

# Cryptanalysis of Curl-P and Other Attacks on the IOTA Cryptocurrency

Ethan Heilman<sup>1</sup>, Neha Narula<sup>2</sup>, Garrett Tanzer<sup>3</sup>, James Lovejoy<sup>2</sup>, Michael Colavita<sup>3</sup>, Madars Virza<sup>2</sup> and Tadge Dryja<sup>2</sup>

<sup>1</sup> Boston University, Boston, Massachusetts, United States of America

[heilman@bu.edu](mailto:heilman@bu.edu)

<sup>2</sup> Massachusetts Institute of Technology Media Lab, Cambridge, Massachusetts, United States of America

[madars@mit.edu](mailto:madars@mit.edu), [tdryja@media.mit.edu](mailto:tdryja@media.mit.edu), [jlovejoy@mit.edu](mailto:jlovejoy@mit.edu), [narula@media.mit.edu](mailto:narula@media.mit.edu)

<sup>3</sup> Harvard University, Cambridge, Massachusetts, United States of America

[gtanzer@college.harvard.edu](mailto:gtanzer@college.harvard.edu), [mcolavita@college.harvard.edu](mailto:mcolavita@college.harvard.edu)

**Abstract.** We present attacks on the cryptography formerly used in the IOTA blockchain, including under certain conditions the ability to forge signatures. We developed practical attacks on IOTA’s cryptographic hash function Curl-P-27, allowing us to quickly generate short colliding messages. These collisions work even for messages of the same length. Exploiting these weaknesses in Curl-P-27, we broke the EU-CMA security of the former IOTA Signature Scheme (ISS). Finally, we show that in a chosen-message setting we could forge signatures and multi-signatures of valid spending transactions (called bundles in IOTA).

**Keywords:** cryptocurrencies · signature forgeries · cryptographic hash functions · cryptanalysis

## 1 Introduction

This paper presents attacks on the signature scheme used to authorize payments in a cryptocurrency known as IOTA. IOTA is a cryptocurrency designed for use in the Internet of Things (IoT) and automotive ecosystems. As of February 20, 2020, it had a market capitalization of \$736 million US dollars [Coi19]. The attacks we describe here work by exploiting a cryptographic weakness in IOTA’s hash function, Curl-P-27. Importantly, our attacks were disclosed and patched in August 2017, and thus no longer impact the security of IOTA’s signature scheme [HNDV17].

IOTA uses cryptographic signatures to authorize payments by users. The IOTA Signature Scheme (ISS) is based on Winternitz One-Time Signatures [Mer89], but unlike traditional Winternitz, in IOTA users sign the hash of a message. Thus, the security of ISS relies on its cryptographic hash function, which was Curl-P-27. Using a differential cryptanalysis attack, we are able to quickly create messages of the same length which hash to the same value with Curl-P-27, breaking the function’s collision resistance. We find an upper bound on the average number of queries to Curl-P-27 to generate a collision.

Using this collision attack, we can generate signature forgeries in IOTA. Our attacks on the IOTA signature scheme function in a *chosen-message* setting, where an attacker creates two payments—a benign payment and a malicious payment—such that a signature on the benign payment is also a valid signature on the malicious payment. Our analysis is just on the IOTA signature scheme and does not include the security of the IOTA network as a whole. These attacks apply to both normal and multi-signature IOTA payments. Spending

from a multi-signature address requires one user to produce a payment for another user to sign, which fits exactly in the chosen-message setting of our attack. We detail how to apply our attack to IOTA payments which spend from multi-signature addresses, and provide a tool for creating collisions in single-signature and multi-signature IOTA payments. We also evaluate the resources required to perform the attack, and show that using 80 cores, we can create colliding IOTA payments in less than twenty seconds on average. We have open sourced and published the software used in these attacks.<sup>1</sup>

We follow up our signature forgery attacks with an analysis of the *known-message-attack* security of the IOTA Signature Scheme. We examine the IOTA Signature Scheme's susceptibility to generic attacks on the underlying hash function. The result of this analysis is an upper bound on the resistance of the IOTA Signature Scheme to these attacks. The security bounds given in this section do not present a risk of *known-message attacks*. However, they do show a reduction in the security parameters of ISS to generic attacks.

A chief contribution of this paper is that it is a real-world case study on how Winternitz One-Time Signatures-like schemes fail when the underlying cryptographic hash function succumbs to cryptanalysis. For instance the literature on Winternitz One-Time Signatures *e.g.*, [BDE<sup>+</sup>11], often analyzes its security properties under the assumption that the signatures will be performed directly on the message itself. However the example of ISS provides strong evidence that in practice implementers of such schemes are unlikely to sign a message directly. This is because signing the message directly results in a signature length which is proportional to the message length. Instead, as is the case with ISS, implementers will sign the hash of the message. This makes the collision resistance of the underlying hash function far more relevant to the security of the signature scheme. Additionally this paper must contend with the limitations and realities of using a specific set of cryptanalytic weaknesses to create signature forgeries resulting in validly formatted, but semantically different, cryptocurrency payments.

## 1.1 Vulnerability Status and Impact

On July 14, 2017, some of the authors began a disclosure process with the IOTA developers. We negotiated a timeline for them to patch the vulnerability and a date after which we could publish our results. On August 7, 2017, the IOTA developers deployed a backwards-incompatible upgrade to mitigate this vulnerability by removing the use of Curl-P-27 to generate signatures in IOTA, and replacing it with another hash function [Søn17]. In order to perform the upgrade, deposits and withdrawals were halted on Bitfinex for approximately three days [Bit17]. All users who held IOTA directly (not via an exchange) were encouraged to upgrade their wallets and addresses. On September 7, 2017 we published our vulnerability report describing the nature of our attack [HNDV17].

Our vulnerability report included example Curl-P-27 collisions and signature forgeries on validly formatted IOTA payment messages, as well as software to validate these examples. In this paper, we expand upon our previous results by detailing the cryptanalytic techniques we used to break Curl-P-27's collision resistance and providing software to generate said signature forgeries. We also extend these techniques to develop an attack against IOTA's multi-signature scheme when Curl-P-27 is used—the multi-signature setting is particularly well-suited to chosen-message attacks. We did not send any of these forged signatures to the IOTA network or interfere in the IOTA network in any way. As Curl-P-27 is no longer used for ISS, the signature forgery attacks presented in this paper do not impact present-day IOTA. This includes our multi-signature attack in Section 5.2. Curl-P-27 is still used in other parts of IOTA [IOT17]. We do not present attacks on these uses.

Our results in Section 6 are relevant to present day IOTA, but they have no immediate security impact for those who use the default security settings. We privately disclosed <sup>2</sup>

<sup>1</sup><https://github.com/mit-dci/tangled-curl>

<sup>2</sup>When disclosing our results on normalization IOTA co-founder Sergey Ivanchevlo disputed the novelty

them to the IOTA developers on Aug 8, 2018 and then publicly disclosed them March 31 2019.

## 2 Related Work

After the release of our initial report, Colavita and Tanzer [CT18] independently reproduced and implemented our cryptanalysis, as well as proved some new results about Curl-P’s round function—namely, that it is a permutation and that it diffuses differentials across rounds in accordance with a particular closed-form expression. They are now collaborating on this paper. [DRUH18] explores replay attacks against IOTA.

Differential cryptanalysis techniques were first published in 1991 by Biham and Shamir [BS91] (researchers at IBM had discovered similar techniques in 1974 but chose not to disclose them publicly [Cop94]). In this paper, we present a very simple application of differential cryptanalysis on a balanced ternary cryptographic hash function. Apart from our initial vulnerability report and [CT18] we are aware of no prior research on the differential cryptanalysis of balanced ternary-based cryptographic hash functions. However, there is work designing and analyzing a ternary-based cryptographic pseudo-random sequence generator [GJ05]. In response to our attacks on Curl-P, [KTDB19] proposed Troika, a ternary cryptographic hash function with differential and linear cryptanalysis resistance, intended as a drop-in replacement for Curl-P. [LI19] studies preimage attacks against a reduced round Troika.

Exploiting our cryptanalysis of Curl-P-27, we present a chosen-message attack on ISS’s unforgeability. Although the danger of broken collision resistance—and the chosen message attack model—may not be immediately apparent, we see a cautionary tale in the work on the MD5 hash function. In 2004, Wang et al. released the first complete collision for MD5 [WFLY04], and soon after published a generic procedure for generating random collisions [WY05]. In 2005, Lenstra joined Wang to apply this cryptographic vulnerability to X.509 certificates, a cornerstone of the public key infrastructure that enables protocols like HTTPS, and was able to construct pairs of colliding certificates [LWdW05]. Amidst doubts that a certificate authority would sign such suspicious certificates, or that they would even be exploitable once issued because they lacked “meaningful” structure, Stevens joined Lenstra et al. in 2007 to extend the original random collision attack on MD5 to a chosen-prefix collision attack [SLdW07]. This work culminated in 2009, when Stevens et al. announced that they had managed to forge a X.509 certificate with certificate authority privileges that passed verification on all major browsers [SSA<sup>+</sup>09], causing vendors to immediately obsolete MD5.

In October 2017, after the IOTA developers transitioned from using Curl-P-27 to using Kerl—based on Keccak—as the hash function in the IOTA Signature Scheme, an unrelated vulnerability called the 13 or M attack was discovered [Pin18]. This exploit relies on the fact that in IOTA’s signature scheme—which signs the message’s hash in chunks with values in  $[-13, 13]$ —a signature for the number 13 (also represented as ‘M’) reveals as plaintext a derivative of the private key that can be used to forge all subsequent chunks. The IOTA Foundation patched this vulnerability by requiring that if a message hash to be signed includes a 13, then the user must alter the message until no 13s are present in the digest. As an additional remediation step, the IOTA developers transferred potentially compromised funds to addresses under their control, providing a process for users to later apply to the IOTA Foundation in order to reclaim their funds [Rot18]. We also present the first analysis of the resistance of ISS to generic attacks. Our results in this area show that,

---

of our normalization analysis. He argues that based on an earlier statement he made in [blo18] that the weakness in the normalization scheme was both intentional and that they had told us of the existence of this weakness. As far we are aware there is no preceding work providing concrete and correct analysis of the degree to which normalization continues to weaken ISS (IOTA Signature Scheme).

because of a process IOTA performs on the hash of a message called normalization, the IOTA security scheme provides a lower level of security than claimed in public materials.

