is about 179 signatures.

## E Signature subversion

In theory, the **Setup**, **KeyGen**, **Sign** algorithms of a signature scheme can be subverted to leak secret information. However, in practice most blockchain platforms do not generate the setup parameters themselves; instead, widely trusted setup parameters, such as in ED25519, are adopted. Therefore, we don't consider **Setup** algorithm in this work. In terms of **KeyGen** algorithms, they are usually based on some one-way function, and it is possible to leak  $O(\log \lambda)$  bits through rejective sampling. Nevertheless, for most signature schemes, this would not be sufficient to allow the adversary to forge a signature. See *leakage resilient signatures* in [62, 63] for more discussion. Therefore, this work focuses on the subversion of the **Sign** algorithm. As a result, we adopt the following modified definition of undetectability from [47].

**Public/Secret Undetectability.** The undetectability is used to model the fact that normal users cannot distinguish if a signature is produced by a subverted signing algorithm or the genuine one.

**Definition 3.** Let  $\mathcal{S} = (\text{Setup}, \text{KeyGen}, \text{Sign}, \text{Verify})$  be a signature scheme. Let  $\mathcal{M}$  be the message space. We say a subverted  $\text{Sign}^*$  algorithm is secretly undetectable w.r.t.  $\mathcal{S}$  if for all PPT adversary  $\mathcal{A}$  we have any  $\{(\text{PK}_i, \text{SK}_i)\}_{i=1}^n$  output by  $\text{KeyGen}(\text{param})$  for any integer  $\lambda \in \mathbb{N}$ , any  $n = \text{poly}(\lambda)$ , any  $\text{param} \leftarrow \text{Setup}(1^\lambda)$ , any  $\{(\text{PK}_i, \text{SK}_i)\}_{i=1}^n$  output by  $\text{KeyGen}(\text{param})$ , and any  $\ell \in [n]$ , we have:

$$\text{Adv}_{\mathcal{A}}^{\text{SU}}(1^\lambda) = \left| \Pr \left[ \text{Expt}_{\mathcal{A}}^{\text{SU}}(1^\lambda) \right] - \frac{1}{2} \right| = \text{negl}(\lambda)$$

w.r.t. the following game/experiment:

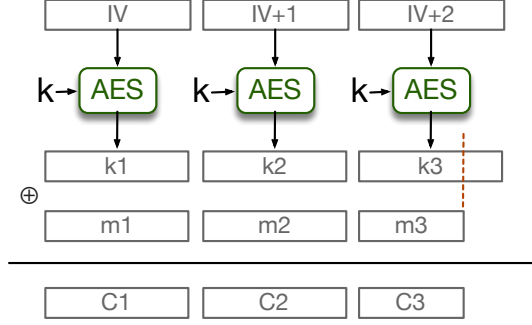
$\text{Expt}_{\mathcal{A}}^{\text{SU}}(1^\lambda)$

1. Pick  $b \leftarrow \{0, 1\}$ ;
2. Send  $(\{\text{PK}_i\}_{i=1}^n, \text{SK}_\ell)$  to  $\mathcal{A}$ ;
3. **For**  $j \in \{1, \dots, k\}$ ,  $\mathcal{A}$  queries  $m_j \in \mathcal{M}$  and obtains
  - $\sigma_j \leftarrow \text{Sign}(\{\text{PK}_i\}_{i=1}^n, \text{SK}_\ell, \ell, m_j)$  if  $b = 0$ ;
  - $\sigma_j^* \leftarrow \text{Sign}^*(\{\text{PK}_i\}_{i=1}^n, \text{SK}_\ell, \ell, m_j)$  if  $b = 1$ ;
4.  $\mathcal{A}$  outputs a bit  $b'$ ;
6. **Return**  $b \stackrel{?}{=} b'$ ;

We say a subverted  $\text{Sign}^*$  algorithm is publicly undetectable w.r.t.  $\mathcal{S}$  if in step 1 of the above game  $\mathcal{A}$  only receives  $\{\text{PK}_i\}_{i=1}^n$ .

