

# Subversion-Resilient Signatures: Definitions, Constructions and Applications

Giuseppe Ateniese, Bernardo Magri, and Daniele Venturi

*Department of Computer Science, Sapienza University of Rome*

October 30, 2015

## Abstract

We provide a formal treatment of security of digital signatures against *subversion attacks* (SAs). Our model of subversion generalizes previous work in several directions, and is inspired by the proliferation of software attacks (e.g., malware and buffer overflow attacks), and by the recent revelations of Edward Snowden about intelligence agencies trying to surreptitiously sabotage cryptographic algorithms. The main security requirement we put forward demands that a signature scheme should remain unforgeable even in the presence of an attacker applying SAs (within a certain class of allowed attacks) in a fully-adaptive and *continuous* fashion. Previous notions—e.g., the notion of security against algorithm-substitution attacks introduced by Bellare *et al.* (CRYPTO '14) for symmetric encryption—were non-adaptive and non-continuous.

In this vein, we show both positive and negative results for the goal of constructing subversion-resilient signature schemes.

- **Negative results.** As our main negative result, we show that a broad class of randomized signature schemes is unavoidably insecure against SAs, even if using just a single bit of randomness. This improves upon earlier work that was only able to attack schemes with larger randomness space. When designing our new attack we consider undetectability as an explicit adversarial goal, meaning that the end-users (even the ones knowing the signing key) should not be able to detect that the signature scheme was subverted.
- **Positive results.** We complement the above negative results by showing that signature schemes with *unique* signatures are subversion-resilient against all attacks that meet a basic undetectability requirement. A similar result was shown by Bellare *et al.* for symmetric encryption, who proved the necessity to rely on *stateful* schemes; in contrast unique signatures are *stateless*, and in fact they are among the fastest and most established digital signatures available. As our second positive result, we show how to construct subversion-resilient identification schemes from subversion-resilient signature schemes. We finally show that it is possible to devise signature schemes secure against arbitrary tampering with the computation, by making use of an un-tamperable cryptographic reverse firewall (Mironov and Stephens-Davidowitz, EUROCRYPT '15), i.e., an algorithm that “sanitizes” any signature given as input (using only public information). The firewall we design allows to successfully protect so-called re-randomizable signature schemes (which include unique signatures as special case).

As an additional contribution, we extend our model to consider multiple users and show implications and separations among the various notions we introduced. While our study is mainly theoretical, due to its strong practical motivation, we believe that our results have important implications in practice and might influence the way digital signature schemes are selected or adopted in standards and protocols.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>	<b>5.3</b>	Random-Message Attacks . . . . .	18
1.1	Our Results and Techniques . . . . .	2	<b>5.4</b>	Subversion-Resilient Identification Schemes . . . . .	20
1.2	Impact . . . . .	4	<b>6</b>	<b>Reverse Firewalls for Signatures</b>	<b>22</b>
1.3	Related Work . . . . .	5	6.1	Properties . . . . .	23
<b>2</b>	<b>Preliminaries</b>	<b>7</b>	6.2	Necessity of Self-Destruct . . . . .	24
2.1	Notation . . . . .	7	6.3	Patching Re-Randomizable Signatures . . . . .	25
2.2	Signature Schemes . . . . .	7	<b>7</b>	<b>The Multi-User Setting</b>	<b>27</b>
2.3	Pseudorandom Functions . . . . .	8	7.1	Multi-User Impersonation . . . . .	27
<b>3</b>	<b>Subverting Signatures</b>	<b>8</b>	7.2	Multi-User Public/Secret Undetectability . . . . .	28
3.1	Impersonation . . . . .	9	7.3	Impersonation Relations . . . . .	29
3.2	Public/Secret Undetectability . . . . .	10	7.4	Undetectability Relations . . . . .	30
<b>4</b>	<b>Mounting Subversion Attacks</b>	<b>11</b>	<b>8</b>	<b>Mounting Multi-User SAs</b>	<b>32</b>
4.1	Attacking Coin-Injective Schemes . . . . .	11	8.1	Attacking Coin-Injective Schemes (Multi-User Version) . . . . .	32
4.2	Attacking Coin-Extractable Schemes . . . . .	14	8.2	Attacking Coin-Extractable Schemes (Multi-User Version) . . . . .	34
<b>5</b>	<b>Security of Unique Signatures</b>	<b>16</b>			
5.1	The Verifiability Condition . . . . .	16			
5.2	Chosen-Message Attacks . . . . .	17			

## 1 Introduction

Balancing national security interests with the rights to privacy of lawful citizen is always a daunting task. It has been particularly so in the last couple of years after the revelations of Edward Snowden [PLS13, BBG13, Gre14] that have evidenced a massive collection of metadata and other information perpetrated by several intelligence agencies. It is now clear that intelligence operators were not just interested in collecting and mining information but they also actively deployed malware, exploited zero-day vulnerabilities, and carried out active attacks against standard protocols. In addition, it appears some cryptographic protocol specifications were modified to embed backdoors.

Whether this activity was effective or even allowed by the constitution is open to debate and it is indeed being furiously discussed among policy makers, the public, and the intelligence community. Ultimately, a balance between security and privacy must be found for a free and functioning society.

The ability of substituting a cryptographic algorithm with an altered version was first considered formally by Young and Yung (extending previous works of Simmons on subliminal channels [Sim83, Sim84]), who termed this field *kleptography* [YY96, YY97]. The idea is that the attacker surreptitiously modifies a cryptographic scheme with the intent of subverting its security. This research area has recently been revitalized by Bellare *et al.* [BPR14] who considered encryption algorithms with the possibility of mass surveillance under the algorithm-substitution attack. They analyzed the possibility of an intelligence agency substituting an encryption algorithm with the code of an alternative version that undetectably reveals the secret key or the plaintext. What they uncovered is that any randomized and stateless encryption scheme would fall to generic algorithm-substitution attacks. The only way to achieve a meaningful security

guarantee (CPA-security) is to use a nonce-based encryption that must keep state. Unfortunately, only stateless schemes are deployable effectively with the current network technology and indeed all deployed encryption algorithms are in this class.

In this paper we analyze digital signature schemes under the so-called *subversion attacks* (SAs), that in particular include algorithm-substitution and kleptographic attacks as a special case, but additionally cover more general malware and virus attacks (see below). Unlike encryption, we show positive results and truly efficient schemes that provide the strongest security guarantee and can thus be deployed within real systems. We stress that our intention is not to propose schemes that can be abused by criminals to avoid monitoring. We are motivated by pure scientific curiosity and aspire to contribute to an active field of research.

## 1.1 Our Results and Techniques

We introduce a new and generic framework and definitions for subversions of digital signatures. In the standard black-box setting, a signature scheme should remain unforgeable even against an adversary able to obtain signatures on (polynomially many) chosen messages. Our security definitions empower the adversary with the ability of *continuously* subverting the signing algorithm within a class  $\mathcal{A}$  of allowed SAs. For each chosen subversion in the class, the adversary can access an oracle that answers (polynomially many) signature queries using the subverted signature algorithm. Importantly, the different subversions can be chosen in a fully-adaptive manner possibly depending on the target verification key of the user.

We believe our model is very general and flexible, as it nicely generalizes previous models and definitions. First off, when the class  $\mathcal{A}$  consists of a set of algorithms containing a secretly embedded backdoor, and in case the adversary is restricted to non-adaptively choose only a single subversion algorithm from this class, we obtain the setting of algorithm-substitution and kleptographic attacks as a special case. However, we note that the above definition is far more general as it covers (fully-adaptive and continuous) *tampering with the computation* performed by the signing algorithm (within the class  $\mathcal{A}$ ). This models, for instance, a machine running a signature software infected by a malware (e.g., via a buffer overflow attack [One96, Fry00, PB04]); we also obtain memory and randomness tampering (see Section 1.3) as a special case. We refer the reader to Section 3.1 (where we introduce our model formally) for a more comprehensive discussion.

Clearly, without making any restriction on the class  $\mathcal{A}$  (or without making additional assumptions) there is no hope for security: An arbitrary subverted signature algorithm could, for instance, just ignore all inputs and output the secret key. In this paper we investigate two approaches to tackle attacks of this sort and obtain positive results.

- **Limiting the adversarial power.** We consider a setting where the adversarial goal is to subvert the signature algorithm in a way that is *undetectable* to the end-user (or at least allows to maintain plausible deniability). For instance the simple attack above—where the subversion outputs the secret key—is easily detectable given only public information. As we show in Section 5, requiring that the class  $\mathcal{A}$  satisfies a basic undetectability requirement already allows for interesting positive results.
- **Using a Reverse Firewall.** In Section 6 we show that security against *arbitrary* tampering with the computation can be achieved, by making the additional assumption of an un-tamperable cryptographic reverse firewall (RF) [MS15]. Roughly, a RF takes as input a message/signature pair and is allowed to “sanitize” the input signature using only *public information*.

A more detailed description of our techniques follows.

**Negative results.** We define what it means for a class  $\mathcal{A}$  of SAs to be (efficiently) undetectable; roughly this means that a user, given polynomially many queries, cannot distinguish the output of the genuine signature algorithm from the output of the subverted algorithm. See Section 3.2 for a precise definition. Our definitions of undetectability are similar in spirit to the ones put forward by [BPR14] for the setting of symmetric encryption. Importantly we distinguish the case where the user (trying to detect the attack) knows only public or private information (i.e., it knows the secret key).<sup>1</sup>

Next, we explore the possibility of designing classes of SAs that are (even secretly) undetectable and yet allow for complete security breaches. This direction was already pursued by Bellare *et al.*, who showed that it is possible to stealthily bias the random coins of sufficiently randomized symmetric encryption schemes in a way that allows to extract the secret key after observing a sufficient number of (subverted) ciphertexts. As a first negative result, we explain how to adapt the “biased randomness attack” of [BPR14] to the case of signature schemes.

The above generic attack requires that the signature scheme uses a minimal amount of randomness (say, 7 bits). This leaves the interesting possibility that less randomized schemes (such as the Katz-Wang signature scheme [KW03], using only one bit of randomness) might be secure. In Section 4, we present a new attack showing that this possibility is vacuous: Our attack allows to stealthily bias the randomness in a way that later allows to extract the signing key—regardless of the number of random bits required by the scheme—assuming that the targeted signature scheme is *coin-extractable*. The latter roughly means that the random coins used for generating signatures can be extracted efficiently from the signature itself; as we discuss in more detail in Section 4.2 many real schemes (including Katz-Wang) are coin-extractable.