### 3 Background

In this section, we provide the necessary preliminaries to understand our attacks. We start with a short review of some of IOTA's uncommon design features and terminology. We also provide an overview of the Curl-P hash function and the IOTA Signature Scheme (ISS).

#### 3.1 IOTA Design

IOTA currently has several uncommon design features. First, IOTA uses balanced ternary instead of binary; second, payments in IOTA are known as bundles; third, IOTA uses a new data structure called a tangle rather than a traditional chain of blocks; and fourth, IOTA employs a trusted party called a coordinator to checkpoint the state and approve payments.

IOTA's data structures use balanced ternary, or base three; instead of bits in  $\{0, 1\}$ , it uses *trits* in  $\{-1, 0, 1\}$ , and instead of bytes of eight bits, it uses *trytes* of three trits. A tryte is represented as an integer in  $[-13, 13]$ . IOTA often serializes trytes as the letters A-Z and the number 9.

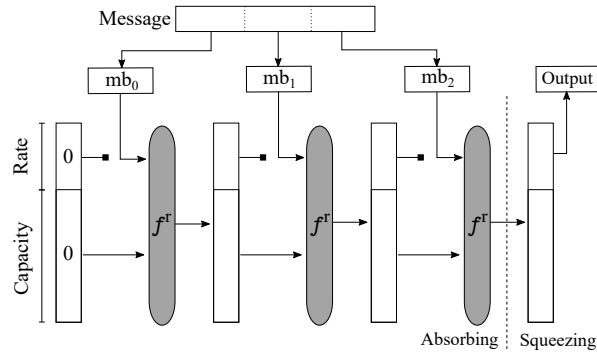
A payment in IOTA is represented by a data structure called a *bundle*. Bundles are composed of multiple transactions, but IOTA transactions are not like transactions in other cryptocurrencies; they are buffers which store inputs or outputs. IOTA transactions include, among others, address, signature, value, and tag fields. We provide a detailed description of the IOTA bundle and transaction format when describing our attacks in Section 5.

IOTA is built upon the concept of a *tangle* [Pop16]. This is similar to a Directed Acyclic Graph-chain (DAG-chain), where each block can reference more than one block parent [SZ15]. In IOTA's case, however, there are no blocks to aggregate multiple payments. Instead, each transaction must have a nonce proving Proof-of-Work (PoW) and include pointers to two other transactions. In order to add transactions to the tangle, a user selects two *tip* transactions from the tangle to reference in her transaction. Once created and signed, the user performs sufficient Proof-of-Work and broadcasts the transaction (or transactions, in the case of a bundle) to the IOTA network.

In IOTA as it is currently deployed, a bundle must also be approved by the coordinator to be accepted. The coordinator is a trusted party run by the IOTA developers that approves and checkpoints the state of the tangle by signing it. This has led to concerns that IOTA is centralized, or under the control of the IOTA developers [Wal17]. The IOTA developers argue IOTA is not centralized and that the coordinator is a temporary measure. The source code for the coordinator is not publicly available. Since we did not interact with the IOTA network we cannot confirm how the coordinator would impact our proposed attacks, but we are not aware of any mechanism in the coordinator that would prevent the attacks presented in this paper.

#### 3.2 IOTA's Signature Scheme (ISS)

IOTA uses a signature scheme inspired by the Winternitz One-Time Signatures (W-OTS) [Mer89]. W-OTS is an optimization of Lamport signatures [Lam79], operating on multiple bits (in IOTA, trits) at once to trade computational cost for decreased public key size.



**Figure 1:** The Curl-P construction.

ISS differs in several important aspects from W-OTS. First, ISS operates on the hashes of messages rather than on the messages directly, as in traditional W-OTS.<sup>3</sup> Second, rather than use a checksum, ISS performs a technique dubbed *normalization* on the hash of the message. As we will show in Section 6, the use of normalization instead of a checksum reduces the security parameters of ISS when compared with W-OTS.

ISS has three security levels. The first security level only signs the first third of the message hash. The second security level signs the first two thirds of the message hash. Finally, the third security level signs the entire message hash. Because our attack works against the highest security level, it also works against any of the lower security levels. For this reason, when we talk about ISS we will implicitly assume that security level three is used. In Section 6 we analyze the strength of the different security levels.

### 3.3 Curl-P

In this section, we describe the Curl-P hash function. Curl-P (sometimes referred to as Curl) is a cryptographic hash function designed specifically for use in IOTA. It has been used for a number of purposes in IOTA, including creating transaction addresses, creating message digests, Proof-of-Work (PoW), and hash-based signatures. At a high level, Curl-P follows the pattern of a Sponge Construction [BDPVA08, GJMG11], but it differs in some key areas. As the IOTA project has not provided any formal specification or analysis of Curl-P, we base our description on the open source implementation of Curl-P made available by the IOTA developers.

Unlike most cryptographic hash functions, Curl-P operates on trits in balanced ternary. For clarity, we represent individual trits with lowercase letters such as  $a, b, c, x, y, z$  and sequences of trits as uppercase letters such as  $S, N, X, Y$ , unless we are referring a particular trit within a sequence of trits, where we will use subscript notation such as  $S_i$ . Following IOTA's convention, the R in Curl-P-R denotes the number of rounds used (*e.g.*, Curl-P-27 denotes 27-round Curl-P).

As shown in Figure 1, Curl-P operates as follows: (1) Curl-P initializes an all zero state  $S$  of length 729 trits. (2) The message is broken into message blocks  $mb_0 \dots mb_n$ , each 243 trits. Curl-P employs no message padding; instead, if a message is not a multiple of 243 in length, the last message block is allowed to be less than 243 trits.<sup>4</sup> (3) In turn, each

<sup>3</sup>One can remark that modern signature schemes, like XMSS, LMS, and SPHINCS+ sign a *randomized* hash digest, and thereby achieve multi-target security, and relax conditions placed on the internal (compressing) hashes (*i.e.*, they remain secure if internal hashes are second preimage-resistant; down from collision-resistance). However, our attack is not against the W-OTS internals, as used in ISS, but against this “outer” hash invocation, therefore even though ISS should have considered a randomized variant, it would still not prevent attacks if the digest hash is completely broken, as is the case here.

<sup>4</sup>In some implementations of Curl-P an error is thrown if the message length is not an even multiple of

message block  $mb_0 \cdots mb_n$  is copied into the first third of the state  $S$  and then that state  $S$  is transformed by the function  $f^r$ . (4) Finally, when no more message blocks remain Curl-P returns the first third of the final state as the hash output. For a more detailed description, see Algorithm 1.

---

**Algorithm 1:** The sponge-like construction used by Curl-P.

**Function CurlHash(msg):**

```

S ← {0}^{729};
for p ← 0; p < |msg|; p ← p + 243 do
  if p + 243 < |msg| then
    mb ← msg[p, p + 243];
  else
    mb ← msg[p, |msg| - 1];
  S[0, |mb|] ← mb;
  S ← f^r(S);
return S[0, 243];

```

---

Now let's turn our attention to the function  $f^r$ , which is used to transform the state  $S$ . The transformation function  $f^r$  is actually just the function  $f$  recursively called on the state  $S$  for  $r$  rounds, *e.g.*,  $f^3(S) = f(f(f(S)))$ . Curl-P-27 is the Curl-P hash function which uses  $f^{27}$  as its transformation function.

Each round of  $f^r$  generates a new state from the current state by calling  $f$ . As described in Algorithm 2, each trit in the new state is determined by applying the simple function  $g$  to a pair of trits in the current state. Each trit in the current state is used twice, once as the first parameter to  $g$  (represented by  $a$ ) and once as the second parameter (represented by  $b$ ). In Table 1, we give  $g$  as a substitution box or s-box.

**Table 1:** S-box used by Curl-P: takes two trits  $a$ ,  $b$  and returns a trit  $c$ .

$$c = g(a, b)$$

	$b = -1$	$b = 0$	$b = 1$
$a = -1$	1	1	-1
$a = 0$	0	-1	1
$a = 1$	-1	0	0

## 4 Cryptanalysis of Curl-P

In this section, we apply common differential cryptanalysis methods to engineer meaningful full-state collisions in Curl-P-27. Our attack constructs two messages of the same length which differ at only a single trit position and hash to the same value under Curl-P-27. Our technique lets us have a large degree of control over the content of the colliding messages, including arbitrary message prefixes and suffixes. In the next section, we will exploit this control over Curl-P-27 to forge signatures on valid IOTA payments.

We were unable to find a formal specification or documentation of Curl-P or Curl-P-27

---

243 preventing trivial collisions.

---

**Algorithm 2:** The transform function  $f(S)$  called by the Curl Hash Function.

```

i ← 0;
for pos ← 0; pos < 729; pos ← pos + 1 do
  j ← i;
  if i < 365 then
    i ← i + 364;
  else
    i ← i - 365;
  N[pos] ← g(Sj, Si);
return N;

```

---

beyond the source code published as part of the IOTA open source project<sup>5</sup>. Furthermore, in our correspondence with the IOTA developers, they have stated that Curl-P-27 is designed to collide for specific sets of inputs [blo18]. In fact, Curl-P-27 is clearly non-random. As explored in detail by [CT18], Curl-P-27’s non-random behavior can be observed in messages of the same length; collisions and second preimages are trivial to generate for messages of different lengths. Thus, to ensure we have truly broken Curl-P-27 we show that our collision attack meaningfully breaks a security property of Curl-P-27 on which the IOTA Signature Scheme (ISS) depends (see Section 5).

At a high level, our attack works as follows. We choose two messages of at least three message blocks in length which differ at only a single trit. To decrease the difficulty of our attack, we choose these messages such that they satisfy certain constraint equations (explained in Section 4.2). Once we arrive at the message blocks that differ between the two messages we need to ensure that a collision occurs. To do this, we randomly modify a set of trits in both messages. This set of altered trits is limited to the differing message block in each of the two messages. The idea is that after running the transform function  $f^{27}$  on the differing message blocks, the only differences are in the first third of the resultant states.

$$f^{27}(S)[243, 729] = f^{27}(S')[243, 729]$$

Because Curl-P replaces the first third of the state with the next message block, these differences are erased, causing a full state collision. We exploit the differential properties of Curl-P-27 to brute force a 1-trit difference for many of the rounds of the transformation function such that it is unlikely these differences will diffuse beyond the first third of the state by the final round. Finding two messages that maintain a 1-trit difference across a sufficient number of rounds to generate a collision is upper bounded by 7.6 million or  $2^{22.87}$  queries to Curl-P-27.