**Reverse firewall based stego-resistance.** One of our countermeasures utilize miners as active online watchdogs, a.k.a. reverse firewalls. The reverse firewall based stego-resistance definition is slightly different from the above one. More specifically, we say a signature scheme  $\mathcal{S} = (\text{Setup}, \text{KeyGen}, \text{Sign}, \text{Verify}, \text{Wrap})$  is reverse firewall based stego-resistant if for all PPT adversary and  $\text{Sign}^*$  we have

$$\text{Adv}_{\mathcal{A}}^{\text{SU}}(1^\lambda) = \left| \Pr \left[ \text{Expt}_{\mathcal{A}}^{\text{SR}}(1^\lambda) \right] - \frac{1}{2} \right| = \text{negl}(\lambda)$$



**Fig. 22.** Ciphertext Stealing (CTR mode)

w.r.t. the following game/experiment:

**Expt** $_{\mathcal{A}}^{\text{SR}}(1^\lambda)$

1. Pick  $b \leftarrow \{0, 1\}$ ;
2. Send  $(\{\text{PK}_i\}_{i=1}^n, \text{SK}_\ell)$  to  $\mathcal{A}$ ;
3. **For**  $j \in \{1, \dots, k\}$ ,  $\mathcal{A}$  queries  $m_j \in \mathcal{M}$ :
  - o If  $b = 0$ :
    - Compute  $\sigma_j \leftarrow \text{Sign}(\{\text{PK}_i\}_{i=1}^n, \text{SK}_\ell, \ell, m_j)$ ;
    - Return  $\tilde{\sigma}_j \leftarrow \text{Wrap}(\{\text{PK}_i\}_{i=1}^n, m_j, \sigma_j)$  to  $\mathcal{A}$ ;
  - o If  $b = 1$ :
    - Compute  $\sigma_j^* \leftarrow \text{Sign}^*(\{\text{PK}_i\}_{i=1}^n, \text{SK}_\ell, \ell, m_j)$ ;
    - Return  $\tilde{\sigma}_j^* \leftarrow \text{Wrap}(\{\text{PK}_i\}_{i=1}^n, m_j, \sigma_j^*)$  to  $\mathcal{A}$ ;
4.  $\mathcal{A}$  outputs a bit  $b'$ ;
6. **Return**  $b \stackrel{?}{=} b'$ ;

## F Ciphertext Stealing Technique

In our attack, the leaked information is encrypted by a semantically secure symmetric-key encryption scheme. To minimize the number of lines of code to be changed, we need to adopt a readily implemented encryption algorithm. In our experiment, both Bytecoin and Monero wallets already have AES-128 algorithm, which; therefore, can be used as a building block of the semantically secure encryption. However, the message length is usually not a perfect multiple of 128 bits. To maximize the subversion channel capacity, one option is to adopt the concept of Ciphertext Stealing (CTS) [64]. For any given message length, with ciphertext stealing techniques, the ciphertext length is exactly the same as the message length (besides the IV). Our generic attack, in Sec. 3, uses CTR-mode based ciphertext stealing technique, where the last encryption block is truncated to fit the message length. The encryption algorithm  $\text{CTS-Enc}_k(\text{IV}, m)$  is depicted in Fig. 22. We refer interested readers to [64] for more operation modes with CTS, such as CBC.

## G Summary of Existing Countermeasures

In the following we summarize and discuss the current practices and techniques that can deter arbitrary content insertion on blockchains.

### G.1 ASA-Resistant Techniques

**Deterministic Cryptographic Primitives.** As identified by many researchers [29, 30, 32, 33], the root of such class of subversion attacks is the *uncontrolled randomness*. As such, the first intuitive solution is to use deterministic signatures. However, this approach has limited usage in the blockchain context.

**Signature Synthetic Randomness.** In a similar manner to deterministic DSA proposed by RFC 6979 [42], signatures can be made deterministic using synthetic randomness. More specifically, assume a signature algorithm consumes  $\ell$  random coins, denoted as  $r_1, \dots, r_n$ . Without loss of generality, suppose the signing algorithm takes as input the signing key  $s$  and the message  $m$ . We can generate the needed random coins deterministically as  $r_i := \text{hash}(s, m, i)$ . Based on heuristics property and onewayness of the hash function,  $r_i$  is unpredictable due to the entropy of  $s$ . On the other hand, this tweak allows offline watchdogs (verification algorithms) to compare and test an implementation with its specification.