**Positive results.** As a first positive result we show that fully deterministic schemes with *unique*<sup>2</sup> signatures are existentially unforgeable under chosen-message attacks against the class of SAs that satisfies the so-called verifiability condition.<sup>3</sup> This means that—for all values in the message space—signatures produced by the subverted signature algorithm should (almost always) verify correctly under the target verification key (note that both attacks mentioned above fall into this category).

Clearly, the assumption that the verifiability condition should hold for all messages is quite a strong one. Hence, we also relax the verifiability condition to hold for all but a negligible fraction of the messages. However, we are not able to prove that unique signatures achieve existential unforgeability under chosen-message attacks against the class of SAs that satisfies relaxed verifiability.<sup>4</sup> Instead, as our second positive result, we show that unique signatures are existentially unforgeable under random-message attacks (where the adversary can only see potentially subverted signatures of random messages) against the class of SAs that satisfies relaxed verifiability. Interestingly, this weaker security flavor is still useful for applications, e.g. to construct subversion-resilient identification schemes.

As our third positive result, we provide a way how to achieve the ambitious goal of protecting signature schemes against *arbitrary* SAs, relying on a cryptographic reverse firewall. The latter

---

<sup>1</sup>As we show, secret and public undetectability are *not* equivalent, in that there exist natural classes of SAs that are publicly undetectable but secretly detectable.

<sup>2</sup>A signature scheme is unique if for a honestly generated verification key there is a single valid signature for each message.

<sup>3</sup>One might ask whether a similar result holds for all deterministic schemes where signatures are not unique; the answer to this question is negative as our attacks also apply to certain types of deterministic schemes (e.g., de-randomized schemes—see the proof of Theorem 9 in Section 7.4).

<sup>4</sup>In fact, as shown very recently by Degabriele *et al.* [DFP15] for the case of symmetric encryption, it is not hard to show that such limitation is inherent: No (even deterministic) scheme can achieve security under chosen-message attacks against the class of SAs that meets relaxed verifiability. See Section 1.3 for more details.

primitive was recently introduced in [MS15] (see also [DMS15]) to model security of arbitrary two-party protocols run on machines possibly corrupted by a virus. On a high level a RF for a signature scheme is a piece of software taking as input a message/signature pair  $(m, \sigma)$  and some *public* state, and outputting a “patched” signature  $(m, \sigma')$ ; the initial state of the firewall is typically a function of the verification key  $vk$ . A good RF should maintain functionality, meaning that whenever the input is a valid message/signature pair the patched signature (almost always) verifies correctly under the target verification key. Moreover, we would like the firewall to preserve unforgeability; this means that patched signatures (corresponding to signatures generated via the subverted signing algorithm) should not help an adversary to forge on a fresh message.

We prove that every signature scheme that is re-randomizable (as defined in [HJK12]) admits a RF that preserves unforgeability against arbitrary SAs. Re-randomizable signatures admit an efficient algorithm *ReRand* that takes as input a tuple  $(m, \sigma, vk)$  and outputs a signature  $\sigma'$  that is distributed uniformly over the set of all valid signatures on message  $m$  (under  $vk$ ); unique signatures, for instance, are re-randomizable. Upon input a pair  $(m, \sigma)$  our firewall uses the public state to verify  $(m, \sigma)$  is valid under  $vk$ , and, in case the test passes, it runs *ReRand* on  $(m, \sigma)$  and outputs the result. Otherwise the firewall simply returns an invalid symbol  $\perp$  and *self-destructs*, i.e., it stops processing any further query.<sup>5</sup> The latter is a requirement that we prove to be unavoidable: No RF can at the same time maintain functionality and preserve unforgeability of a signature scheme without the self-destruct capability.

We remark that our results and techniques for the setting of RFs are incomparable to the ones in [MS15]. The main result of Mironov and Stephens-Davidowitz is a compiler that takes as input an arbitrary two-party protocol and outputs a functionally equivalent (but different) protocol that admits a RF preserving both functionality and security. Instead, we model directly security of RFs for signatures schemes in the game-based setting; while our goal is more restricted (in that we only design RFs for signatures), our approach results in much more efficient and practical solutions.

**Multi-user setting.** Our discussion so far considered a single user. In Section 7 we discuss how our models and results can be extended to the important (and practically relevant) multi-user scenario. In particular, similarly to [BPR14], we generalize our undetectability and security notions to a setting with  $u > 1$  users, where each user has a different signing/verification key.

As we argue, security in the single-user setting already implies security in the multi-user setting (by a standard hybrid argument). This does not hold for undetectability, as there exists classes of SAs that are undetectable by a single user but can be efficiently detected by more than one user. However, as we show in Section 8, the concrete attacks analysed in Section 4 can be modified to remain undetectable even with multiple users.

## 1.2 Impact

Our study has strong implications in practice and might influence the way digital signature schemes are selected or adopted in standards and protocols. A subverted signature scheme is arguably even more deceitful and dangerous in practice than subverted encryption. Indeed, it is well-known that authenticated encryption must involve digital certificates that are signed by Certification Authorities (CAs). If a CA is using a subverted signature scheme, it is reasonable to expect the signing key will eventually be exposed. With knowledge of the signing key, it is possible to impersonate any user and carry out elementary man-in-the-middle attacks. This

---

<sup>5</sup>This can be implemented, for instance, by having the public state include a single one-time writable bit used to signal a self-destruct took place.

renders the use of any type of encryption utterly pointless and underlines the important role played by signatures in the context of secure communications.

Unfortunately, signature schemes currently employed to sign digital certificates, or used in protocols such as OTR, TLS/SSL, SSH, etc., are all susceptible to a subversion attack and their use should possibly be discontinued. The positive news however is that there already exist signature schemes that are subversion-resilient and they are efficient and well-established. This is in contrast with encryption where *good* schemes are not deployable in all contexts since they require retention of state information (see [BPR14]).

### 1.3 Related Work

Sabotage of cryptographic primitives before and during their deployment has been the focus of extensive research over the past years. We briefly review the main results below.

**Subliminal channels and backdoored implementations.** After their introduction, the potential of subliminal channels has been explored in several works (e.g., [Des88a, Des88b, BDI<sup>+</sup>99]); this line of research lead for instance to the concept of divertible protocols, that are intimately related to reverse firewalls.

The setting of backdoored implementations has also been the focus of extensive research. This includes, in particular, the realm of kleptography and SETUP attacks (see [YY04] for a survey). In recent work, Dodis *et al.* [DGG<sup>+</sup>15] provide a formal treatment of trapdoored pseudorandom generators (building on previous work of Vazirani and Vazirani [VV83]); this setting is of particular importance, given the potential sabotage of the NIST Dual EC PRG [NIS07]. Russell *et al.* [RTYZ15] consider the setting of *complete* subversion, where all algorithms (including for instance the key generation algorithm) are subject to kleptographic attacks, and show how to build (trapdoor) one-way functions in this model.

We refer the reader to [SFKR15] for a taxonomy of these (and more) types of attacks.

**Stateless subversion attacks** Bellare *et al.* [BJK15] introduced a stronger definition of undetectability where the user trying to detect a subverted scheme can query messages to the signing oracle and receive as output a pair  $(\sigma, \tau)$  containing the signature and the state; since the original signature scheme is stateless the state  $\tau$  will always be empty (represented by the empty string  $\epsilon$ ) when querying the original signing algorithm. This way, any stateful subversion can be easily detected by checking whether  $\tau = \epsilon$ .

Moreover, the “biased randomness attack” from [BPR14] is modified into a stateless version, so it can remain undetectable under this stronger undetectability definition. We note that the attack from Fig. 1 (also shown in [AMV15]) can also be made stateless by using the same techniques.

**Input-triggered subversions.** In a very recent paper, Degabriele, Farshim and Poettering (DFP) [DFP15] pointed out some shortcomings of the Bellare-Patterson-Rogaway (BPR) [BPR14] security model for subversion resilience of symmetric encryption schemes. Consider the class of SAs that upon input a secret (trapdoor) message  $\bar{m}$  outputs the secret key, but otherwise behaves like the genuine signature algorithm. Clearly this class of SAs will be undetectable by the users, as without knowing the trapdoor there is only a negligible chance to query the secret message  $\bar{m}$  and check if the signature algorithm was subverted (at least if the message space is large enough). Yet, an adversary mounting a chosen-message attack can recover the signing key by asking a signature for message  $\bar{m}$ .

As a consequence, it is impossible to prove existential unforgeability under-chosen message attacks against such “input-triggered” subversions (in the BPR model). Note however that, for the case of signatures, one can still prove a positive result by restricting the adversary to only see signatures of random messages (i.e., in case of a random-message attack). Indeed, input-triggered subversions meet our notion of relaxed verifiability (see Section 1.1) and thus our positive results for unique signatures apply to such case.

The solution proposed by DFP is to modify the definition of undetectability so that the adversary (and not the user) specifies the input messages to the (potentially subverted) encryption algorithm, whereas the goal of the user is to detect the attack given access to the transcript of all queries made by the adversary (and answers to these queries). Hence, a scheme is said to be subversion-resilient if there exists a fixed polynomial-time test algorithm such that either a subversion attack cannot be detected efficiently but it does not leak any useful information, or it is possible to efficiently detect that the system was subverted.<sup>6</sup>

It is possible to make a similar change as in [DFP15] and adapt the DFP model to signature schemes in order to achieve security under chosen-message attacks. The end result would share some similarities with our approach using cryptographic RFs;<sup>7</sup> however, our framework provides notable advantages. First, note that the DFP model does not provide any guarantee against SAs that are efficiently detectable, whereas our RF model explicitly accounts for the actions to be taken after an attack is detected; this is particularly relevant for signature schemes where our generic attack uncovered the necessity of a self-destruct capability. Second, the polynomial-time detection test in DFP is performed directly by the user since it requires knowledge of the secret key. This is problematic in practice since often the user’s machine is completely compromised; instead, in our framework, a cryptographic RF for a signature scheme relies only on public information and could easily be located on a (untrusted) external proxy.

**Tampering attacks.** A related line of research analyzes the security of cryptosystems against tampering attacks. Most of these works are restricted to the simpler setting of memory tampering (sometimes known as related-key security), where only the secret key of a targeted cryptoscheme is subject to modification. By now we know several concrete primitives that remain secure against different classes of memory-tampering attacks, including pseudorandom functions and permutations [BK03, Luc04, BC10, AFPW11, BCM11], pseudorandom generators and hard-core bits [GL10], hash functions [GOR11], public-key encryption [AHI11, Wee12], identification and digital signature schemes [KKS11, DFMV13]. Elegant generic compilers are also available, relying on so-called tamper-resilient encodings and non-malleable codes (see, among others, [GLM<sup>+</sup>04, DPW10, LL12, FMNV14, FMVW14, ADL14, JW15, DLSZ15, AGM<sup>+</sup>15, FMNV15, DFMV15]).