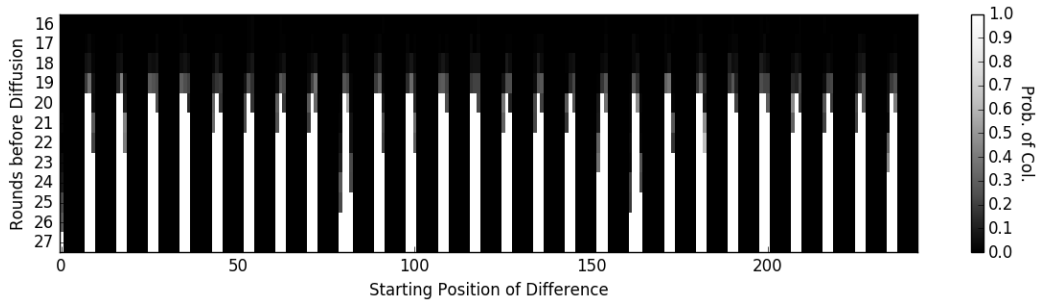
In Figure 2, we visualize our experimental results for how differences diffuse over rounds in Curl-P-27. We shade the graph by the probability that a collision occurs given a particular starting position of a 1-trit difference (x-coordinate) and a lower bound on the number of rounds that a 1-trit difference does not diffuse (y-coordinate). To experimentally generate this dataset we performed 100 samples per position and round depth ( $11 \cdot 243 \cdot 100 = 267300$  samples in total). Each sample was initialized to a random state with random difference injected at the anticipated position and round.

In related work, Colavita and Tanzer [CT18, Lemma 3] make the following observation: two message blocks with 1-trit difference in position  $p$ , when evolved over  $\ell$  rounds, can only have differences in a contiguous modular region of size at most  $2^\ell$ . Moreover, the starting point  $p'_\ell$  of this region depends only on  $p$  and  $\ell$ . Later, their attacks use exhaustive

---

<sup>5</sup>The IOTA developers have asked us to note their statement that “it was widely communicated that IOTA was utilizing a prototype hash function since inception”.





**Figure 2:** Our experimental estimate of the prob. of a full state collision for any two message blocks differing by one trit: x-axis is the position of the trit that differs and y-axis is number of rounds before the two internal states differ by more than one trit.

search techniques to find a 1-trit difference does not diffuse.

Our collision attack uses a 1-trit difference at position 17 in the differing message blocks. Using the results of this experiment, we can calculate the probability of a collision if a 1-trit difference, starting at position 17, is maintained for certain number of rounds. By maintained, we mean that for each of these rounds the difference between both states is always 1-trit. For position 17, the probability of a collision is 1.0 for 20 rounds. Thus, if we prevent diffusion of a 1-trit difference starting at position 17 for at least 20 rounds we should find a collision. This attack should work for some of the other positions in the input message block (as shown in [CT18]). Note that since this estimate is slightly pessimistic since it does not count all the ways a collision could result from a 1-trit difference. Instead it only estimates the probability of collision where 1-trit difference is maintained for at least 20 rounds *i.e.*, does not diffuse.

#### 4.1 Differential Properties of Curl-P Transformation Function $f^r$

In this section, we show how to find states that maintain a 1-trit difference for at least 20 rounds. This involves analyzing the differential properties of Curl-P's transformation function  $f$ .

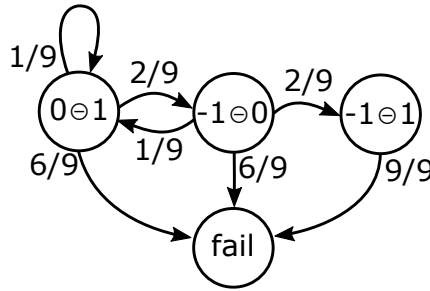
Differential cryptanalysis is concerned with studying the propagation patterns of differences between two or more sets of inputs. The most common technique is the discovery of *differential trails*. A differential trail is a probabilistic bias of how a set of differences will propagate to another set of differences through many rounds of a cryptographic function. Here we only work with a specific differential trail, namely a 1-trit difference between the two states  $S, S'$  under repeated applications of the transformation function  $f$ . We show that Curl-P has a strong bias toward maintaining a 1-trit difference across rounds (*i.e.*, applications of  $f$ ).

Let's first introduce some necessary terminology. Since Curl-P operates on trits in  $\{-1, 0, 1\}$  instead of bits in  $\{0, 1\}$ , we must use new notation for ternary differentials. To represent the difference between two trits,  $x$  and  $x'$  we use  $\ominus$  (*e.g.*,  $0 \ominus -1$ , which means either  $x = 0$  and  $x' = -1$  or  $x = -1$  and  $x' = 0$ ). By the term diffusion, we indicate that after an application of  $f$  the number of differences between the two states has increased (*i.e.*, the differences have diffused).

Our attack is built around the fact that the s-box  $g$  does not always propagate differences. For example, consider two sets of inputs and outputs to  $g$ :  $a, b, c$  and  $a', b', c'$  such that  $g(a, b) = c$  and  $g(a', b') = c'$ . We make the following observations:

1. For all possible values, if  $a \neq a'$  and  $b = b'$  then it will always be the case that  $c \neq c'$ .





**Figure 3:** The Markov chain represents the probability of starting from a single difference of a particular type and ending in a single difference.

2. If  $a = a'$  and  $b \neq b'$  then both  $c = c'$  and  $c \neq c'$  are possible (e.g.,  $a = a' = 1, b = 0$  and  $b' = -1$  then  $c = c' = 0$ ).

Each round of  $f^r$  i.e., each application of  $f$  in  $f^r$ , is called to update the state. As discussed in Section 3.3, each trit in the updated state depends on the output of applying  $g$  to two trits in the prior state. Each trit in the state is plugged into the s-box  $g$  twice, once as the first parameter  $a$  and once as the second parameter  $b$ . This means that a 1-trit difference will always propagate to the next round, since when it is the first parameter  $a$  to the s-box  $g$ , the output of  $g$  will differ based on a difference in  $a$  (as shown in observation 1). Thus if you apply  $f$  to two states  $S, S'$  which have 1-trit difference the updated states  $f(S), f(S')$  will either differ by 1 trit or 2 trits. It will never result in a 0-trit difference.

We model the probability that a 1-trit difference will remain a 1-trit difference across  $k$  rounds of Curl-P. As shown in Figure 3, by enumerating all possible inputs to  $g$  we develop a Markov model of the possible difference states after an application of  $f$  starting from a 1-trit difference. For instance, if the 1-trit difference in the current round is  $0 \ominus 1$ , then with probability  $1/9$  the difference stays the same (i.e.,  $0 \ominus 1$ ) in the next round, with probability  $2/9$  the difference becomes  $0 \ominus -1$  in the next round, or with  $6/9$  the number of differences increases from 1 to 2 (marked as the fail state as it fails to maintain a 1-trit difference).

As shown below, we convert the Markov model to a state transition matrix.

$$\begin{bmatrix} 1/9 & 2/9 & 0 & 6/9 \\ 1/9 & 0 & 2/9 & 6/9 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix}^k$$

The top row represents the state transitions probabilities of  $0 \ominus 1$ , the second row  $-1 \ominus 0$ , third row  $-1 \ominus 1$ , and fourth row that a 1-trit difference diffuses to a 2-trit difference (the fail state). Using this matrix, we compute a lower bound on the probability that after  $k$  applications of  $f$ , the number of differences remains at 1. This is a lower bound as our analysis does not count transitions which increase to a difference of 2 trits or more and then later become a 1-trit difference.

Thus, starting from a 1-trit  $0 \ominus 1$  difference we calculate a lower bound on the probability of a 1-trit difference by raising the matrix to the number of rounds we wish to investigate. For example, if we raise it to the power 3, the transition probabilities in the matrix represent the probability that you arrive at that difference after three rounds. Thus, we can measure the probability we don't fail after  $k$  rounds.

Earlier, we experimentally verified that if a 1-trit difference starting at position 17 is maintained for 20 rounds of Curl-P-27 (i.e., 20 applications of  $f$ ), then the probability of a collision is 1.0. Using our state transition matrix we calculate for 20 rounds, our attack

has a per query success probability lower bounded by  $2^{-42}$ . That is, we need to make on the order of  $2^{42}$  queries to Curl-P-27 with different message pairs before finding a pair that will maintain a 1-trit difference for 20 rounds. This message pair will result in a collision. In the next section, we will show we can significantly reduce the necessary number of queries to Curl-P-27.

## 4.2 Solving for a 1-Trit Difference

In this section, we show how to reduce the number of queries to Curl-P-27 by selecting messages with particular properties. We first show how to constrain two states  $S$  and  $S'$ , which differ by 1-trit, such that for at least 9 applications of  $f$  a 1-trit difference will be maintained (*i.e.*, there is no diffusion of differences). To do this, we represent  $f$  as a system of equations and solve for particular values of trits in states  $S$  and  $S'$ .

We can represent the transformation function  $f^t(S)$  as a series of equations. For example, a single call to  $f$  can be written as

$$f(S)_0 = g(S_0, S_{364}), f(S)_1 = g(S_{364}, S_{728}), \dots, f(S)_{728} = g(S_{365}, S_0)$$

where  $f(S)_0$  is the trit in position 0 of the updated state after  $f$  is applied. Since each round is just the recursive application of  $f$ , we can write the value of a particular trit after a number of rounds of  $f$  in terms of some of the initial values of the state  $S$ . We use superscript to denote the number of rounds of  $f$ . For example, with

$$f^2(S)_6 = g(g(S_{366}, S_1), g(S_{184}, S_{548}))$$

we specify the equation for the trit in position 6 after two rounds of  $f$ .