Note that no *probabilistic polynomial time* black-box verification mechanism can ensure an implementation exactly matches its specification. This is because a malicious functionality may be triggered by a specific input, and it is impossible to verify that an implementation behaves as expected for all inputs. For instance, our attack can be modified so that the signing algorithm behaves honestly for all the inputs, except when the input message  $m = m^*$  the signing algorithm switches to our attack version, where  $m^*$  is the hidden trigger that has high entropy. The elimination of such hidden triggers is discussed below.

Instead, the offline implementation verification is only required to check polynomially many randomly sampled inputs (together with randomly sampled explicit randomness) and compare the corresponding outputs of the implementation with its specification. The most recommended approach to achieve automatic verification is to use so-called executable specifications [65]. We emphasize that the offline implementation verification algorithm must be trusted and certified. Nevertheless, it is universal and only needs simple comparison functionality, which makes it easy to ensure subversion freeness.

**Verifiable Random Function (VRF).** VRFs, proposed by Micali et al. [43], are random functions that non-interactively prove the correctness of their outputs. However, due to their randomness, VRF outputs are still susceptible to rejection sampling attacks. Hence, although their use may decrease the steganography throughput, VRFs can not completely deter steganography and subversion attacks in the context of public blockchains. Similarly, the work presented in [66] is vulnerable to rejection-sampling subversion attacks.

**Split Model.** The idea was initially proposed by Russell et al. in [44] to clip the power of subversion by hashing the output of randomized primitives. This idea

was developed further in subsequent work [45, 46] to decompose each randomness-generation function  $\text{RG}$  into two separate components  $\text{RG}_0$  and  $\text{RG}_1$ , execute them independently, and compose their output using a deterministic function  $\phi$  to generate the final random number. The authors of [46] emphasize that practical implementation of such splitting necessitates the independent execution of the separate components, which can be achieved by executing them in separate environments, e.g. by running them in different virtual machines or containers like Docker [67].

## G.2 Preventative Trust-Based Countermeasures

**Reverse Firewalls (RF).** Ateniese et al. [47] proposed the idea of a trusted un-tamperable reverse firewall (RF). RF is an algorithm that, using public information, re-randomizes, and hence sanitizes, that output of possibly subverted randomized signature algorithms. Although this technique is effective against subversion attacks on randomized signatures; however, it requires an active trusted firewall and works only on re-randomizable signatures. Therefore, it might not be suitable for most signature schemes in the context of blockchains.

**Self-Guarding Protocols.** Recently Fischlin and Mazaheri [48] have proposed a novel technique that proactively defends against ASA’s assuming *temporary trust*, i.e. ASA happens after a period of an honest initial phase. Namely, they provided constructions for homomorphic public-key encryption, symmetric-key encryption, signature schemes and PUF-based key exchange. In general, their constructions are divided into two phases; a sampling phase, and a challenge phase. In the sampling phase, or honest initial phase, a sample of ciphers, in the case of encryption, or signatures, in the case of signature schemes, is honestly generated. This sample is stored and used in the second phase to obfuscate ciphers/signatures and detect possible ASA attacks.