The setting of randomness tampering, where the random coins of a cryptographic algorithm are subject to tampering, has also been considered. For instance Austrin *et al.* [ACM<sup>+</sup>14] consider so-called  $p$ -tampering attacks, that can efficiently tamper with each bit of the random tape with probability  $p$ . In this setting they show that some cryptographic tasks (including commitment schemes and zero-knowledge protocols) are impossible to achieve, while other tasks (in particular signature and identification schemes) can be securely realized.

Yet another related setting is that of tampering attacks against gates and wires in the computation of a cryptographic circuit, and the design of tamper-proof circuit compilers [IPSW06, FPV11, DK12, KT13, DK14, GIP<sup>+</sup>14].

<sup>6</sup>For instance, in case of the attack outlined above, the polynomial-time test could simply decrypt the ciphertext and check the outcome matches the input message.

<sup>7</sup>On a high level, one can interpret the polynomial-time test as playing the role of the reverse firewall.

## 2 Preliminaries

### 2.1 Notation

For a string  $x$ , we denote its length by  $|x|$ ; if  $\mathcal{X}$  is a set,  $|\mathcal{X}|$  represents the number of elements in  $\mathcal{X}$ . When  $x$  is chosen randomly in  $\mathcal{X}$ , we write  $x \leftarrow_{\$} \mathcal{X}$ . When  $\mathbf{A}$  is an algorithm, we write  $y \leftarrow \mathbf{A}(x)$  to denote a run of  $\mathbf{A}$  on input  $x$  and output  $y$ ; if  $\mathbf{A}$  is randomized, then  $y$  is a random variable and  $\mathbf{A}(x; r)$  denotes a run of  $\mathbf{A}$  on input  $x$  and randomness  $r$ . An algorithm  $\mathbf{A}$  is *probabilistic polynomial-time* (PPT) if  $\mathbf{A}$  is randomized and for any input  $x, r \in \{0, 1\}^*$  the computation of  $\mathbf{A}(x; r)$  terminates in at most  $\text{poly}(|x|)$  steps.

We denote with  $\kappa \in \mathbb{N}$  the security parameter. A function  $\text{negl} : \mathbb{N} \rightarrow \mathbb{R}$  is negligible in the security parameter (or simply negligible) if it vanishes faster than the inverse of any polynomial in  $\kappa$ , i.e.  $\text{negl}(\kappa) = \kappa^{-\omega(1)}$ .

The statistical distance between two random variables  $\mathbf{A}$  and  $\mathbf{B}$  defined over the same domain  $\mathcal{D}$  is defined as  $\mathbb{SD}(\mathbf{A}; \mathbf{B}) = \frac{1}{2} \sum_{x \in \mathcal{D}} |\mathbb{P}[\mathbf{A} = x] - \mathbb{P}[\mathbf{B} = x]|$ . We rely on the following lemma (which follows directly from the definition of statistical distance):

**Lemma 1.** *Let  $\mathbf{A}$  and  $\mathbf{B}$  be a pair of random variables, and  $E$  be an event defined over the probability space of  $\mathbf{A}$  and  $\mathbf{B}$ . Then,*

$$\mathbb{SD}(\mathbf{A}; \mathbf{B}) \leq \mathbb{SD}(\mathbf{A}; \mathbf{B} | \neg E) + \mathbb{P}[E].$$

### 2.2 Signature Schemes

A signature scheme is a triple of algorithms  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  specified as follows: (i)  $\text{KGen}$  takes as input the security parameter  $\kappa$  and outputs a verification/signing key pair  $(vk, sk) \in \mathcal{VK} \times \mathcal{SK}$ , where  $\mathcal{VK} := \mathcal{VK}_{\kappa}$  and  $\mathcal{SK} := \mathcal{SK}_{\kappa}$  denote the sets of all verification and secret keys produced by  $\text{KGen}(1^{\kappa})$ ; (ii)  $\text{Sign}$  takes as input the signing key  $sk \in \mathcal{SK}$ , a message  $m \in \mathcal{M}$  and random coins  $r \in \mathcal{R}$ , and outputs a signature  $\sigma \in \Sigma$ ; (iii)  $\text{Vrfy}$  takes as input the verification key  $vk \in \mathcal{VK}$  and a pair  $(m, \sigma)$ , and outputs a decision bit that equals 1 iff  $\sigma$  is a valid signature for message  $m$  under key  $vk$ .

Correctness of a signature scheme says that verifying honestly generated signatures always works (with overwhelming probability over the randomness of all involved algorithms).

**Definition 1** (Correctness). Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a signature scheme. We say that  $\mathcal{SS}$  satisfies  $\nu_c$ -correctness if for all  $m \in \mathcal{M}$

$$\mathbb{P}[\text{Vrfy}(vk, (m, \text{Sign}(sk, m))) = 1 : (vk, sk) \leftarrow \text{KGen}(1^{\kappa})] \geq 1 - \nu_c,$$

where the probability is taken over the randomness of  $\text{KGen}$ ,  $\text{Sign}$ , and  $\text{Vrfy}$ .

The standard notion of security for a signature scheme demands that no PPT adversary given access to a signing oracle returning signatures for arbitrary messages, can forge a signature on a “fresh” message (not asked to the signing oracle).

**Definition 2** (Existential Unforgeability). Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a signature scheme. We say that  $\mathcal{SS}$  is  $(t, q, \varepsilon)$ -existentially unforgeable under chosen-message attacks ( $(t, q, \varepsilon)$ -EUF-CMA in short) if for all PPT adversaries  $\mathbf{A}$  running in time  $t$  it holds:

$$\mathbb{P}[\text{Vrfy}(vk, (m^*, \sigma^*)) = 1 \wedge m^* \notin \mathcal{Q} : (vk, sk) \leftarrow \text{KGen}(1^{\kappa}); (m^*, \sigma^*) \leftarrow \mathbf{A}^{\text{Sign}(sk, \cdot)}(vk)] \leq \varepsilon,$$

where  $\mathcal{Q} = \{m_1, \dots, m_q\}$  denotes the set of queries to the signing oracle. Whenever  $\varepsilon(\kappa) = \text{negl}(\kappa)$  and  $q = \text{poly}(\kappa)$ , we simply say that  $\mathcal{SS}$  is EUF-CMA.



**Unique signatures.** For our positive results we rely on so called *unique* signatures, that we define next. Informally a signature scheme is unique if for any message there is a single signature that verifies w.r.t. a honestly generated verification key.

**Definition 3** (Uniqueness). Let  $\mathcal{SS}$  be a signature scheme. We say that  $\mathcal{SS}$  satisfies  $\nu_u$ -uniqueness if  $\forall m \in \mathcal{M}$  and  $\forall \sigma_1, \sigma_2$  s.t.  $\sigma_1 \neq \sigma_2$

$$\mathbb{P}[\text{Vrfy}(vk, (m, \sigma_1)) = \text{Vrfy}(vk, (m, \sigma_2)) = 1 : (vk, sk) \leftarrow \text{KGen}(1^\kappa)] \leq \nu_u,$$

where the probability is taken over the randomness of the verification and key generation algorithms.

Full Domain Hash signatures with trapdoor permutations, for instance RSA-FDH [BR96], are unique. Sometimes unique signatures are also known under the name of *verifiable unpredictable functions* (VUFs).<sup>8</sup> Known constructions of VUFs exist based on strong RSA [MRV99], and on several variants of the Diffie-Hellman assumption in bilinear groups [Lys02, Dod03, DY05, ACF14, Jag15].

### 2.3 Pseudorandom Functions

Let  $F : \{0, 1\}^\kappa \times \mathcal{X} \rightarrow \mathcal{Y}$  be an efficient keyed function, where  $\mathcal{X}$  and  $\mathcal{Y}$  denote the domain and the range of  $F$ . Denote by  $\mathcal{F}$  the set of all functions mapping  $\mathcal{X}$  into  $\mathcal{Y}$ .

**Definition 4** (Pseudorandom function). A function  $F : \{0, 1\}^\kappa \times \mathcal{X} \rightarrow \mathcal{Y}$  is a  $(t, q, \varepsilon)$ -secure pseudorandom function (PRF), if for all adversaries  $D$  running in time at most  $t$  we have

$$\left| \mathbb{P}_{s \leftarrow \mathcal{S}} \left[ \mathbb{D}^{F_s(\cdot)}(1^\kappa) = 1 \right] - \mathbb{P}_{f \leftarrow \mathcal{F}} \left[ \mathbb{D}^{f(\cdot)}(1^\kappa) = 1 \right] \right| \leq \varepsilon,$$

where  $D$  asks at most  $q$  queries to its oracle.

## 3 Subverting Signatures

We proceed to define what it means for an adversary  $B$  to subvert a signature scheme  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$ . We model subversion as the ability of the adversary to replace the genuine signing algorithm with a different algorithm within a certain class  $\mathcal{A}$  of Subversion Attacks (SAs). A subversion of  $\mathcal{SS}$  is an algorithm  $\tilde{A} \in \mathcal{A}$ , specified as follows.

- Algorithm  $\tilde{A}(\cdot, \cdot; \cdot)$  takes as input a signing key  $sk \in \mathcal{SK}$ , a message  $m \in \mathcal{M}$ , random coins  $r \in \mathcal{R}$ , and outputs a subverted signature  $\tilde{\sigma} \in \Sigma$ , where  $\tilde{\sigma} := \tilde{A}(sk, m; r)$ . Notice that algorithm  $\tilde{A}$  is completely arbitrary, with the only restriction that it maintains the same input-output interfaces as the original signing algorithm.

In particular, algorithm  $\tilde{A}$  can hard-wire arbitrary auxiliary information chosen by the adversary, which we denote by a string  $\alpha \in \{0, 1\}^*$ . In general we also allow algorithm  $\tilde{A}$  to be *stateful*, even in case the original signing algorithm is not, and we denote the corresponding state by  $\tau \in \{0, 1\}^*$ ; the state is only used internally by the subverted algorithm and never outputted to the outside.

In Section 3.1 we define what it means for a signature scheme to be secure against a certain class of SAs. In Section 3.2 we define what it means for a class of SAs to be *undetectable* by a user. Some of our definitions are similar in spirit to the ones put forward in [BPR14], except that our modelling of subversion is more general (see below for a more detailed comparison).