Using this representation, we find the equations guaranteeing that a 1-trit  $0\oplus 1$  difference is maintained for 9 rounds. We then find a message prefix that satisfies these equations. To do this we wrote a simple program which modifies trits until we find a set of values that satisfies the equations. This program takes less than a second to find a solution (see Section 5.3 for our performance evaluation). The number of values in the two colliding messages that must be fixed to satisfy the equation scales with number of rounds *i.e.*, more rounds means more trits in the message whose values can't be changed. To ensure the later stages of the attack have plenty of flexibility to change the messages we limit our approach to 9 rounds. Thus, given a particular message template, we only have to change a small set of trits in two message blocks to transform it into a satisfactory message.

## 4.3 Finding Collisions

We now combine our two methods to generate collisions for Curl-P-27. We refer to the technique shown in Section 4.2 of choosing a message prefix such that a 1-trit difference at a particular position will not diffuse across 9 rounds as the *constraint phase* of our attack. We call our method in Section 4.1 of trying different messages to increase the number of rounds for which a 1-trit difference is maintained the *brute-force phase*.

Our collision attack requires at least three message blocks:  $mb_a$ ,  $mb_b$ , and  $mb_c$ , where  $mb_b$  contains the difference. Any number of message blocks can exist before  $mb_a$ . Any number of message blocks can exist between  $mb_a$  and  $mb_b$ .  $mb_c$  is always the next message after  $mb_b$  and overwrites the differences in the first third of the state created by  $mb_b$ . The actual value of  $mb_c$  has no impact on the attack and can be anything.

Our full attack works as follows. First, in the constraint phase of our attack we find a suitable message prefix by altering trits in parts of  $mb_a$  and  $mb_b$ , such that they guarantee a 1-trit difference across 9 rounds of  $f$ . Next, in the brute-force phase we randomly alter trits in particular positions in  $mb_b$  with the objective of finding two messages such that a 1-trit difference starting in position 17 is maintained across 20 rounds. As the constraint phase

ensures that each attempt in the brute-force phase also maintains a 1-trit difference across 9 rounds, the attack complexity of the brute-force phase is reduced from 20 rounds to 11 rounds. As a result, the estimate of the success rate per query is reduced to approximately  $2^{-22.87}$  or one out of 7.6 million.

As is typical in differential cryptanalysis, our probability calculations make the simplifying assumption that actual values of the differing inputs are uniformly random. Due to the low diffusion rate of Curl-P-27 and the non-random properties of Curl-P this assumption may not always hold. However, as seen in Section 5.3, the estimates given in this section are reasonably close to the actual results.

## 5 Exploiting Collisions in Curl-P to Forge Signatures

In this section, we show how our collision attack against Curl-P-27 can be used to perform a signature forgery attack against the IOTA Signature Scheme (ISS). Continuing from the previous section, we show how to create two valid IOTA bundles (*i.e.*, payments), which differ in at most two trits and have the same Curl-P-27 hash. Then, we will describe the setting of our attack which exploits these colliding bundles to forge signatures. Finally, we will show how to perform this attack against multi-signatures (multisigs) as used by ISS.

### 5.1 Chosen-Message Attack on ISS

Our attack is a chosen-message attack, which means that a malicious user Eve tricks a user Alice by asking Alice to sign a bundle,  $b_1$ , and then later producing a different bundle,  $b_2$ , which also verifies under the signature Alice provided. In more detail:

1. Alice generates a key pair (PK, SK).
2. Eve uses our collision attack on Curl-P-27 to produce two bundles  $b_1, b_2$  such that  $b_1 \neq b_2$  and  $\text{CurlHash}(b_1) = \text{CurlHash}(b_2)$ .
3. Eve sends  $b_1$  to Alice and asks Alice to sign it. Alice inspects  $b_1$  and confirms that it is benign.
4. Alice sends Eve a signature  $\sigma$  on  $b_1$ , *i.e.*,  $\text{Sign}(\text{SK}, b_1) \rightarrow \sigma$ .
5. Eve produces a signature, bundle pair  $(\sigma, b_2)$  such that  $b_1 \neq b_2$ ,  $b_2$  is a valid bundle, and  $b_2$  verifies under Alice's PK even though Alice has never seen  $b_2$ .

In Section 4.3, we introduced the general format of our attack, which requires at least three message blocks  $\text{mb}_a$ ,  $\text{mb}_b$ , and  $\text{mb}_c$ . To perform the first phase of the attack, we set certain trits in  $\text{mb}_a$  and  $\text{mb}_b$  to particular values. In the brute-force phase, we change other trits in  $\text{mb}_b$  each attempt and check to see if we have achieved a collision. However, the bundles must pass the validity checks in the IOTA software in order for them to be accepted in IOTA as valid bundles, which limits the trits we can modify to perform our attack.

A bundle's hash is computed by hashing the concatenation of the address, value, tag, timestamp, current index, and last index fields of each transaction in the bundle. Each transaction supplies two message blocks. Recall from Section 3.1 that "transactions" in IOTA are more like inputs and outputs; a valid payment requires multiple transactions. The format of a transaction is shown in Figure 4. Most of these fields are constrained in well-formatted bundles—for example, the values in a bundle cannot sum to a negative number, the timestamp must be within a certain range, and the indexes must line up with the transactions in the bundle. Tags do not impact the semantics or validity of the bundle

Signature Fragment (6561)				
Address (243)				
Value (81)	Tag (81)	TS (27)	Current (27)	Last (27)
Bundle Hash (243)				
Trunk Transaction Hash (243)				
Branch Transaction Hash (243)				
Nonce (243)				

**Figure 4:** IOTA transaction format. Field sizes are in trits. The shaded fields are used to calculate the bundle hash of the bundle in which the transaction is included. A bundle must have multiple transactions in order to be a valid payment.

and can contain arbitrary trits. Thus, for both the constraint phase and for each attempt in the brute-force phase of our attack, we only change the trits in the tags.

Another important question is where Eve can place the collision to cause damage. In our initial vulnerability report, we demonstrated colliding bundles for two different styles of attack: one which places the collision in the address field so Alice unwittingly signs a bundle which burns funds that were originally intended for Eve, allowing Eve to claim Alice made a mistake, and a second which places two collisions in two different value fields in a bundle so that Alice unwittingly signs a bundle which pays Eve more than intended. In the following section, we describe in detail the latter attack style for bundles which require multiple signatures, which fits our chosen message setting.

## 5.2 Multi-signature Attack

One criticism of the signature forgery attacks presented in our vulnerability report [HNDV17] is that they are chosen-message attacks, that is, Eve must ask Alice to sign a bundle. To help demonstrate the importance of chosen-message security, we now extend our attacks to the IOTA multi-signature (multisig) scheme [Sch18]. In multisig, funds can be spent only by signatures from multiple parties. To spend, one party creates a bundle and asks the other party to sign it, which is exactly a chosen-message attack. The IOTA Foundation encourages exchanges deploying a hot storage/cold storage solution<sup>6</sup> to use multisig for securely storing funds [Fou18a]. One of the main reasons multisig is used in a cryptocurrency context is that it requires that an attacker must compromise more than one party to steal funds. Our attack removes this security benefit of multisig. We will consider a simple case of a 2-of-2 multisig where two parties both sign to spend funds; however, our attack generalizes to more complex settings.

Consider two parties—Eve and Alice—each holding a pair of ISS keys— $(PK_E, SK_E)$  and  $(PK_A, SK_A)$ —and funds which can only be spent by both a signature from Eve’s secret key and a signature from Alice’s secret key. This implies that Eve and Alice previously entered into a 2-of-2 multisig and are now spending those funds jointly. Our attack will work as follows: Eve will compute two colliding bundles, one which pays funds to Alice and one of which pays funds to Eve. Eve will sign and send to Alice the bundle that pays Alice. Once she has Alice’s signature, Eve will use it on the colliding bundle to create a valid bundle which Alice never saw or authorized, and will broadcast this bundle.<sup>7</sup> In this setting, Eve is either malicious or has been compromised by a malicious party.

<sup>6</sup>This is cryptocurrency terminology for structuring the control of an account holding funds such that any transfer of funds requires the consent of two secret keys. One of the secret keys is “cold,” meaning it is kept in a location not connected to the internet such as an airgapped device.

<sup>7</sup>Note that if Alice is a cold wallet, she relies on Eve to broadcast the transaction.

Txn	MB	Message block contents				
0	0	Address				
	1	Value	Tag	TS	Current	Last
1	2	Address				
	3	Value	Tag	TS	Current	Last
2	4	Address				
	5	Value	Tag	TS	Current	Last
3	6	Address				
	7	Value	Tag	TS	Current	Last
4	8	Address				
	9	Value	Tag	TS	Current	Last
...						

**Figure 5:** A bundle portioned into its message blocks. We target two collisions: One in the 17th position of message block 3 by manipulating the trits in the tag portions of message blocks 1 and 3, and one in the 17th position of message block 7 by manipulating the tags in 5 and 7. The red arrows indicate the collisions.

In order to construct such bundles, Eve places the collisions in certain value fields in certain transactions. Figure 5 shows the first four transactions in such a bundle divided into message blocks. The highlighted fields are the trits relevant to our attack. Eve causes a collision in trit 17 of the value field in the second transaction (message block 3) by manipulating the trits in the tag fields both before and after the collision. By doing this, Eve can produce two bundles with different values in the second transaction that have the same bundle hash. Eve creates a second collision later on in the bundle in the fourth transaction (message block 7), this time arranging the collision so that the values still sum to zero in both colliding bundles. This serves to change what amounts are paid to whom in the transaction.

Generating these collisions essentially requires running the attack twice, sequentially. In our current collision tool, we require one transaction between the two transactions where we collide in the value fields. Other than this requirement, and the requirement that the collisions are not in the first or last transactions, we can handle bundles with different numbers of transactions. That our tool can only cause collisions in the 17th trit of a message block is a limitation of the tool’s current implementation, not of the cryptanalysis techniques described in Section 4. Our tool does not depend on the specific addresses and values in the transactions to generate collisions, but the collision trits in the values that are changed must be different in order to produce valid bundles. For example, if trit 17 is zero in both Alice and Eve’s output values in  $b_1$ , then flipping trit 17 in Eve’s output to one in  $b_2$  will cause the values in  $b_2$  to not sum to zero. In  $b_1$  Alice’s output value’s trit 17 should be one and Eve’s should be zero.

In Appendix B we show the contents of two example bundles we created using this technique. In this example, the bundles are spending a multisig input of 500,000,000 IOTA controlled by Alice and Eve. Alice signs a bundle which pays Eve 1 IOTA and the remainder to other addresses. In the colliding bundle, Eve receives 129,140,164 IOTA, at the expense of Alice’s address.