To further illustrate this technique, Fig. 23 presents a simplified pseudo-code for the construction of a self-guarding signature scheme  $\mathcal{S}^{sg} = (\text{KeyGen}^{sg}, \text{Sign}^{sg}, \text{Verify}^{sg})$  from a deterministic signature scheme  $\mathcal{S} = (\text{KeyGen}, \text{Sign}, \text{Verify})$  [48]. As seen in this figure, the key generation algorithm  $\text{KeyGen}^{sg}$  is given a security parameter  $(1^\lambda)$ , generates a list of  $\lambda$  key pairs  $(sk_i, pk_i)$ , and sets the private key  $sk^{sg} := (sk_1, \dots, sk_\lambda)$  and the public key  $pk^{sg} := (pk_1, \dots, pk_\lambda)$ . As part of the *trusted* sampling phase,  $\mathcal{S}^{sg}.\text{Sample}(sk^{sg})$  is executed to generate a queue of  $\lambda$  pairs of randomly generated messages  $m_{r,i}$  and their corresponding signatures  $\sigma_{r,i}$ . When signing a message  $m$  with the *possibly subverted* algorithm  $\mathcal{S}.\text{Sign}$ , the self-guarding signing algorithm  $\mathcal{S}^{sg}.\text{Sign}$  does the following. For each  $i = 1, \dots, \lambda$ , it randomly picks  $b_i \xleftarrow{\$} \{0, 1\}$ , and execute the original signing algorithm  $\mathcal{S}.\text{Sign}$  twice, once to sign  $m_{r,i}$  and another time to sign  $(m_{r,i} \oplus [m || \sigma_{r,i}])$ . The order of which is signed first is determined by the value of  $b_i$ , and each time the generated signature for  $m_{r,i}$  is compared with its previously and *honestly* generated signature  $\sigma_{r,i}$ . This is done to ensure detecting any possible subversion with probability close to  $1/2$ . If any discrepancy is detected, abort, otherwise return the signature  $\sigma := (m_{r,1}, \sigma_{r,1}, \sigma_1, \dots, m_{r,\lambda}, \sigma_{r,\lambda}, \sigma_\lambda)$ . The verification  $\mathcal{S}^{sg}.\text{Verify}$  is

$\mathcal{S}^{sg}(\text{KeyGen, Sample, Sign, Verify})$

$\mathcal{S}^{sg}.\text{KeyGen}^{sg}(1^\lambda)$ :

- For  $i = 1 \dots \lambda$ :  $(sk_i, pk_i) \xleftarrow{\$} \mathcal{S}.\text{KeyGen}(1^\lambda)$ ;
- $(sk^{sg}, pk^{sg}) := ((sk_1, \dots, sk_\lambda), (pk_1, \dots, pk_\lambda))$ ;
- Return  $(sk^{sg}, pk^{sg})$ ;

$\mathcal{S}^{sg}.\text{Sample}(sk^{sg})$ :

- Empty  $\mathcal{Q}$ ;
- For  $i = 1 \dots \lambda$ :  
 $m_{r,i} \xleftarrow{\$} \{0, 1\}^\ell$ ;  $\sigma_{r,i} := \mathcal{S}.\text{Sign}(sk_i, m_{r,i})$ ;  
 Store  $(m_{r,i}, \sigma_{r,i})$  in  $\mathcal{Q}$ ;
- Return  $\mathcal{Q}$ ;

$\mathcal{S}^{sg}.\text{Sign}^{sg}(sk^{sg}, m, \mathcal{Q})$ :

- if  $\mathcal{Q}$  is empty: Return  $\perp$ ;
- For  $i = 1 \dots \lambda$ :  
 Retrieve  $(m_{r,i}, \sigma_{r,i})$  from  $\mathcal{Q}$ ;  
 $b_i \xleftarrow{\$} \{0, 1\}$ ;  
 if  $(b = 0)$ :  $(m^0, m^1) := (m_{r,i}, m_{r,i} \oplus [m || \sigma_{r,i}])$ ;  
 else:  $(m^0, m^1) := (m_{r,i} \oplus [m || \sigma_{r,i}], m_{r,i})$ ;  
 $\sigma^0 := \mathcal{S}.\text{Sign}(sk_i, m^0)$ ;  $\sigma^1 := \mathcal{S}.\text{Sign}(sk_i, m^1)$ ;  
 if  $(\sigma^{b_i} \neq \sigma_{r,i})$ : Return  $\perp$ ;  
 $\sigma_i := \sigma^{(1-b_i)}$ ;
- $\sigma := (m_{r,1}, \sigma_{r,1}, \sigma_1, \dots, m_{r,\lambda}, \sigma_{r,\lambda}, \sigma_\lambda)$ ;
- Return  $\sigma$ ;