<sup>8</sup>Strictly speaking, VUFs satisfy a stronger requirement—namely the uniqueness property holds even for maliciously generated verification keys; the weak variant above is sufficient for the results of this paper.

### 3.1 Impersonation

We consider two security definitions, corresponding to different adversarial goals.

**Indistinguishability.** In the first definition, it is required that an adversary  $\mathsf{B}$  having access to polynomially many subversion oracles chosen adaptively (possibly depending on the user's verification key), cannot distinguish signatures produced via the standard signing algorithm from subverted signatures.

**Definition 5** (Indistinguishability against SAs). Let  $\mathcal{SS} = (\mathsf{KGen}, \mathsf{Sign}, \mathsf{Vrfy})$  be a signature scheme, and  $\mathcal{A}$  be some class of SAs for  $\mathcal{SS}$ . We say that  $\mathcal{SS}$  is  $(t, n, q, \varepsilon)$ -indistinguishable w.r.t. *continuous*  $\mathcal{A}$ -SAs if for all PPT adversaries  $\mathsf{B}$  running in time  $t$ , we have  $|\mathbb{P}[\mathsf{B} \text{ wins}] - \frac{1}{2}| \leq \varepsilon(\kappa)$  in the following game:

1. The challenger runs  $(vk, sk) \leftarrow \mathsf{KGen}(1^\kappa)$ , samples  $b \leftarrow_{\$} \{0, 1\}$ , and gives  $vk$  to  $\mathsf{B}$ .
2. The adversary  $\mathsf{B}$  can ask the following two types of queries; the queries can be specified adaptively and in an arbitrary order:
  - Choose an algorithm  $\tilde{\mathsf{A}}_j \in \mathcal{A}$ , for  $j \in [n]$ , and give it to the challenger.
  - Forward a pair  $(j, m_{i,j})$  to the challenger, where  $i \in [q]$  and  $j \in [n]$ . The answer to each query depends on the value of the secret bit  $b$ . In particular, if  $b = 1$ , the output is  $\sigma_{i,j} \leftarrow \mathsf{Sign}(sk, m_{i,j})$ ; if  $b = 0$ , the output is  $\tilde{\sigma}_{i,j} \leftarrow \tilde{\mathsf{A}}_j(sk, m_{i,j})$ .
3. Finally,  $\mathsf{B}$  outputs a value  $b' \in \{0, 1\}$ ; we say that  $\mathsf{B}$  wins iff  $b' = b$ .

Whenever  $\varepsilon(\kappa) = \mathit{negl}(\kappa)$ ,  $q = \mathit{poly}(\kappa)$ , and  $n = \mathit{poly}(\kappa)$  we simply say that  $\mathcal{SS}$  is indistinguishable against continuous  $\mathcal{A}$ -SAs.

**Impersonation under chosen-message attacks.** We also consider an alternative (strictly weaker—cf. Section 7.3) definition, where the goal of the adversary is now to forge a signature on a “fresh” message (not asked to any of the oracles).

**Definition 6** (EUF-CMA against SAs). Let  $\mathcal{SS} = (\mathsf{KGen}, \mathsf{Sign}, \mathsf{Vrfy})$  be a signature scheme, and  $\mathcal{A}$  be some class of SAs for  $\mathcal{SS}$ . We say that  $\mathcal{SS}$  is  $(t, n, q, \varepsilon)$ -EUF-CMA w.r.t. *continuous*  $\mathcal{A}$ -SAs if for all PPT adversaries  $\mathsf{B}$  running in time  $t$ , we have  $\mathbb{P}[\mathsf{B} \text{ wins}] \leq \varepsilon(\kappa)$  in the following game:

1. The challenger runs  $(vk, sk) \leftarrow \mathsf{KGen}(1^\kappa)$ , and gives  $vk$  to  $\mathsf{B}$ .
2. The adversary  $\mathsf{B}$  is given oracle access to  $\mathsf{Sign}(sk, \cdot)$ . Upon input the  $i$ -th query  $m_i$ , this oracle returns  $\sigma_i \leftarrow \mathsf{Sign}(sk, m_i)$ ; let  $\mathcal{Q} = \{m_1, \dots, m_q\}$  be the set of all queried messages.
3. For each  $j \in [n]$ , the adversary  $\mathsf{B}$  can adaptively choose an algorithm  $\tilde{\mathsf{A}}_j \in \mathcal{A}$ . For each algorithm,  $\mathsf{B}$  is given oracle access to  $\tilde{\mathsf{A}}_j(sk, \cdot)$ . Upon input a message  $\tilde{m}_{i,j}$ , the oracle returns  $\tilde{\sigma}_{i,j} \leftarrow \tilde{\mathsf{A}}_j(sk, \tilde{m}_{i,j})$ ; let  $\tilde{\mathcal{Q}}_j = \{\tilde{m}_{1,j}, \dots, \tilde{m}_{q,j}\}$  be the set of all queried messages to the oracle  $\tilde{\mathsf{A}}_j$ .
4. Finally,  $\mathsf{B}$  outputs a pair  $(m^*, \sigma^*)$ ; we say that  $\mathsf{B}$  wins iff  $\mathsf{Vrfy}(vk, (m^*, \sigma^*)) = 1$  and  $m^* \notin \mathcal{Q} \cup \tilde{\mathcal{Q}}$ , where  $\tilde{\mathcal{Q}} := \bigcup_{j=1}^n \tilde{\mathcal{Q}}_j$ .

Whenever  $\varepsilon(\kappa) = \mathit{negl}(\kappa)$ ,  $q = \mathit{poly}(\kappa)$ , and  $n = \mathit{poly}(\kappa)$  we simply say that  $\mathcal{SS}$  is EUF-CMA against continuous  $\mathcal{A}$ -SAs.

**Remarks.** Some remarks on the above definitions are in order.

- First, note that it is impossible to prove that a signature scheme  $\mathcal{SS}$  satisfies Definition 5 (and consequently Definition 6) for an *arbitrary* class  $\mathcal{A}$ , without making further assumptions.<sup>9</sup> To see this, consider the simple algorithm that ignores all inputs and outputs the secret key.<sup>10</sup>
- We observe that continuous  $\mathcal{A}$ -SAs security, implies security against continuous tampering attacks with the secret key. This can be seen by considering a class of algorithms  $\mathcal{A}_{\text{key}} = \{\tilde{A}_f\}_{f \in \mathcal{F}}$ , where  $\mathcal{F}$  is a class of functions such that each  $f \in \mathcal{F}$  has a type  $f : \mathcal{SK} \rightarrow \mathcal{SK}$ , and for all  $f \in \mathcal{F}$ ,  $m \in \mathcal{M}$  and  $r \in \mathcal{R}$  we have that  $\tilde{A}_f(\cdot, m; r) := \text{Sign}(f(\cdot), m; r)$ .<sup>11</sup>
- It is useful to compare Definition 5 to the security definition against algorithm-substitution attacks given in [BPR14] (for the case of symmetric encryption). In the language of Bellare *et al.* [BPR14], a subversion of a signature scheme would be a triple of algorithms  $\widetilde{\mathcal{SS}} = (\widetilde{\text{KGen}}, \widetilde{\text{Sign}}, \widetilde{\text{Vrfy}})$ , where in the security game  $\widetilde{\text{KGen}}$  is run by the challenger in order to obtain a trapdoor  $\alpha \in \{0, 1\}^*$  and some initial state  $\tau \in \{0, 1\}^*$  which are both hard-wired in the algorithm  $\widetilde{\text{Sign}} := \text{Sign}_{\alpha, \tau}$  (and given to B).<sup>12</sup>

The above setting can be cast in our framework by considering the class of SAs  $\mathcal{A}_{\text{BPR14}} := \{\tilde{A}_{\alpha, \tau} : (\alpha, \tau) \leftarrow \widetilde{\text{KGen}}(1^\kappa)\}$ , and by setting  $n = 1$  in Definition 5. Our definition is more general, as it accounts for arbitrary classes of SAs and moreover allows B to subvert a user’s algorithm continuously and in a fully-adaptive fashion (possibly depending on the target verification key).

**Multi-user setting.** For simplicity Definition 5 and 6 consider a single user. We provide an extension to the more general setting with  $u \geq 2$  users, together with a complete picture of the relationships between different notions, in Section 7.

### 3.2 Public/Secret Undetectability

By undetectability, we mean the inability of ordinary users to tell whether signatures are computed using the subverted or the genuine signing algorithm. We will distinguish between the case where a subversion is *publicly* or *secretly* undetectable. Roughly speaking, public undetectability means that no user can detect subversions using the verification key  $vk$  only (i.e., without knowing the signing key  $sk$ ); secret undetectability means that no user, even with knowledge of the signing key  $sk$ , can detect subversions.

A formal definition follows. While reading it, bear in mind that the challenger plays the role of the “bad guy” trying to sabotage the signature scheme without being detected.

**Definition 7** (Public/Secret Undetectability). Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a signature scheme, and  $\mathcal{A}$  be some class of SAs for  $\mathcal{SS}$ . We say that  $\mathcal{A}$  is *secretly*  $(t, q, \varepsilon)$ -undetectable

<sup>9</sup>Looking ahead, one of our positive results achieves security w.r.t. arbitrary SAs assuming the existence of a cryptographic reverse firewall. See Section 6.

<sup>10</sup>In case the secret key is too long, one can make the algorithm stateful so that it outputs a different chunk of the key at each invocation. Alternatively, consider the class of algorithms  $\{\tilde{A}_{\tilde{m}}\}_{\tilde{m} \in \mathcal{M}}$  that always outputs a signature  $\tilde{\sigma}$  on  $\tilde{m}$ ; obviously this subversion allows to forge on  $\tilde{m}$  without explicitly querying the message to any of the oracles.

<sup>11</sup>It is worth noting that already for  $n = 1$  Definition 6 implies *non-adaptive* key tampering, as the subverted algorithm can hard-wire (the description of) polynomially many pre-set tampering functions.