Generating colliding single-signature bundles operates in much the same way; our vulnerability report demonstrated a signature forgery on a bundle which paid out to three addresses. In the benign bundle  $b_1$  Alice receives 50,000 and 810,021,667 IOTA to two addresses she controls and pays 100 IOTA to Eve. In the malicious bundle  $b_2$  Eve changes this so that she receives 129,140,263 instead of 100, at the expense of Alice’s funds. We

**Table 2:** Run time of the constraint phase, brute-force phase, and the entire multisig attack which involves running each phase twice. Measurements over 5000 iterations.

Part	Average	Min	Max
Constraint phase	1.1 s	0.27 s	6.1 s
Brute-force phase	7.2 s	0.04 s	74 s
Multisig	15.2 s	1.4 s	74 s

have not investigated the effects of placing collisions in fields besides the value and address. Other attacks might be possible.

### 5.3 Performance Analysis

We ran these attacks on a 80-core Intel machine with 8 2.4GHz 10-core Intel chips and 256 GB of RAM, running 64-bit Linux 4.9.74. Our attacks use all of the CPU but a negligible amount of the RAM. As described in Section 4.3, finding a collision consists of two phases: the constraint phase calculates the set of constraints, and the brute-force phase generates randomness in the tags to find collisions.

The constraint phase generates and solves eighteen equations, two for each of the first nine rounds of Curl-P-27. The constraint phase is implemented in Python, and runs on a single core. We did not try to optimize the first phase. Table 2 shows the average, minimum, and maximum times of running the first phase 5000 times, when colliding on the 17th trit.

Table 2 also shows measurements for the brute-force phase, which uses the trits and template generated from the first phase to brute-force one collision. This is implemented in Go and parallelizes well, so we use all 80 cores of our server. On average it only takes 7.2 seconds to find one collision using the output of the first phase. On average, it takes 5.2M attempts to find one collision, with the minimum and maximum attempts over 5000 runs 1279 and 53M, respectively. This corroborates our analysis in Section 4.3.

In order to perform the multisig attack described in Section 5.2, we must run both the constraint phase and the brute-force phase twice, sequentially, to find two collisions. Using our collision tool, it takes on average 15.2 seconds to produce two multisig bundles which differ in two places. Table 2 shows the average, minimum, and maximum times for 5000 runs with the same starting bundle.

## 6 Security Against Generic Attacks

In this section, we show upper bounds on the security offered by ISS against generic brute force attacks. These attacks are independent of the underlying hash function (currently instantiated as Kerl, a wrapper around Keccak-384), which we will model here as a random oracle, *i.e.*, an ideal cryptographic hash function. The attacks described in this section are unrelated to the attacks on Curl-P and apply to the live IOTA network. However they do not present a critical vulnerability for default security parameters of ISS.<sup>8</sup> Generic attacks consider the strength of a function to brute-force attacks *e.g.*, trying values blindly until guessing the preimage. The generic attacks presented here focus on a part of the IOTA Signature Scheme (ISS) called *normalization*. Put simply, when signing a message IOTA composes the output of Kerl or Curl-P with a normalization function. The output space of the normalization function is smaller than the input space. Therefore even if two

<sup>8</sup>We disclosed these results to the IOTA developers on Aug 6 2018 and made these results public March 31 2019.

messages do not collide in the hash function, collisions or preimages could be found after the normalization function is applied.

To review informally, a cryptographic hash function is an efficiently computable function  $\mathcal{H}$  that maps an input of arbitrary length to an output digest (or hash) with constant length  $\lambda$ , and that satisfies the following security properties [KL14]:

**Collision Resistance:** It is hard to find  $x_1, x_2$  such that  $x_1 \neq x_2$  and  $\mathcal{H}(x_1) = \mathcal{H}(x_2)$ .

**Second Preimage Resistance:** Given  $x_1$ , it is hard to find  $x_2 \neq x_1$  such that  $\mathcal{H}(x_1) = \mathcal{H}(x_2)$ .

**Preimage Resistance:** Given  $y$ , it is hard to find  $x$  such that  $\mathcal{H}(x) = y$ .

In concrete implementations, “hard” means that breaking these guarantees should take a number of queries to the hash function exponential in  $\Theta(\lambda)$ . Preimage resistances are typically expected to have  $\lambda$  bits of security, while collision resistance has  $\lambda/2$  bits of security (due to the generic birthday attack) [Wie05].

However, ISS does not act directly on  $\mathcal{H}(\text{msg})$ , but rather on a modified hash value  $\text{Norm}(\mathcal{H}(\text{msg}))$ , so we must analyze the above security parameters for the composed hash function, which has a reduced output space and thus decreased difficulty of brute force attacks. Rather than sign the message hash directly, ISS modifies the hash of the message in two ways: First, it sets the last trit in the digest to 0; this is done to deal with mismatched buffer sizes when converting from binary to ternary. Second, it applies the normalization procedure described in Algorithm 3. This is done to prevent an attack that is typically mitigated in W-OTS by the addition of a checksum. As part of normalization, if any of the digest trytes have the value 13, normalization will fail and the signature generation retries the hashing and signature process with a modified message; this is to mitigate the attack where signing the value 13 will leak part of the signer’s private key [Pin18].

The normalization process splits the message digest into three chunks of 27 trytes each. Then in each chunk, starting from the left, it increments or decrements each tryte within the range  $[-13, 13]$  until the sum of all the trytes in the digest is 0. For example, if the last 26 trytes in the chunk sum to 12 and the first tryte is 5, the normalization procedure will decrement the first tryte to  $-12$ . Once the normalization process computes the normalized hash (and the hash does not fail the 13 check, mentioned above), either one, two, or three of these chunks will be output and signed, corresponding to level 1, 2, or 3 security in ISS. We describe the normalization process as pseudocode in Alg. 3.

This means that to forge a level 1 signature on a message  $\text{msg2}$  given a valid signature on  $\text{msg1}$ , only the first thirds of  $\text{Norm}(\mathcal{H}(\text{msg1}))$  and  $\text{Norm}(\mathcal{H}(\text{msg2}))$  must collide. IOTA claims that these levels provide 81, 162, and 243 trits of security (128.3, 256.7, and 385.1 bits) respectively [Fou18b], as they are equivalent to using hash functions with smaller output size and thus smaller  $\lambda$ . Security level 2 is the default. Security level 1, while supported and once recommended for small value amounts [Søn16], is no longer recommended for use [Fou18d].

In order to compute the impact normalization has on preimage and collision resistance, we do the following: First, we determine the number of outputs that result from normalization. Second, we compute the number of inputs that do not fail normalization. Most valid inputs which are passed to normalization will fail causing the signature method to fail and requiring a change to the message to be signed. Finally, since the output of the normalization function is not uniform, to determine the preimage and collision resistance we must compute the Shannon entropy. Shannon entropy takes into account the differences in probabilities between outputs of the normalization function. Intuitively, some outputs are more likely because the normalization algorithm adjusts trytes starting from the left until the sum is zero. This means that normalized chunks that have large tryte values



---

**Algorithm 3:** ISS's deterministic hash normalization procedure.

```

Function Norm(hmsg, level):
  hmsgTrytes ← ToTrytes(hmsg);
  for  $i \leftarrow 0; i < \text{level}; i \leftarrow i + 1$  do
    chunk ← hmsgTrytes[ $i * 27 : i * 27 + 27$ ];
    for  $j \leftarrow 0; \text{sum}(\text{chunk}) \neq 0$  do
      if  $\text{sum}(\text{chunk}) < 0$  then
        if  $\text{chunk}[j] == 13$  then
           $j \leftarrow j + 1;$ 
        else
           $\text{chunk}[j] \leftarrow \text{chunk}[j] + 1;$ 
      if  $\text{sum}(\text{chunk}) > 0$  then
        if  $\text{chunk}[j] == -13$  then
           $j \leftarrow j + 1;$ 
        else
           $\text{chunk}[j] \leftarrow \text{chunk}[j] - 1;$ 
    for  $n \leftarrow 0; n < 27; n \leftarrow n + 1$  do
      if  $\text{chunk}[n] = 13$  then
        return Fail;
    nmsg[ $i * 27 : i * 27 + 27$ ] ← chunk;
  return ToTrits(nmsg)

```

---

later in the string will have more preimages than those with small tryte values and thus a random input is more likely to be normalized to one of these outputs.

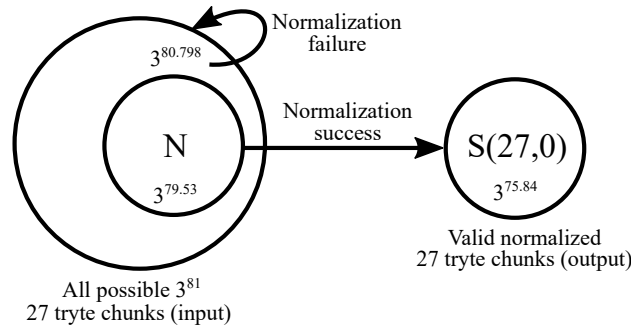
To compute the size of all possible valid normalization outputs, we introduce the following recurrence which we will use throughout this section. Let  $S(n, u)$  denote the number of  $n$ -tryte chunks not containing a 13 with a sum of  $u$ . This produces the following recurrence:

$$S(n, u) = \begin{cases} 1 & \text{if } n = 0 \wedge u = 0 \\ 0 & \text{if } n = 0 \wedge u \neq 0 \\ \sum_{k=u-12}^{u+13} S(n-1, k) & \text{otherwise .} \end{cases}$$

The two base cases,  $(n = 0, u = 0)$  and  $(n = 0, u \neq 0)$ , correspond to tryte strings of length 0 which should only sum to the value 0, so there are no valid tryte strings that sum to a value other than 0. For the recurrent case (when  $n > 0$ ), we consider every possible value of the first tryte: all integers in the range  $[-13, 12]$ . We do not count 13 because the number 13 is not allowed in normalized outputs. When the first tryte has value  $v$ , in order for the total of the string to be  $u$ , the remaining  $n - 1$  trytes must sum to  $u - v$ . Therefore there are  $S(n - 1, u - v)$  such tryte chunks for each value of  $v$ . Letting  $k = u - v$  and summing over  $v \in [-13, 12]$ , we obtain the equation above.