$\mathcal{S}^{sg}.\text{Verify}^{sg}(pk^{sg}, m, \sigma)$ :

- result := true;
- For  $i = 1 \dots \lambda$ :  
 result := result  $\wedge$   $\mathcal{S}.\text{Verify}(pk_i, m_{r,i}, \sigma_{r,i})$ ;  
 result := result  $\wedge$   $\mathcal{S}.\text{Verify}(pk_i, (m_{r,i} \oplus [m || \sigma_{r,i}]), \sigma_i)$ ;
- Return result;

**Fig. 23.** Self-guarding signature  $\mathcal{S}^{sg}$  based on a signature scheme  $\mathcal{S} = (\text{KeyGen, Sign, Verify})$ .  $\mathcal{Q}$  represents a queue of pairs of random messages  $m_{r,i}$  and their honestly-generated signatures  $\sigma_{r,i}$ , and  $\ell$  is the message and signature space length. Other sanity checks have been omitted for simplicity. For more details refer to [48].

carried out by re-constructing the string  $(m_{r,i} \oplus [m || \sigma_{r,i}])$  for each  $i = 1, \dots, \lambda$ , and calling the original verification algorithm  $\mathcal{S}.\text{Verify}$  twice, to verify each  $\sigma_i$  and  $\sigma_{r,i}$ . If all signatures are valid, return true, and false otherwise.



### G.3 Blockchain-based techniques

**Light Blockchains.** To solve issues related with blockchain size and scalability, new blockchain designs have emerged. For example, PascalCoin [49] is a cryptocurrency that does not keep the full history of transactions but rather stores the last 100 blocks in its ledger, and actual account balances are stored in a another cryptographic structure called the SafeBox. A very similar approach is used in the mini-blockchain scheme [50] that is implemented by Cryptonite [68] which stores the actual balances in a structure called the *account tree* which is updated by the transactions in the blockchain. Because new transactions reference the *account tree* and not previous transactions in the blockchain, transactions in older blocks can be discarded. Note, older block headers are still kept in the mini-blockchain. Although these solutions are mainly proposed solve scalability issues, these new designs can deter permanent storage of malicious content.

**Redactable Blockchains.** Redactable blockchains have been proposed in [5] to rewrite, remove, and insert new blocks in the blockchain. In their technique, which is based on the use of Chameleon hashes [69], the redaction could be performed by a trusted central node, or a group of nodes who posses the Chameleon hash trapdoor. Similarly,  $\mu$ chain [4] proposes a mutable blockchain that is based on consensus. Hence, if malicious content is identified, mutable blockchains can effectively deter the persistent storage of such content on the blockchain.

**Content Filters.** *Content filters* target human readable strings to detect and reject unwanted content, e.g. rejecting a transaction if its 20-Byte destination address has 18 printable characters [3].

**Increasing Transaction Fees.** Although increasing the transaction fees is not advisable for promoting blockchain among innocent users, and can unfairly penalize users who relay on large transactions, e.g. exchange services, minimum mandatory fees has been proposed as a countermeasure in [3] to render content insertion economically infeasible for large transactions.

**Self-verifying Addresses.** The goal of this technique is to deter content insertion in Bitcoin by using arbitrary addresses. The authors of [3] suggested that instead of sending an address  $a$ ,  $c_a$  is sent in the transaction, where  $c_a = (G^a, r, \text{Sign}(G^a || r, a))$ ,  $r = \text{CRC32}(t_1 || \dots || t_i)$ , and  $t_i$  is the transaction corresponding to the  $i^{\text{th}}$  input. A similar approach is to limit the address Space. For example, PascalCoin [49] has a finite address space, and accounts are limited but can be associated with any public key. Although may not be intended to stop content-insertion, this practice can deter the arbitrary manipulation of transactions' addresses.

## H Non-interactive Zero-knowledge

Here we briefly introduce non-interactive zero-knowledge (NIZK) schemes in the Random Oracle (RO) model. Let  $\mathcal{R}$  be an efficiently computable binary relation.