<sup>12</sup>The algorithm  $\widetilde{\text{Vrfy}}$  is not explicitly part of the definitions in [BPR14]—in fact, a secure scheme implicitly excludes that any  $\widetilde{\text{Vrfy}}$  algorithm exists—and can be considered as part of the adversary itself.

w.r.t.  $\mathcal{SS}$  if for all PPT users  $U$  running in time  $t$ , there exists an efficient challenger such that  $|\mathbb{P}[U \text{ wins}] - \frac{1}{2}| \leq \varepsilon(\kappa)$  in the following game:

1. The challenger runs  $(vk, sk) \leftarrow \text{KGen}(1^\kappa)$ , chooses an algorithm  $\tilde{A} \in \mathcal{A}$  (possibly depending on  $vk$ ), samples  $b \leftarrow_s \{0, 1\}$  and gives  $(vk, sk)$  to  $U$ .
2. The user  $U$  can ask queries  $m_i \in \mathcal{M}$ , for all  $i \in [q]$ . The answer to each query depends on the secret bit  $b$ . In particular, if  $b = 1$ , the challenger returns  $\sigma_i \leftarrow \text{Sign}(sk, m_i)$ ; if  $b = 0$ , the challenger returns  $\tilde{\sigma}_i \leftarrow \tilde{A}(sk, m_i)$ .
3. Finally,  $U$  outputs a value  $b' \in \{0, 1\}$ ; we say that  $U$  wins iff  $b' = b$ .

We say that  $\mathcal{A}$  is *publicly* undetectable w.r.t.  $\mathcal{SS}$  if in step 1. of the above game,  $U$  is only given the verification key. Moreover, whenever  $\varepsilon(\kappa) = \text{negl}(\kappa)$  and  $q = \text{poly}(\kappa)$  we simply say that  $\mathcal{A}$  is secretly/publicly undetectable w.r.t.  $\mathcal{SS}$ .

Our definition of undetectability is similar to the corresponding definition considered by Bellare *et al.* [BPR14] for the case of symmetric encryption. One key difference is that, in the definition above, the challenger is allowed to choose the subversion algorithm possibly depending on the verification key of the user.<sup>13</sup> While one could in principle define even stronger forms of undetectability, e.g. by requiring that continuous and fully-adaptive SAs remain undetectable, we do not pursue this direction here. The reason for this is that the attacks we analyze in Section 4 are non-adaptive and only require to use a single subversion.

**Secret vs. public undetectability.** While secret undetectability clearly implies public undetectability, the converse is not true. In particular, in Section 7.4 we show that there exists a signature scheme  $\mathcal{SS}$  and a set of subversions  $\mathcal{A}$  of it such that  $\mathcal{A}$  is publicly undetectable w.r.t.  $\mathcal{SS}$  but it is secretly detectable w.r.t.  $\mathcal{SS}$ .

**Multi-user setting.** For simplicity Definition 7 considers a single user. We provide an extension to the more general setting with  $u \geq 2$  users, together with a complete picture of the relationships between different notions, in Section 7.

## 4 Mounting Subversion Attacks

In Section 4.1 we show that the biased-randomness attack of [BPR14] (adapted to the case of signatures), satisfies secret undetectability as per Definition 7 while allowing to recover the user's signing key with overwhelming probability. This attack allows to break all signature schemes using a sufficient amount of randomness; in Section 4.2 we present a new attack allowing to surreptitiously subvert even signature schemes using only little randomness (say 1 bit), provided that the targeted scheme satisfies an additional property.

### 4.1 Attacking Coin-Injective Schemes

We start by recalling an information-theoretic lemma from [BPR14]. Suppose  $g : \mathcal{R} \rightarrow \mathcal{R}'$  where  $\mathcal{R}, \mathcal{R}' \subseteq \{0, 1\}^*$ ,  $f : \{0, 1\}^* \rightarrow \{0, 1\}$ , and  $\rho = |\mathcal{R}|$ . For  $b \in \{0, 1\}$  consider the following *biased* distribution:

$$\tilde{\mathcal{R}}^{f,g}(b, \mathcal{R}) = \{r \in \mathcal{R} : f(g(r)) = b\}. \quad (1)$$

The lemma below roughly says that if a value  $r$  is chosen at random from the real distribution  $\mathcal{R}$ , the probability that  $r$  is also in the biased distribution  $\tilde{\mathcal{R}}$  is high if  $|\mathcal{R}|$  is large enough.

<sup>13</sup>Looking ahead, our new attack (cf. Section 4.2) will rely on this feature in the multi-user setting.

**SA class  $\mathcal{A}_{\text{bias}}^F$**

Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a randomized signature scheme with randomness space  $\mathcal{R}$ , and  $F : \{0, 1\}^\kappa \times \{0, 1\}^* \rightarrow \{0, 1\}$  be a pseudorandom function. The class  $\mathcal{A}_{\text{bias}}^F$  consists of a set of algorithms  $\{\tilde{\mathbf{A}}_{s,\tau}\}_{s \in \{0,1\}^\kappa, \tau=1}$ , where each algorithm in the class behaves as follows:

$\tilde{\mathbf{A}}_{s,\tau}(sk, m)$ :

- For  $|sk| = \ell$ , let  $i := \tau \bmod \ell$ .
- Define the function  $g(\cdot) := \text{Sign}(sk, m; \cdot) \parallel \tau$  and sample a random element  $\tilde{r}$  from the distribution

$$\tilde{\mathcal{R}}^{F(s,\cdot),g(\cdot)}(sk[i], \mathcal{R}) := \{r \in \mathcal{R} : F(s, g(r)) = sk[i]\}. \quad (2)$$

- Return the signature  $\sigma := \text{Sign}(sk, m; \tilde{r})$ , and update the state  $\tau \leftarrow \tau + 1$ .

**Extracting the signing key.** Given as input a vector of signatures  $\vec{\sigma} = (\sigma_1, \dots, \sigma_\ell)$ , for each signature  $\sigma_i \in \vec{\sigma}$  try to extract the  $i$ -th bit of the signing key by defining  $sk'[i] := F(s, \sigma_i \parallel i)$ . Return the signing key  $sk' := (sk'[1], \dots, sk'[\ell])$ .

**Figure 1:** Attacking coin-injective schemes

**Lemma 2** (Lemma 1 of [BPR14]). *Let  $f, g, b, \mathcal{R}$ , and  $\tilde{\mathcal{R}} = \tilde{\mathcal{R}}^{f,g}(b, \mathcal{R})$  be as defined above. Then, if  $g$  is injective and  $f$  is drawn at random, for all  $r \in \mathcal{R}$  we have*

$$\mathbb{P}_{\tilde{r} \leftarrow \tilde{\mathcal{R}}} [r = \tilde{r}] = (1 - 2^{-\rho})/\rho.$$

The following attack is based on the biased-randomness attack from [BPR14]. Roughly, what it does is to embed a trapdoor—a key for a pseudorandom function—in the subverted signing algorithm and to “bias” the randomness in a way that it becomes possible to any party that knows the trapdoor to leak one bit of the signing key for each signed message under that signing key. Hence, if the adversary can obtain at least  $|sk|$  signed messages then it can later extract the entire signing key in full.

For the analysis, which relies on Lemma 2, we will need to assume the signing function is injective w.r.t. its random coins—a notion which we define below.

**Definition 8** (Coin-injective). We say that  $\mathcal{SS}$  is *coin-injective* if for all  $m \in \mathcal{M}$ , and for all  $(vk, sk) \leftarrow \text{KGen}(1^\kappa)$ , we have that  $\text{Sign}(sk, m; \cdot)$  is injective.

**Theorem 1.** *Let  $F : \{0, 1\}^\kappa \times \{0, 1\}^* \rightarrow \{0, 1\}$  be a  $(t_{\text{prf}}, q_{\text{prf}}, \varepsilon_{\text{prf}})$ -secure PRF. For a randomized, coin-injective signature scheme  $\mathcal{SS}$  with randomness space of size  $\rho = |\mathcal{R}|$ , consider the class of SAs  $\mathcal{A}_{\text{bias}}^F$  described in Fig. 1. Then,*

- (i)  $\mathcal{A}_{\text{bias}}^F$  is secretly  $(t, q, \varepsilon)$ -undetectable for  $t \approx t_{\text{prf}}$ ,  $q \approx q_{\text{prf}}$  and  $\varepsilon \leq q \cdot 2^{-(\rho+1)} + \varepsilon_{\text{prf}}$ .
- (ii) Each  $\tilde{\mathbf{A}} \in \mathcal{A}_{\text{bias}}^F$  recovers the signing key of the user with probability at least  $(1 - (1/2 + \varepsilon_{\text{prf}})^\rho)^\ell$  where  $\ell$  is the size of the signing key.

*Proof.* (i) Let  $\mathbf{G}$  be the game described in Definition 7, where the challenger picks  $\tilde{\mathbf{A}} \leftarrow \mathcal{A}_{\text{bias}}^F$  (independently of the user’s verification key). Consider the game  $\mathbf{G}_0$ , an identical copy of game  $\mathbf{G}$  when  $b = 0$ , and consider the game  $\mathbf{G}_1$ , an identical copy of game  $\mathbf{G}$  when  $b = 1$ . For the first part of the proof the objective is to show that  $\mathbf{G}_0 \approx \mathbf{G}_1$ .

Now consider game  $\mathbf{G}'_0$  an identical copy of game  $\mathbf{G}_0$  except that  $\mathbf{G}'_0$  utilizes the distribution from Eq. (1) instead of the distribution from Eq. (2).

**Claim 1.**  $|\mathbb{P}[\text{U wins in } \mathbf{G}_0] - \mathbb{P}[\text{U wins in } \mathbf{G}'_0]| \leq \varepsilon_{\text{prf}}$ .

*Proof.* We assume that there exists a user  $\text{U}$  that distinguishes between games  $\mathbf{G}_0$  and  $\mathbf{G}'_0$ , and we build a distinguisher  $\text{D}$  (using  $\text{U}$ ) that breaks the pseudorandomness of the PRF  $F$ . Distinguisher  $\text{D}$  is described below below.

Distinguisher  $\text{D}$ :

- Run  $(vk, sk) \leftarrow \text{KGen}(1^\kappa)$ , and return  $(vk, sk)$  to  $\text{U}$ .
- For each query  $m_i \in \mathcal{M}$  asked by  $\text{U}$ , do:
  1. Pick a random  $r \leftarrow_s \mathcal{R}$  and compute  $x_i = \text{Sign}(sk, m_i; r) \parallel \tau$ .
  2. Forward  $x_i$  to the target oracle, which answers with  $y_i = f(x_i)$  if  $b = 0$  or with  $y_i = F(s, x_i)$  if  $b = 1$  (for a hidden bit  $b$ ).
  3. If  $y_i = sk[i]$ , then forward  $\sigma_i = \text{Sign}(sk, m_i; r)$  as an answer to the query of  $\text{U}$ , otherwise return to step (1).<sup>14</sup>
- Output whatever  $\text{U}$  outputs.