Using this recurrence, we can count how many 27-tryte chunks ( $n = 27$ ) sum to 0 ( $u = 0$ ) and do not contain a 13. The number of all possible normalization outputs for one chunk is  $S(27, 0) \approx 3^{75.84}$ , meaning that the normalized distribution has a max entropy of approximately 75.84 trits (120.21 bits). However, as we stated earlier, the output of normalization is not uniformly distributed. Therefore the actual entropy (and security) is less than this upper bound.

We will now compute the value  $N$ , the number of inputs (unnormalized chunks) which normalize successfully; that is, after normalization none of the trytes are 13. Figure 6 shows how the input space of all possible 27-tryte chunks maps to the output space of



**Figure 6:** Mapping of input space of unnormalized chunks to the output space of valid, normalized chunks.

$S(27, 0)$  valid normalized chunks.  $N$  is the total number of unnormalized chunks in the input space that will end up completing the normalization procedure without failure. Let  $\ell = 27$  be the length of the tryte string to be normalized. The following is the expression for  $N$ :

$$N = \sum_{k=-12}^{12} 27 \cdot S(\ell - 1, k) + \sum_{p=1}^{\ell-1} \sum_{b=-12}^{12} 27^p \cdot (14 - b) \cdot S(\ell - 1 - p, 13p - b)$$

To understand the derivation for  $N$ , we note that the normalization will attempt to make the tryte string sum to 0 by incrementing or decrementing the first tryte until it reaches 13 or  $-13$ , proceeding onto the second tryte if it hits one of these limits while the sum is still not 0. Note that if a tryte reaches 13, normalization will fail and abort. Therefore, if normalization modifies the second tryte, the first tryte must have been modified to  $-13$ . Similarly, if normalization modifies the  $n$ th tryte for  $n > 1$ , the previous  $n - 1$  trytes must have been modified to  $-13$ .

The value  $p$  represents the number of trytes that are modified to  $-13$ . The first term in  $N$  is almost all the tryte strings which only require modifying the first tryte (or none) in normalization. Here,  $p = 0$ , meaning that at the end of normalization either no trytes were changed, or only the first tryte was changed and it did not reach a value of  $-13$ . Further, suppose the first tryte takes on a final (post-normalization) value  $j$  in the range  $[-12, 12]$ . As only the first tryte is potentially modified, our output string must contain the final  $\ell - 1$  trytes unmodified. If these final  $\ell - 1$  trytes sum to  $-j$ , then the first tryte will take on a final (post-normalization) value of  $j$ , regardless of its initial value. There are  $S(\ell - 1, -j)$  such values for the final  $\ell - 1$  trytes. As there are 27 potential values of the first tryte for each of these cases, we obtain the first summation of  $N$ .

Next, suppose that  $p > 0$ , meaning that  $p$  trytes have been normalized to  $-13$  at the end of the procedure. Consider the final value  $b$  of the tryte following the  $p$  trytes with value  $-13$  (see Figure 7). This tryte can take on any value  $b$  in the range  $[-12, 12]$ . Note that the first  $p$  trytes will be reduced to  $-13$  and the next tryte will be reduced to  $b$  if and only if the final  $\ell - 1 - p$  trytes sum to  $13p - b$ . There are  $S(\ell - 1 - p, 13p - b)$  such tryte strings. Furthermore, note that the initial value of the tryte following the first  $p$  trytes must be greater than or equal to than its final value  $b$ , as it is never increased by the normalization procedure. This is because to obtain  $-13$  values in the first  $p > 0$  trytes, the original sum must be greater than or equal to 0, and thus normalization will never increase tryte values. There are exactly  $14 - b$  tryte values greater than or equal to  $b$ . Finally, again we recognize that the initial values for the first  $p$  trytes do not matter, as they are unconditionally normalized to  $-13$  by the procedure. This yields the  $27^p$  term. Combining and iterating over  $p$ , which can take on any value between 1 and  $\ell - 1$ , and

**Table 3:** 1st and 2nd preimage resistance, claimed and actual in trits (bits).

Level	Claimed Security	1st & 2nd Preimage Resistance Upper Bound
1	81 trits (128.38 bits)	73.12 trits (115.90 bits)
2	162 trits (256.76 bits)	146.25 trits (231.80 bits)
3	243 trits (385.15 bits)	218.37 trits (346.11 bits)

$b$  which varies in the range  $[-12, 12]$ , we obtain the latter term of  $N$ . This value of  $N$  implies that 19.9% of unnormalized chunks will complete normalization without failure.

$$\vec{t} = \left( \underbrace{(-13, -13, \dots, -13)}_{\text{first } p \text{ trytes}}, b, \underbrace{t_{p+2}, t_{p+3}, \dots, t_\ell}_{\substack{\text{remaining } \ell - 1 - p \text{ trytes} \\ \text{must sum to } 13p - b}} \right)$$

**Figure 7:** Adjusted hash value when normalization procedure stops after normalizing first  $p$  trytes. Letting  $b$  be the value of  $p + 1$ -th tryte, the remaining  $\ell - 1 - p$  trytes must sum to  $13p - b$ , as the first  $p$  trytes were normalized to  $-13$ . We conclude that this case happens for  $S(\ell - 1 - p, 13p - b)$  strings.

In order to analyze the entropy of the resulting normalized distribution, we must compute the probability that a random input produces each normalized output. We begin by restricting our analysis to those inputs which correctly normalize (of which there are  $N$ ). The first summation captures each input string which results in a normalized string beginning with no trytes of value  $-13$ . Each of these output strings is produced by input strings matching in all  $\ell - 1$  final trytes with any value for the first tryte. Thus the probability that each one is produced is  $\frac{27}{N}$ . Per the analysis above, there are  $S(\ell - 1, k)$  such strings for each value  $k$  the initial tryte can take on in the normalized output.

In the latter summation, note that each normalized string beginning with  $p$  trytes of value  $-13$ , followed by a tryte of value  $b$ , is produced by exactly  $27^p(14 - b)$  input strings. For each value of  $p$  and  $b$ , there are  $S(\ell - 1 - p, 13 - p - b)$  such strings. Using this analysis, we can construct a vector  $\vec{p}$  of length  $S(27, 0)$  in which each element corresponds to the probability of obtaining a given output string. This probability vector can then be used to compute relevant entropies and collision probabilities.

Using this probability vector, we compute the Shannon entropy of the distribution to obtain 73.12 trits (115.90 bits) of security, lower than the 81 trits (128.38 bits) claimed by IOTA. Furthermore, the normalized distribution has a considerably lower min-entropy of 40.56 trits (64.29 bits). As each chunk in the hash is independent under the random oracle model, this entropy scales linearly with the security level *i.e.*, level 2 has twice the entropy of level 1.

When generalizing to the third security level, we note that the last trit of the digest is zeroed before normalization. This may affect the entropy to a very minor degree. Our estimates do not account for this, but the effect can be simulated by modifying the base cases of  $S$ . Without this complication, we obtain the following breakdown of the three security levels against preimage attacks:

These results translate directly to upper bounds for ISS's EU-RMA security. If an attacker Eve sees a signature on a message `msg1` chosen by Alice, it is sufficient to find a second preimage `msg2` such that  $\text{Norm}(\mathcal{H}(\text{msg1})) = \text{Norm}(\mathcal{H}(\text{msg2}))$  in the relevant chunks for the security level. This brute forced message can have meaningful structure, as long as there are sufficiently large unspecified regions, like the Tag region of IOTA transactions described in Section 5, in which to probe the search space.

**Table 4:** Collision resistance (CR) upper bound, claimed and actual upper bound in trits (bits). Third column, marked valid, shows impact of rejecting collisions that fail the normalization output check.

Level	Claimed Security	CR Upper Bound	CR Upper Bound (valid)
1	40.5 trits (64.19 bits)	29.74 trits (47.14 bits)	34.15 trits (54.14 bits)
2	81 trits (128.38 bits)	58.85 trits (93.27 bits)	63.26 trits (100.27 bits)
3	121.5 trits (192.57 bits)	87.42 trits (138.56 bits)	91.83 trits (145.56 bits)

We also calculate an upper bound for collision resistance of the composed hash function and apply this result to ISS’s EU-CMA security. Consider a generic brute force attack, in which we aim to find `msg1` and `msg2` such that  $\text{Norm}(\mathcal{H}(\text{msg1})) = \text{Norm}(\mathcal{H}(\text{msg2}))$  for the chunks used in our given security level. Furthermore, suppose `msg1` and `msg2` are drawn from a sufficiently large uniform distribution. Let  $X$  denote the previously described distribution of the normalization function’s image, and  $M$  denote the number of samples drawn from  $X$  before finding a collision. We can compute the quantity  $\beta^{-1} = \|X\|_2^2$ , where  $\beta^{-1}$  represents the probability that two independent samples from  $X$  are equal, as  $\sum_i p_i^2$  using the probability vector defined above. Therefore, we can compute an upper bound on level one collision resistance using  $E[M] \leq 2\sqrt{\beta}$  [Wie05]:

$$\beta^{-1} \approx 1.67 \cdot 10^{-28}$$

$$E[M] \leq 2\sqrt{\beta} \approx 1.55 \cdot 10^{14} \approx 2^{47.14}$$

Thus, assuming we draw uniformly from the preimage of valid outputs of the normalization function, we can stage a brute force attack on level 1 keys using only  $2^{47.14}$  queries. Repeating this analysis for levels 2 and 3 yields the complexities in Table 4.

However, note that only approximately 0.788% of inputs are in the preimage of valid normalization outputs. Thus, our attack time is increased by a factor of approximately  $2^7$ . Upper bounds for collision resistance translate directly to upper bounds for ISS’s EU-CMA security. Eve can generate two messages that collide under the composed hash function, and Alice’s signature for one will be valid for the other. Using standard techniques, these brute-forced collisions can have distinct and prespecified structure at the cost of 2x overhead.

While these brute-force attacks largely do not have practical query complexities—with the notable exception of the level 1 EU-CMA attack—these bounds are not only lower than stated figures, but they also represent just one attack vector against an idealized ISS.

## 7 Discussion

The IOTA developers have made several statements discussing the impact and the cause of these vulnerabilities. We summarize these statements and address some of the concerns.