For pairs  $(x, w) \in \mathcal{R}$  we call  $x$  the statement and  $w$  the witness. Let  $\mathcal{L}_{\mathcal{R}}$  be the NP language consisting of statements in  $\mathcal{R}$ , i.e.  $\mathcal{L}_{\mathcal{R}} = \{x \mid \exists w \text{ s.t. } (x, w) \in \mathcal{R}\}$ . An NIZK scheme includes following algorithms: a PPT algorithm  $\text{Prov}$  that takes as input  $(x, w) \in \mathcal{R}$  and outputs a proof  $\pi$ ; a polynomial time algorithm  $\text{Verify}$  takes as input  $(x, \pi)$  and outputs 1 if the proof is valid and 0 otherwise.

**Definition 4 (NIZK Proof of Membership in the RO Model).**  $\text{NIZK}_{\mathcal{R}}^{\text{RO}}.\{\text{Prov}, \text{Verify}, \text{Sim}, \text{Ext}\}$  is an NIZK Proof of Membership scheme for the relation  $\mathcal{R}$  if the following properties hold:

- *Completeness:* For any  $(x, w) \in \mathcal{R}$ ,

$$\Pr \left[ \zeta \leftarrow \{0, 1\}^\lambda; \pi \leftarrow \text{Prov}^{\text{RO}}(x, w; \zeta) : \text{Verify}^{\text{RO}}(x, \pi) = 1 \right] \geq 1 - \text{negl}(\lambda).$$

- *Zero-knowledge:* If for any PPT distinguisher  $\mathcal{A}$  we have

$$\left| \Pr[\mathcal{A}^{\text{RO}, \mathcal{O}_1}(1^\lambda) = 1] - \Pr[\mathcal{A}^{\text{RO}, \mathcal{O}_2}(1^\lambda) = 1] \right| \leq \text{negl}(\lambda).$$

The oracles are defined as follows:  $\mathcal{O}_1$  on query  $(x, w) \in \mathcal{R}$  returns  $\pi$ , where  $(\pi, aux) \leftarrow \text{Sim}^{\text{RO}}(x)$ ;  $\mathcal{O}_2$  on query  $(x, w) \in \mathcal{R}$  returns  $\pi$ , where  $\pi \leftarrow \text{Prov}^{\text{RO}}(x, w; \zeta)$  and  $\zeta \leftarrow \{0, 1\}^\lambda$ .

- *Soundness:* For all PPT adversary  $\mathcal{A}$ ,

$$\Pr \left[ \begin{array}{c} (x, \pi) \leftarrow \mathcal{A}^{\text{RO}}(1^\lambda) : \\ x \notin \mathcal{L}_{\mathcal{R}} \wedge \text{Verify}^{\text{RO}}(x, \pi) = 1 \end{array} \right] \leq \text{negl}(\lambda).$$

**Definition 5 (NIZK Proof of Knowledge in the RO Model).**  $\text{NIZK}_{\mathcal{R}}^{\text{RO}}.\{\text{Prov}, \text{Verify}, \text{Sim}, \text{Ext}\}$  is an NIZK Proof of Knowledge scheme for the relation  $\mathcal{R}$  if the completeness, zero-knowledge, and extraction properties hold, where the extraction is defined as follows.

- *Extractability:* For all PPT adversary  $\mathcal{A}$ ,

$$\Pr \left[ \begin{array}{c} (x, \pi) \leftarrow \mathcal{A}^{\text{RO}}(1^\lambda); \\ w \leftarrow \text{Ext}^{\text{RO}}(x, \pi) : \\ \text{Verify}^{\text{RO}}(x, \pi) = 1 \vee (x, w) \in \mathcal{R} \end{array} \right] \geq 1 - \text{negl}(\lambda).$$

We need non-interactive zero-knowledge proofs/arguments of knowledge and non-interactive zero-knowledge proofs/arguments of membership. For simplicity, we will drop RO from the superscript if the context is clear.

We use  $\text{NIZK}_{\mathcal{R}_i}^{\text{RO}}.\text{Verify}$  and  $\text{NIZK}_{\mathcal{R}_i}^{\text{RO}}.\text{Sim}$  to denote the corresponding verification algorithm and simulator, respectively.