Notice that the probability that  $\text{D}$  aborts in step (3) of the reduction is the same probability that in game  $\mathbf{G}_0$  and  $\mathbf{G}'_0$  the subverted signing algorithm fails to sample from the set  $\tilde{\mathcal{R}}$ . It follows that in case  $b = 0$  distinguisher  $\text{D}$  perfectly emulates the distribution of  $\mathbf{G}_0$ , whereas in case  $b = 1$  it perfectly emulates the distribution of  $\mathbf{G}'_0$ . The claim follows.  $\square$

**Claim 2.**  $|\mathbb{P}[\text{U wins in } \mathbf{G}'_0] - \mathbb{P}[\text{U wins in } \mathbf{G}_1]| \leq q \cdot 2^{-(\rho+1)}$ .

*Proof.* Abusing notation, let us write  $\mathbf{G}'_0$  and  $\mathbf{G}_1$  for the distribution of the random variables corresponding to  $\text{U}$ 's view in games  $\mathbf{G}'_0$  and  $\mathbf{G}_1$  respectively. For an index  $i \in [0, q]$  consider the hybrid game  $\mathbf{H}_i$  that answers the first  $i$  signature queries as in game  $\mathbf{G}'_0$  while all the subsequent queries are answered as in  $\mathbf{G}_1$ . We note that  $\mathbf{H}_0 = \mathbf{G}_1$  and  $\mathbf{H}_q = \mathbf{G}'_0$ .

We claim that for all  $i \in [q]$ , we have  $\text{SD}(\mathbf{H}_{i-1}, \mathbf{H}_i) \leq 2^{-(\rho+1)}$ . To see this, fix some  $i \in [q]$  and denote with  $\mathbf{R}$  (resp.  $\tilde{\mathbf{R}}$ ) the random variable defined by sampling an element from  $\mathcal{R}$  (resp.  $\tilde{\mathcal{R}}$ ) uniformly at random. Clearly,

$$\begin{aligned} \text{SD}(\mathbf{H}_{i-1}, \mathbf{H}_i) &\leq \text{SD}(\mathbf{R}, \tilde{\mathbf{R}}) = \frac{1}{2} \cdot \sum_{r \in \mathcal{R}} \left| \mathbb{P}[\mathbf{R} = r] - \mathbb{P}[\tilde{\mathbf{R}} = r] \right| \\ &= \frac{1}{2} \cdot \sum_{r \in \mathcal{R}} \left| \frac{1}{\rho} - \frac{1 - 2^{-\rho}}{\rho} \right| \\ &= \frac{1}{2} \cdot 2^{-\rho} = 2^{-(\rho+1)}, \end{aligned} \tag{3}$$

where Eq. (3) follows by Lemma 2.

The claim now follows by the triangle inequality, as

$$\text{SD}(\mathbf{G}_1, \mathbf{G}'_0) \leq \sum_{i=1}^q \text{SD}(\mathbf{H}_{i-1}, \mathbf{H}_i) \leq q \cdot 2^{-(\rho+1)}.$$

$\square$

---

<sup>14</sup>In case  $|\mathcal{R}|$  is exponential  $\text{D}$  simply aborts after polynomially many trials.

The two claims above finish the proof of statement (i).

(ii) For the second part of the proof we show that the attack of Fig. 1 fails to recover the secret key with probability at most  $e_1 + e_2 + \dots + e_\ell$ , where  $e_j := \mathbb{P}[sk'[j] \neq sk[j]]$ . In the analysis, we replace for simplicity the function  $F$  with a truly random function  $f$ ; a generalization accounting for the negligible error due to the use of a pseudorandom function is straightforward. Note that all applications of  $f$  are independent because we append the value  $\tau$  to each query.

Now if  $g$  is injective and  $f$  is a random function that outputs one bit, then for each element  $r \in \mathcal{R}$  we have  $\mathbb{P}[f(g(r)) = sk[j]] = 1/2$ . Extending to the entire set  $\mathcal{R}$  of size  $\rho$  we have that

$$e_j := \mathbb{P}\left[\tilde{\mathcal{R}}^{f \circ g}(sk[j], \mathcal{R}) = \emptyset\right] = 2^{-\rho},$$

is the error probability for each bit of the secret key. Therefore the probability of recovering the key is at least  $(1 - 2^{-\rho})^\ell$ .  $\square$

Notice that for the attack to be undetectable with high probability, the underlying signature scheme needs to rely on a *minimal* amount of randomness, say  $\rho \geq 2^7$ .

**Making the attack stateless** Note that the attack of Fig. 1 requires the subverted signature algorithm to maintain a state of logarithmic size (the counter  $\tau$ ). At first sight this might seem a strong assumption, since the original signing algorithm is typically stateless.

However, if we assume that the message space is polynomial and that the adversary can control the input messages, one can easily adapt the attack to be completely stateless by letting the input message play the role of the counter  $\tau$ . Namely, algorithm  $\tilde{\mathbf{A}}$  interprets the input message  $m$  as an integer  $i \in [\ell]$  and proceeds as before. The above adaptation still allows to recover the signing key and it is undetectable under the *strong undetectability* definition put forward by [BJK15].

## 4.2 Attacking Coin-Extractable Schemes

The attack on Section 4.1 allows to break all sufficiently randomized schemes. This leaves the interesting possibility to show a positive result for schemes using less randomness, e.g., the Katz-Wang signature scheme [KW03] that uses a single bit of randomness. In this section we present a simple attack (cf. Fig. 2) ruling out the above possibility for all signature schemes that are *coin-extractable*, a notion which we define next.

**Definition 9** (Coin-extractable). Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a signature scheme. We say that  $\mathcal{SS}$  is  $\nu_{ext}$ -coin-extractable if there exists a PPT algorithm CExt such that for all  $m \in \mathcal{M}$

$$\mathbb{P}[\sigma = \text{Sign}(sk, m; r) : (vk, sk) \leftarrow \text{KGen}(1^\kappa); \sigma = \text{Sign}(sk, m); r \leftarrow \text{CExt}(vk, m, \sigma)] \geq 1 - \nu_{ext}.$$

We point that many existing signature schemes are coin-extractable:

- All *public-coin* signature schemes [Sch12], where the random coins used to generate a signature are included as part of the signature. Concretely, the schemes in [GHR99, CS00, NPS01, CL02, Fis03, CL04, BB08, HW09a, HW09b, HK12], and the Unstructured Rabin-Williams scheme [Ber08], are all public-coin.
- The Katz-Wang scheme [KW03], where the signature on a message  $m$  is computed as  $\sigma = f^{-1}(H(m||r))$  such that  $f$  is a trapdoor permutation,  $H$  is a hash function, and  $r$  is random bit. Given a pair  $(m, \sigma)$  the extractor simply sets  $r = 1$  iff  $f(\sigma) = H(m||1)$ .
- The PSS signature scheme [BR96, Cor02].

**SA class  $\mathcal{A}_{\text{cext}}$**

Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a coin-extractable, randomized signature scheme with randomness space  $\mathcal{R}$  of size  $\rho = 2^d$ . For simplicity assume that  $d|\ell$ , where  $\ell$  is the size of the signing key (a generalization is straightforward). The class  $\mathcal{A}_{\text{cext}}$  consists of a set of algorithms  $\{\tilde{\mathbf{A}}_{s,\tau}\}_{s \in \{0,1\}^\ell, \tau=0}$ , where each algorithm in the class behaves as follows:

$\tilde{\mathbf{A}}_{s,\tau}(sk, m)$ :

- If  $\tau \geq \ell$  output a honestly generated signature  $\sigma := \text{Sign}(sk, m; r)$ .
- Else,
  - for each value  $j \in [d]$  compute the biased random bit  $\tilde{r}[j] := s[\tau + j] \oplus sk[\tau + j]$ ;
  - return the signature  $\sigma := \text{Sign}(sk, m; \tilde{r})$ , and update the state  $\tau \leftarrow \tau + d$ .

**Extracting the signing key.** Given as input a vector of signatures  $\vec{\sigma} = (\sigma_1, \dots, \sigma_{\ell/d})$ , parse the trapdoor  $s$  as  $\ell/d$  chunks of  $d$  bits  $s = \{s_1, \dots, s_{\ell/d}\}$ . For each signature  $\sigma_i \in \vec{\sigma}$  try to extract the  $d$ -bit chunk  $sk'_i$  of the signing key as follows.

- Extract the randomness from the  $i$ -th signature  $\tilde{r} \leftarrow \text{CExt}(vk, m_i, \sigma_i)$ .
- For each value  $j \in [d]$  compute the secret key bit  $sk'_i[j] := \tilde{r}[j] \oplus s_i[j]$ .

Return the signing key  $sk' := (sk'_1, \dots, sk'_{\ell/d})$ .

**Figure 2:** Attacking coin-extractable schemes

**Theorem 2.** For a randomized,  $\nu_{\text{ext}}$ -coin-extractable, signature scheme  $\mathcal{SS}$  with randomness space  $\mathcal{R}$  of size  $\rho = 2^d$ , consider the class of SAs  $\mathcal{A}_{\text{cext}}$  described in Fig. 2. Then,

- (i)  $\mathcal{A}_{\text{cext}}$  is secretly  $(t, q, 0)$ -undetectable for  $t, q \in \mathbb{N}$ .
- (ii) Each  $\tilde{\mathbf{A}} \in \mathcal{A}_{\text{cext}}$  recovers the signing key of the user with probability at least  $(1 - \nu_{\text{ext}})^{\ell/d}$ , where  $\ell$  is the size of the key.

*Proof.* (i) Let  $\mathbf{G}$  be the game described in Definition 7, where the challenger picks  $\tilde{\mathbf{A}} \leftarrow_{\$} \mathcal{A}_{\text{cext}}$  uniformly at random (and independently of the user's verification key). Consider the game  $\mathbf{G}_0$ , an identical copy of game  $\mathbf{G}$  when  $b = 0$ , and consider the game  $\mathbf{G}_1$ , an identical copy of game  $\mathbf{G}$  when  $b = 1$ . For the first part of the proof the objective is to show that  $\mathbf{G}_0 \approx \mathbf{G}_1$ .

**Claim 3.**  $|\mathbb{P}[\text{U wins in } \mathbf{G}_0] - \mathbb{P}[\text{U wins in } \mathbf{G}_1]| = 0$ .

*Proof.* Abusing notation, let us write  $\mathbf{G}_0$  and  $\mathbf{G}_1$  for the distribution of the random variables corresponding to U's view in games  $\mathbf{G}_0$  and  $\mathbf{G}_1$  respectively. For an index  $i \in [0, q]$  consider the hybrid game  $\mathbf{H}_i$  that answers the first  $i$  signature queries as in game  $\mathbf{G}_0$  while all the subsequent queries are answered as in  $\mathbf{G}_1$ . We note that  $\mathbf{H}_0 \equiv \mathbf{G}_1$  and  $\mathbf{H}_q \equiv \mathbf{G}_0$ .