The IOTA developers have argued that the chosen-message attack model is irrelevant in the context of the complete IOTA network: specifically, that the chosen-message setting is implausible because “in IOTA an attacker doesn’t choose the signed message” [blo18]. In response to the critique on the chosen-message setting we extended our signature-forgery attack to work on payments spending from a multisig address since the multisig protocol explicitly allows one user to choose the message that another user will sign.

The IOTA developers also argued that “even most valid attacks” would fail on the live IOTA network because of unspecified “protection mechanisms” in the closed-source coordinator [blo18, Fou18c]. The attacks presented in the vulnerability report and this paper are against the IOTA Signature Scheme in isolation. We did not analyze these attacks within the context of the complete IOTA system.

Additionally, IOTA developers claimed that the ability to find colliding inputs to Curl-P-27 was intentional and was for the purposes of preventing “scam clones.” It is worth quoting them in full: “The IOTA team made a design decision early on to prevent this possibility [of scam clones] by purposefully introducing the Curl-P hashing function with known practical collisions. This had the express purpose of rendering fraudulent clones of the protocol useless in their application as a DLT protocol, while at the same time guaranteeing the security of the IOTA protocol and network as a whole.” They argue that the closed-source IOTA coordinator would protect the IOTA network from these purposefully introduced flaws, which they refer to as a “copy-protection mechanism” [Fou18c]. Based on this statement, IOTA seems to indicate that our research, in addition to discovering a novel attack on the IOTA Signature Scheme, may have uncovered an intentionally-placed backdoor in the cryptography of IOTA.

## 8 Conclusion

This paper presents chosen-message signature forgery attacks on the IOTA Signature Scheme when using the hash function Curl-P-27. We explain the cryptanalysis methods we used to create full-state collisions on same-length messages which differ in only a single position. We describe how to use these methods to create two valid IOTA bundles which can differ in multiple positions but still hash to the same value, and thus a signature for one is a valid signature for the other. We give examples placing these differences in the value fields of a bundle, and show that an attacker can produce such bundles using easily-accessible hardware in tens of seconds.

## Acknowledgments

The authors would like to thank Andy Sellars, Joi Ito, Vincenzo Iozzo, Sharon Goldberg, and Ward Heilman for feedback and guidance.

The research leading to these results has received funding from: the Ethics and Governance of Artificial Intelligence Fund, funders of the MIT Digital Currency Initiative, and US NSF grant 1350733.

## References

- [BDE<sup>+</sup>11] Johannes Buchmann, Erik Dahmen, Sarah Ereth, Andreas Hülsing, and Markus Rückert. On the security of the winternitz one-time signature scheme. In *International Conference on Cryptology in Africa*, pages 363–378. Springer, 2011.
- [BDPVA08] Guido Bertoni, Joan Daemen, Michael Peeters, and Gilles Van Assche. On the indistinguishability of the sponge construction. *Lecture Notes in Computer Science*, 4965:181–197, 2008.
- [Bit17] Bitfinex. Iota protocol upgrade august 08, 2017, 2017. <https://www.bitfinex.com/posts/215>, archived at <https://web.archive.org/web/20180722235151/https://www.bitfinex.com/posts/215>.
- [blo18] Tangle blog. Full emails of ethan heilman and the digital currency initiative with the iota team leaked, 2018. <http://www.tangleblog.com/wp-content/uploads/2018/02/letters.pdf>, archived at <https://web.archive.org/web/20180228182122/http://www.tangleblog.com/wp-content/uploads/2018/02/letters.pdf>.

- [//www.tangleblog.com/wp-content/uploads/2018/02/letters.pdf](http://www.tangleblog.com/wp-content/uploads/2018/02/letters.pdf),  
<https://archive.is/6imWR>.
- [BR96] Mihir Bellare and Phillip Rogaway. The exact security of digital signatures—how to sign with rsa and rabin. In *International Conference on the Theory and Applications of Cryptographic Techniques*, pages 399–416. Springer, 1996.
- [BS91] Eli Biham and Adi Shamir. Differential cryptanalysis of des-like cryptosystems. *Journal of CRYPTOLOGY*, 4(1):3–72, 1991.
- [Coi19] CoinmarketCap. Coinmarketcap iota nov 8 2019, 2019. <https://coinmarketcap.com/currencies/iota/historical-data/?start=20130428&end=20191109/>, archived at <https://coinmarketcap.com/currencies/iota/historical-data/?start=20130428&end=20191109>.
- [Cop94] Don Coppersmith. The data encryption standard (des) and its strength against attacks. *IBM journal of research and development*, 38(3):243–250, 1994.
- [CT18] Michael Colavita and Garrett Tanzer. A cryptanalysis of IOTA’s curl hash function, 2018. <https://www.boazbarak.org/cs127/Projects/iota.pdf>.
- [DRUH18] Gerard De Roode, Ikram Ullah, and Paul JM Havinga. How to break iota heart by replaying? In *2018 IEEE Globecom Workshops (GC Wkshps)*, pages 1–7. IEEE, 2018.
- [Fou18a] IOTA Foundation. Iota guide – generating secure multisig addresses (hot and coldwallet), 2018. <https://domschiener.gitbooks.io/iota-guide/content/exchange-guidelines/generating-multisignature-addresses.html>, archived at <https://archive.is/087kP>.
- [Fou18b] IOTA Foundation. Iota guide – seeds and accounts, 2018. <https://domschiener.gitbooks.io/iota-guide/content/chapter1/seeds-private-keys-and-addresses.html>, archived at <https://archive.is/yZf16>.
- [Fou18c] IOTA Foundation. Official iota foundation response to the digital currency initiative at the mit media lab – part 4 / 4, 2018. <https://blog.iota.org/official-iota-foundation-response-to-the-digital-currency-initiative-at-the-mit-media-lab-part-4-11fdccc9eb6d>, archived at <http://web.archive.org/web/20180727155405/https://blog.iota.org/official-iota-foundation-response-to-the-digital-currency-initiative-at-the-mit-media-lab-part-4-11fdccc9eb6d?gi=4be3ca82ed48>.
- [Fou18d] IOTA Foundation. Security levels. IOTA Basics 1.0, 2018. <https://docs.iota.org/docs/iota-basics/0.1/references/security-levels>.
- [GBH18] Leon Groot Bruinderink and Andreas Hülsing. “oops, i did it again” – security of one-time signatures under two-message attacks. In *Selected Areas in Cryptography – SAC 2017*, pages 299–322, 2018.
- [GJ05] Guang Gong and Shaoquan Jiang. The editing generator and its cryptanalysis. *International Journal of Wireless and Mobile Computing*, 1(1):46–52, 2005.
- [GJMG11] Bertoni Guido, Daemen Joan, P Michaël, and VA Gilles. Cryptographic sponge functions, 2011.

- [Gol04] Oded Goldreich. *Foundations of Cryptography: Basic Applications*, volume 2. Cambridge University Press, New York, NY, USA, 2004.
- [Han17] Paul Handy. Merged kerl implementation, 2017. <https://github.com/iotaledger/iri/commit/539e413352a77b1db2042f46887e41d558f575e5>, archived at <https://archive.is/jCisX>.
- [HNDV17] Ethan Heilman, Neha Narula, Thaddeus Dryja, and Madars Virza. Iota vulnerability report: Cryptanalysis of the curl hash function enabling practical signature forgery attacks on the iota cryptocurrency, 2017. <https://github.com/mit-dci/tangled-curl/blob/master/vuln-iota.md>.
- [IOT17] IOTAledger. IOTA kerl specification, 2017. <https://github.com/iotaledger/kerl/blob/master/IOTA-Kerl-spec.md>, archived at <https://web.archive.org/web/20180617175320/https://github.com/iotaledger/kerl/blob/master/IOTA-Kerl-spec.md>.
- [KL14] Jonathan Katz and Yehuda Lindell. *Introduction to Modern Cryptography*. CRC Press, second edition edition, 2014.
- [KTDB19] Stefan Kölbl, Elmar Tischhauser, Patrick Derbez, and Andrey Bogdanov. Troika: a ternary cryptographic hash function. *Designs, Codes and Cryptography*, pages 1–27, 2019.
- [Lam79] Leslie Lamport. Constructing digital signatures from a one-way function. Technical report, Technical Report CSL-98, SRI International Palo Alto, 1979.
- [LI19] Fukang Liu and Takanori Isobe. Preimage attacks on reduced troika with divide-and-conquer methods. In *International Workshop on Security*, pages 306–326. Springer, 2019.
- [LWdW05] Arjen K. Lenstra, Xiaoyun Wang, and Benne de Weger. Colliding x.509 certificates. Cryptology ePrint Archive, Report 2005/067, 2005. <https://eprint.iacr.org/2005/067>.
- [Mer89] Ralph C Merkle. A certified digital signature. In *Conference on the Theory and Application of Cryptology*, pages 218–238. Springer, 1989.
- [Pin18] Willem Pinckaers. (lekkertech) IOTA signatures, private keys and address reuse?, 2018. <http://blog.lekkertech.net/blog/2018/03/07/iota-signatures/>, archived at <https://archive.is/CnydQ>.
- [Pop16] Serguei Popov. The tangle. *cit. on*, page 131, 2016.
- [Rot18] Ralf Rottmann. Iota reclaim identification verification process, 2018. <https://blog.iota.org/iota-reclaim-identification-verification-process-e316647e06e6>, archived at <https://web.archive.org/web/20180710000243/https://blog.iota.org/iota-reclaim-identification-verification-process-e316647e06e6?gi=b8190e111e7f>.
- [Sch18] Dominik Schiener. Iota multi-signature scheme, 2017 (accessed February 3, 2018). <https://github.com/iotaledger/wiki/blob/master/multisigs.md>IOTA Multi-Signature Scheme.
- [SLdW07] Marc Stevens, Arjen Lenstra, and Benne de Weger. Chosen-prefix collisions for md5 and colliding x. 509 certificates for different identities. In *Advances in Cryptology – EUROCRYPT 2007*, pages 1–22. Springer, 2007.