We claim that for all  $i \in [q]$ , we have  $\mathbf{H}_{i-1} \equiv \mathbf{H}_i$ . To see this, fix some  $i \in [q]$  and denote with  $\mathbf{R}$  (resp.  $\tilde{\mathbf{R}}$ ) the random variable defined by sampling an element from  $\mathcal{R}$  (resp.  $\tilde{\mathcal{R}}$ ) uniformly at random. It is easy to see that  $\mathbf{R}$  and  $\tilde{\mathbf{R}}$  are identically distributed, as the biased distribution consists of a one-time pad encryption of (part of) the signing key with a uniform key. The claim follows.  $\square$



(ii) For the second part of the proof we note that the attack of Fig. 2 successfully recovers the biased randomness  $\tilde{r}$  of each  $\sigma_i \in \{\sigma_1, \dots, \sigma_{\ell/d}\}$  and computes the chunk  $sk_i$  of the signing key with probability at least  $1 - \nu_{ext}$ . This gives a total probability of recovering the entire signing key of at least  $(1 - \nu_{ext})^{\ell/d}$ .  $\square$

**Making the attack stateless.** Note that the attack of Fig. 2 requires the subverted signature algorithm to maintain a state of logarithmic size (the counter  $\tau$ ). Similarly to the case of the attack of Fig. 1 this assumption can be removed by letting the input message play the role of the counter. (See also the discussion at the end of Section 4.1.)

## 5 Security of Unique Signatures

In this section we prove that signature schemes with unique signatures are subversion-resilient against SAs that meet a basic undetectability requirement, which we call the verifiability condition.

### 5.1 The Verifiability Condition

We say that  $\mathcal{A}$  meets the *verifiability condition* relative to  $\mathcal{SS}$  if for all  $\tilde{\mathbf{A}} \in \mathcal{A}$  and for all  $m \in \mathcal{M}$  the signatures produced using the subverted signing algorithm  $\tilde{\mathbf{A}}$  (almost) always verify under the corresponding verification key  $vk$ .

**Definition 10** (Verifiability). Let  $\mathcal{A}$  be some class of SAs for a signature scheme  $\mathcal{SS}$ . We say that  $\mathcal{A}$  satisfies  $\nu_v$ -verifiability if for all  $\tilde{\mathbf{A}} \in \mathcal{A}$  and for all  $m \in \mathcal{M}$

$$\mathbb{P} \left[ \text{Vrfy}(vk, (m, \tilde{\mathbf{A}}(sk, m))) = 1 : (vk, sk) \leftarrow \text{KGen}(1^\kappa) \right] \geq 1 - \nu_v,$$

where the probability is taken over the randomness of all involved algorithms.

**Public undetectability vs. verifiability.** One might think that verifiability is a special case of public undetectability. However, this is not true and in fact Definition 10 and 7 are incomparable. To see this, consider the class of SAs  $\mathcal{A}_{\text{msg}} = \{\tilde{\mathbf{A}}_{\tilde{m}}\}_{\tilde{m} \in \mathcal{M}}$  that behaves identically to the original signing algorithm, except that upon input  $\tilde{m} \in \mathcal{M}$  it outputs an invalid signature.<sup>15</sup> Clearly,  $\mathcal{A}_{\text{msg}}$  satisfies public undetectability as a user has only a negligible chance of hitting the value  $\tilde{m}$ ; yet  $\mathcal{A}_{\text{msg}}$  does not meet the verifiability condition as the latter is a property that holds for *all* messages.

On the other hand, consider the class of SAs  $\mathcal{A}_{\text{det}}$  that is identical to the original signing algorithm, except that it behaves deterministically on repeated inputs. Clearly,  $\mathcal{A}_{\text{det}}$  meets the verifiability condition relative to any (even randomized) signature scheme  $\mathcal{SS}$ ; yet  $\mathcal{A}_{\text{det}}$  does not satisfy public undetectability for any randomized signature scheme  $\mathcal{SS}$ , as a user can simply query the same message twice in order to guess the value of the hidden bit  $b$  with overwhelming probability.

**Relaxed verifiability.** Clearly, the assumption that the verifiability condition should hold for all values  $m \in \mathcal{M}$  is quite a strong one. A natural relaxation is to require that the probability in Definition 10 is taken also over the choice of the message.

<sup>15</sup>A similar class of attacks—under the name of input-triggered subversion—has been recently considered in [DFP15] for the case of symmetric encryption.

**Definition 11** (Relaxed Verifiability). Let  $\mathcal{A}$  be some class of SAs for a signature scheme  $\mathcal{SS}$ . We say that  $\mathcal{A}$  satisfies *relaxed  $\nu_v$ -verifiability* if for all  $\tilde{\mathbf{A}} \in \mathcal{A}$

$$\mathbb{P} \left[ \text{Vrfy}(vk, (m, \tilde{\mathbf{A}}(sk, m))) = 1 : (vk, sk) \leftarrow \text{KGen}(1^\kappa); m \leftarrow_s \mathcal{M} \right] \geq 1 - \nu_v,$$

where the probability is taken over the choice of the message and over the randomness of all involved algorithms.

We argue that relaxed verifiability is implied by public undetectability (cf. Definition 7) in many interesting cases.

- *Input-triggered subversions.* Whenever public undetectability holds for all algorithms in the class  $\mathcal{A}$ . This is the case, for instance, for the class  $\mathcal{A}_{\text{msg}}$  of input-triggered subversions described above.

To see this, let  $\mathcal{A}$  be a class of SAs that is publicly undetectable for all  $\tilde{\mathbf{A}} \in \mathcal{A}$ . Towards a contradiction, assume that  $\mathcal{A}$  does not satisfy relaxed verifiability. This means that there exists an  $\tilde{\mathbf{A}} \in \mathcal{A}$  and a polynomial  $p(\cdot)$  such that

$$\mathbb{P} \left[ \text{Vrfy}(vk, (m, \tilde{\mathbf{A}}(sk, m))) = 0 : (vk, sk) \leftarrow \text{KGen}(1^\kappa); m \leftarrow_s \mathcal{M} \right] \geq \frac{1}{p(\kappa)},$$

for infinitely many values of  $\kappa \in \mathbb{N}$ . It follows that  $\tilde{\mathbf{A}}$  can be used to break public undetectability with probability  $1/p(\kappa)$ , by simply signing a random message and trying to verify the outcome.

- *Backdoored implementations.* The above implication also holds for the class  $\mathcal{A}_{\text{BPR14}}$  of algorithm-substitution attacks and backdoored implementations (see paragraph “Remarks” in Section 3.1), as long as the winning condition in Definition 7 and Definition 11 is taken also over the choice of the backdoor (i.e., over the random coins of algorithm  $\widetilde{\text{KGen}}$ ).

## 5.2 Chosen-Message Attacks

The theorem below shows that unique signature schemes (cf. Definition 3) achieve indistinguishability (and thus EUF-CMA) against the class of all SAs that meet the verifiability condition (cf. Definition 10).

**Theorem 3.** *Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a signature scheme with  $\nu_c$ -correctness and  $\nu_u$ -uniqueness, and denote by  $\mathcal{A}_{\text{ver}}^{\nu_v}$  the class of all algorithms that satisfy  $\nu_v$ -verifiability relative to  $\mathcal{SS}$ . Then  $\mathcal{SS}$  is  $(t, n, q, \varepsilon)$ -indistinguishable against continuous  $\mathcal{A}_{\text{ver}}^{\nu_v}$ -SAs, for all  $n, q \in \mathbb{N}$  and for  $\varepsilon \leq qn \cdot (\nu_c + \nu_v + \nu_u)$ .*

*Proof.* Let  $\mathbf{G}$  be the game described in Definition 5. Consider the game  $\mathbf{G}_0$ , an identical copy of game  $\mathbf{G}$  when  $b = 0$ , and consider the game  $\mathbf{G}_1$ , an identical copy of game  $\mathbf{G}$  when  $b = 1$ . The objective here is to show that  $\mathbf{G}_0 \approx \mathbf{G}_1$ .

For an index  $k \in [0, n]$ , consider the hybrid game  $\mathbf{H}_k$  that answers each query  $(j, m_{i,j})$  such that  $j \leq k$  as in game  $\mathbf{G}_0$  (i.e., by running  $\text{Sign}(sk, m_{i,j})$ ), while all queries  $(j, m_{i,j})$  such that  $j > k$  are answered as in  $\mathbf{G}_1$  (i.e., by running  $\tilde{\mathbf{A}}_j(sk, m_{i,j})$ ). We note that  $\mathbf{H}_0 \equiv \mathbf{G}_1$  and  $\mathbf{H}_n \equiv \mathbf{G}_0$ . Abusing notation, let us write  $\mathbf{G}_k$  for the distribution of the random variable corresponding to  $\mathbf{B}$ 's view in games  $\mathbf{G}_k$ .

Fix a particular  $k \in [0, n]$ , and for an index  $l \in [0, q]$  consider the hybrid game  $\mathbf{H}_{k,l}$  that is identical to  $\mathbf{H}_k$  except that queries  $(k, m_{i,k})$  with  $i \leq l$  are treated as in game  $\mathbf{G}_0$ , while queries  $(k, m_{i,k})$  with  $i > l$  are treated as in  $\mathbf{G}_1$ . Observe that  $\mathbf{H}_{k,0} \equiv \mathbf{H}_{k-1}$ , and  $\mathbf{H}_{k,q} \equiv \mathbf{H}_k$ .

**Claim 4.** Fix some  $k \in [0, n]$ . For each  $l \in [0, q]$ , we have  $\mathbb{SD}(\mathbf{H}_{k,l-1}, \mathbf{H}_{k,l}) \leq \nu_c + \nu_v + \nu_u$ .