- [Søn16] David Sønstebø. Iota tech update. Technical discussions, 2016. <https://forum.iota.org/t/iota-tech-update/264>.
- [Søn17] David Sønstebø. Upgrades & updates, 2017. <https://blog.iota.org/upgrades-updates-d12145e381eb>, archived at <https://web.archive.org/web/20180722232608/https://blog.iota.org/upgrades-updates-d12145e381eb?gi=51123f82db22>.
- [SSA<sup>+</sup>09] Marc Stevens, Alexander Sotirov, Jacob Appelbaum, Arjen Lenstra, David Molnar, Dag Arne Osvik, and Benne de Weger. Short chosen-prefix collisions for md5 and the creation of a rogue ca certificate. In *Advances in Cryptology – CRYPTO 2009*, pages 55–69. Springer, 2009.
- [SZ15] Yonatan Sompolinsky and Aviv Zohar. Secure high-rate transaction processing in bitcoin. In *International Conference on Financial Cryptography and Data Security*, pages 507–527. Springer, 2015.
- [Wal17] Eric Wall. IOTA is centralized, 2017. <https://medium.com/@ercwl/iota-is-centralized-6289246e7b4d>, archived at <https://web.archive.org/web/20180616231657/https://medium.com/@ercwl/iota-is-centralized-6289246e7b4d>.
- [WFLY04] Xiaoyun Wang, Dengguo Feng, Xuejia Lai, and Hongbo Yu. Collisions for hash functions md4, md5, haval-128 and ripemd. Cryptology ePrint Archive, Report 2004/199, 2004. <https://eprint.iacr.org/2004/199>.
- [Wie05] Michael J. Wiener. Bounds on birthday attack times. Cryptology ePrint Archive, Report 2005/318, 2005. <https://eprint.iacr.org/2005/318>.
- [WY05] Xiaoyun Wang and Hongbo Yu. How to break md5 and other hash functions. In *Advances in Cryptology – EUROCRYPT 2005*, pages 19–35. Springer, 2005.

## A Security Definitions

In this Section we provide a formal definition of the EU-CMA (Existential Unforgeability under Chosen Message Attack) and EU-RMA (Existential Unforgeability under Random Message Attack) for ISS (IOTA Signature Scheme). We use EU-RMA to model *known message attacks* where the known message is chosen at random. We discuss extending these definitions to include *chosen message attacks* which are limited to messages which are also valid IOTA bundles.

We briefly recall the standard definitions of digital signature schemes and their security, mirroring standard literature [Gol04, KL14, GBH18].

**Definition 1.** A digital signature scheme is a triple of algorithms (KeyGen, Sign, Verify) working as follows:

- $\text{KeyGen}(1^\lambda) \rightarrow (\text{PK}, \text{SK})$ : On input  $1^\lambda$ , the *key generator* KeyGen outputs a keypair  $(\text{PK}, \text{SK})$ , consisting of a public key PK and secret key SK.
- $\text{Sign}(\text{SK}, \text{msg}) \rightarrow \sigma$ : On input SK and message msg, the *signing algorithm* Sign outputs a signature  $\sigma$ .
- $\text{Verify}(\text{PK}, \text{msg}, \sigma) \rightarrow b$ : On input PK, msg,  $\sigma$ , the *verification algorithm* Verify outputs a decision bit  $b$ .

We require *perfect completeness*. That is, for any message  $\text{msg}$ , a validly generated signature must always be accepted:

$$\Pr \left[ \text{Verify}(\text{PK}, \text{msg}, \sigma) = 1 \mid \begin{array}{l} (\text{PK}, \text{SK}) \leftarrow \text{KeyGen}(1^\lambda) \\ \sigma \leftarrow \text{Sign}(\text{SK}, \text{msg}) \end{array} \right] = 1 .$$

The setting relevant for our work in Section 5 is a *chosen message attack* where, before outputting a forgery, an adversary gets to learn signatures for messages of his choice. In particular, for one-time signature schemes, the adversary is able to learn a signature for a single message. Other attack models include random message attacks (adversary is able to learn signatures on random message(s)) which we define to model a known message attack. Key-only attacks (adversary is only able to learn the public key). In particular, prior work offered a key-only attack of IOTA [Pin18] succeeding with probability 1 for  $\approx 3\%$  public keys. Our chosen message attack exploiting vulnerabilities in Curl-P-27 succeeds with probability 1 (under reasonable heuristics) for all public keys.

**Definition 2.** A one-time signature scheme ( $\text{KeyGen}, \text{Sign}, \text{Verify}$ ) is secure against chosen message attacks (or, EU-CMA secure), if no polynomial-time stateful adversary  $\mathcal{A} = (\mathcal{A}_1, \mathcal{A}_2)$  can output a fresh forgery, except with negligible probability:

$$\Pr \left[ \begin{array}{l} (\text{msg} \neq \text{msg}') \\ \text{and} \\ \text{Verify}(\text{PK}, \text{msg}', \sigma') = 1 \end{array} \mid \begin{array}{l} (\text{PK}, \text{SK}) \leftarrow \text{KeyGen}(1^\lambda) \\ (\text{msg}, \text{st}) \leftarrow \mathcal{A}_1(\text{PK}) \\ \sigma \leftarrow \text{Sign}(\text{SK}, \text{msg}) \\ (\text{msg}', \sigma') \leftarrow \mathcal{A}_2(\text{st}, \text{msg}, \sigma) \end{array} \right] \approx \text{negl}(\lambda) .$$

Notably, ISS uses a hash-then-sign paradigm. Therefore, ability to find collisions for the underlying hash function directly yields to a chosen-message attack. Given two colliding inputs  $x_1$  and  $x_2$  ( $\mathcal{H}(x_1) = \mathcal{H}(x_2)$ ), the adversary will first output  $\text{msg} = x_1$ , and after receiving a correct signature  $\sigma = \text{Sign}(\text{SK}, \text{msg})$  from the challenger, will output  $\text{msg}' = x_2$  and  $\sigma' = \sigma$ . Because both messages yield the same value when hashed,  $\text{Verify}(\text{PK}, \text{msg}', \sigma')$  always accepts.

One could argue that the practical impact of not achieving EU-CMA security is hard to quantify. Indeed, if other parts of the overall system required messages to have certain structure, yet an attacker was only able to produce forgeries for messages that lack this requisite structure, these cryptographic breaks might not yield attacks exploitable against valid messages.

We rule out this possibility by devising a highly flexible collision-finding algorithm (see Section 5). In particular, our algorithm lets us freely generate collisions with arbitrary known prefixes and suffixes. This is sufficient, for example, to generate two colliding IOTA transactions, or two colliding multi-signature bundles. We explain how to generate such forgeries in Appendix B, and experimentally validate that the IOTA reference implementation from August 6, 2017, which is before the commit to change the hash function from Curl-P to Keccak [Han17], accepts the forged signatures.

Section 6 discusses *existential unforgeability under random message attack* (EU-RMA) against ISS, where the message is chosen randomly *i.e.*, sampled uniformly. Note that for one-time signatures the adversary  $\mathcal{A}$  is no longer stateful as it does not require any interacting with the security game beyond receiving  $(\text{PK}, \text{msg}, \sigma)$  triple for randomly chosen  $\text{msg}$ , and outputting a forgery.

**Definition 3.** A one-time signature scheme ( $\text{KeyGen}, \text{Sign}, \text{Verify}$ ) is secure against random message attacks (or, EU-RMA secure), if no polynomial-time adversary  $\mathcal{A}$ , given a signature on a random message can output a fresh forgery, except with negligible probability:

$$\Pr \left[ \begin{array}{l} (\text{msg} \neq \text{msg}') \\ \text{and} \\ \text{Verify}(\text{PK}, \text{msg}', \sigma') = 1 \end{array} \mid \begin{array}{l} (\text{PK}, \text{SK}) \leftarrow \text{KeyGen}(1^\lambda) \\ \text{msg} \xleftarrow{\$} \{-1, 0, 1\}^{\text{poly}(\lambda)} \\ \sigma \leftarrow \text{Sign}(\text{SK}, \text{msg}) \\ (\text{msg}', \sigma') \leftarrow \mathcal{A}(\text{PK}, \text{msg}, \sigma) \end{array} \right] \approx \text{negl}(\lambda) .$$

Txn	Address	Tag	$b_1$ Value	$b_2$ Value
0	Bob	WJ9JPWIYIVQSTFNYY9HCZUQRVBK	182219672	182219672
1	Eve	CYBLAAX9ZA99Q9ZU9CXIU9DXCCW	1	129140164
2	Carol	GOUKGHTRFTRLRHPOBZRMDLM9QIEM	400000	400000
3	Alice	FZZMZXCXWAI9SZAURCR9C9BXDCW	129140164	1
4	Alice,Eve	99999999999999999999999999999999	-500000000	-500000000
5		99999999999999999999999999999999	0	0
6		99999999999999999999999999999999	0	0
7		99999999999999999999999999999999	0	0
8		99999999999999999999999999999999	0	0
9		99999999999999999999999999999999	0	0
10	Alice,Eve	99999999999999999999999999999999	188240163	188240163

**Figure 8:** An example multisig bundle spending from an address controlled by Alice and Eve. Transaction 4 is the funding transaction, and transaction 10 is the change which goes back to Alice and Eve. The collisions are in transactions 1 and 3. Transactions 5-9 are just to hold the signature fragments of Alice and Eve.

The definitions given in this section are intended as a helpful guide for the security assumptions of ISS. Because these definitions are asymptotic definitions they do not cleanly map onto our practical attacks which deal with concrete computational resources. We hope that follow up research will extend this work to develop concrete definitions for security of ISS such as those given for RSA in [BR96].

## B Example Colliding Bundles

Our previous vulnerability report detailed colliding bundles which execute the steal money and waste money attacks for single-signature bundles. In Section 5.2, we described the structure of these attacks and showed how we targeted collisions in value fields using the previous and following tags. Here, we describe in more detail example multi-signature bundles that have the same hash and spend different amounts to Alice and Eve.

Figure 8 describes two different bundles  $b_1$  and  $b_2$  that differ only in the value fields in two transactions. We give the address the funds are being spent to or from, the tags, and the values. Each bundle consists of 11 transactions: one input spending funds that were in a multisig address controlled by Alice and Eve (4), five outputs, including one change output, spending to Alice, Bob, Carol, Eve, and back to Alice and Eve’s multisig address (0-3,10), and five extra transactions which are present only to hold signature fragments from Alice and Eve authorizing the spend (5-9). Signature-holding transactions have empty addresses and values. The tags in first four transactions are generated using our collision tool so that  $b_1$  and  $b_2$  will have the same hash.