*Proof.* Notice that the only difference between  $\mathbf{H}_{k,l-1}$  and  $\mathbf{H}_{k,l}$  is how the two games answer the query  $(k, m_{l,k})$ : Game  $\mathbf{H}_{k,l-1}$  returns  $\sigma_{l,k} \leftarrow \text{Sign}(sk, m_{l,k})$ , whereas game  $\mathbf{H}_{k,l}$  returns  $\tilde{\sigma}_{l,k} \leftarrow \tilde{\mathbf{A}}_k(sk, m_{l,k})$ . Now let  $E_{l,k}$  be the event that  $\sigma_{l,k} \neq \tilde{\sigma}_{l,k}$ . We can write

$$\mathbb{SD}(\mathbf{H}_{k,l-1}, \mathbf{H}_{k,l}) \leq \mathbb{SD}(\mathbf{H}_{k,l-1}; \mathbf{H}_{k,l} | \neg E_{l,k}) + \mathbb{P}[E_{l,k}] \quad (4)$$

$$\leq \nu_c + \nu_u + \nu_v. \quad (5)$$

Eq. (4) follows by Lemma 1 and Eq. (5) follows by the fact that  $\mathbf{H}_{k,l-1}$  and  $\mathbf{H}_{k,l}$  are identically distributed conditioned on  $E_{l,k}$  not happening, and moreover  $\mathbb{P}[E_{l,k}] \leq \nu_c + \nu_u + \nu_v$ . The latter can also be seen as follows. By the correctness condition of  $\mathcal{SS}$  we have that  $\sigma_{l,k}$  is valid for  $m_{l,k}$  under  $vk$  except with probability at most  $\nu_c$ . By the assumption that  $\tilde{\mathbf{A}}_k \in \mathcal{A}_{\text{ver}}^{\nu_v}$  we have that  $\tilde{\sigma}_{l,k}$  is also valid for  $m_{l,k}$  under  $vk$  except with probability at most  $\nu_v$ . Finally, by the uniqueness property of  $\mathcal{SS}$  we have that  $\sigma_{l,k}$  and  $\tilde{\sigma}_{l,k}$  must be equal except with probability at most  $\nu_u$ . It follows that  $\mathbb{P}[E_{l,k}] \leq \nu_c + \nu_u + \nu_v$ , as desired.  $\square$

The statement now follows by the above claim and by the triangle inequality, as

$$\mathbb{SD}(\mathbf{G}_0, \mathbf{G}_1) \leq \sum_{k=1}^n \mathbb{SD}(\mathbf{H}_{k-1}, \mathbf{H}_k) \leq \sum_{k=1}^n \sum_{l=1}^q \mathbb{SD}(\mathbf{H}_{k,l-1}, \mathbf{H}_{k,l}) \leq qn \cdot (\nu_c + \nu_u + \nu_v).$$

$\square$

Unfortunately, unique signatures do not satisfy EUF-CMA against the class of all SAs that satisfy relaxed verifiability (cf. Definition 11). In fact, it is not hard to show that no signature scheme with large enough message space (no matter if randomized or deterministic) can achieve EUF-CMA against such class of SAs.

This can be seen by looking again at the class of SAs  $\mathcal{A}_{\text{msg}} = \{\tilde{\mathbf{A}}_{\tilde{m}}\}_{\tilde{m} \in \mathcal{M}}$  that behaves identically to the original signing algorithm, except that upon input  $\tilde{m} \in \mathcal{M}$  it outputs the secret key. Clearly,  $\mathcal{A}_{\text{msg}}$  satisfies relaxed verifiability as a randomly chosen message will be different from  $\tilde{m}$  with high probability user has only a negligible chance of hitting the value  $\tilde{m}$ ; yet  $\mathcal{A}_{\text{msg}}$  clearly allows to break EUF-CMA for an adversary knowing  $\tilde{m}$ .

### 5.3 Random-Message Attacks

We show that if we restrict to the case of random-message attacks (RMA), i.e. the adversary can only see signatures of randomly chosen messages, unique signatures achieve unforgeability against the class of SAs that meets relaxed verifiability (cf. Definition 11).

**Definition 12** (EUF-RMA against SAs). Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a signature scheme, and  $\mathcal{A}$  be some class of SAs for  $\mathcal{SS}$ . We say that  $\mathcal{SS}$  is  $(t, n, q, \varepsilon)$ -EUF-RMA w.r.t. *continuous*  $\mathcal{A}$ -SAs if for all PPT adversaries  $\mathbf{B}$  running in time  $t$ , we have  $\mathbb{P}[\mathbf{B} \text{ wins}] \leq \varepsilon(\kappa)$  in the game of Definition 6 with the adaptation that the messages in the sets  $\mathcal{Q}, \mathcal{Q}_1, \dots, \mathcal{Q}_n$  are drawn uniformly at random from the message space  $\mathcal{M}$ .

While the above definition might seem a weak guarantee, it is still useful for applications. In particular, in Section 5.4 we show how to use any signature scheme that is EUF-RMA against a given class of SAs, to construct an identification scheme that is subversion-resilient against the same class of SAs.

**Theorem 4.** Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a  $(t, (q+1) \cdot n, \varepsilon)$ -EUF-CMA signature scheme with  $\nu_c$ -correctness and  $\nu_u$ -uniqueness, and denote by  $\mathcal{A}_{\text{r\_ver}}^{\nu_v}$  the class of all algorithms that satisfy relaxed  $\nu_v$ -verifiability relative to  $\mathcal{SS}$ . Then  $\mathcal{SS}$  is  $(t', n, q, \varepsilon')$ -indistinguishable against continuous  $\mathcal{A}_{\text{rel\_ver}}^{\nu_v}$ -SAs, for  $t' \approx t$ , for all  $n, q \in \mathbb{N}$ , and for  $\varepsilon' \leq \varepsilon + qn \cdot (\nu_c + \nu_v + \nu_u)$ .

*Proof.* Let  $\mathbf{G}$  be the game of Definition 12. Consider the modified game  $\mathbf{H}$  that is identical to  $\mathbf{G}$  except that queries to the subverted signing algorithms are answered as described below:

- For all  $i \in [q]$ ,  $j \in [n]$ , sample  $\tilde{m}_{i,j} \leftarrow_{\$} \mathcal{M}$  and return  $\sigma_{i,j} \leftarrow \text{Sign}(sk, \tilde{m}_{i,j})$ .

**Claim 5.**  $|\mathbb{P}[\mathbf{B} \text{ wins in } \mathbf{G}] - \mathbb{P}[\mathbf{B} \text{ wins in } \mathbf{H}]| \leq qn \cdot (\nu_c + \nu_v + \nu_u)$ .

*Proof.* For an index  $k \in [0, n]$ , consider the hybrid game  $\mathbf{H}_k$  that answers each query to the  $j$ -th subversion oracle, such that  $j \leq k$ , as in game  $\mathbf{G}$ , while all queries with  $j > k$  are answered as in  $\mathbf{H}$ . We note that  $\mathbf{H}_0 \equiv \mathbf{H}$  and  $\mathbf{H}_n \equiv \mathbf{G}$ . Abusing notation, let us write  $\mathbf{H}_k$  for the distribution of the random variable corresponding to  $\mathbf{B}$ 's view in game  $\mathbf{H}_k$ .

We will show that  $\mathbb{SD}(\mathbf{H}_{k-1}, \mathbf{H}_k) \leq q \cdot (\nu_c + \nu_v + \nu_u)$  for all  $k$ . Fix a particular  $k \in [0, n]$ , and for an index  $l \in [0, q]$  consider the hybrid game  $\mathbf{H}_{k,l}$  that is identical to  $\mathbf{H}_k$  except that it answers queries  $(k, i)$  with  $i \leq l$  as in game  $\mathbf{G}$ , while all queries  $(k, i)$  with  $i > l$  are treated as in  $\mathbf{H}$ . Observe that  $\mathbf{H}_{k,0} \equiv \mathbf{H}_{k-1}$ , and  $\mathbf{H}_{k,q} \equiv \mathbf{H}_k$ .

We now argue that for each  $l \in [q]$ , one has that  $\mathbb{SD}(\mathbf{H}_{k,l-1}, \mathbf{H}_{k,l}) \leq \nu_c + \nu_v + \nu_u$ . Notice that the only difference between  $\mathbf{H}_{k,l-1}$  and  $\mathbf{H}_{k,l}$  is how the two games answer the query  $(k, l)$ : Game  $\mathbf{H}_{k,l-1}$  returns  $\sigma_{l,k} \leftarrow \text{Sign}(sk, \tilde{m}_{l,k})$ , whereas game  $\mathbf{H}_{k,l}$  returns  $\tilde{\sigma}_{l,k} \leftarrow \tilde{\mathbf{A}}_k(sk, \tilde{m}_{l,k})$  (where  $\tilde{m}_{l,k} \leftarrow_{\$} \mathcal{M}$ ). Now let  $E_{l,k}$  be the event that  $\sigma_{l,k} \neq \tilde{\sigma}_{l,k}$ . We can write

$$\mathbb{SD}(\mathbf{H}_{k,l-1}, \mathbf{H}_{k,l}) \leq \mathbb{SD}(\mathbf{H}_{k,l-1}; \mathbf{H}_{k,l} | \neg E_{l,k}) + \mathbb{P}[E_{l,k}] \quad (6)$$

$$\leq \nu_c + \nu_u + \nu_v. \quad (7)$$

Eq. (6) follows by Lemma 1 and Eq. (7) follows by the fact that  $\mathbf{H}_{k,l-1}$  and  $\mathbf{H}_{k,l}$  are identically distributed conditioned on  $E_{l,k}$  not happening, and moreover  $\mathbb{P}[E_{l,k}] \leq \nu_c + \nu_u + \nu_v$ . The latter can also be seen as follows. By the correctness condition of  $\mathcal{SS}$  we have that  $\sigma_{l,k}$  is valid for  $\tilde{m}_{l,k}$  under  $vk$  except with probability at most  $\nu_c$ . By the assumption that  $\tilde{\mathbf{A}}_k \in \mathcal{A}_{\text{r\_ver}}^{\nu_v}$  we have that  $\tilde{\sigma}_{l,k}$  is also valid for  $\tilde{m}_{l,k}$  under  $vk$  except with probability at most  $\nu_v$  (this is because  $\tilde{m}_{l,k}$  is chosen at random). Finally, by the uniqueness property of  $\mathcal{SS}$  we have that  $\sigma_{l,k}$  and  $\tilde{\sigma}_{l,k}$  must be equal except with probability at most  $\nu_u$ . It follows that  $\mathbb{P}[E_{l,k}] \leq \nu_c + \nu_u + \nu_v$ , as desired.

The claim now follows by the above argument and by the triangle inequality, as

$$\begin{aligned} \mathbb{SD}(\mathbf{G}, \mathbf{H}) &\leq \sum_{k=1}^n \mathbb{SD}(\mathbf{H}_{k-1}, \mathbf{H}_k) \\ &\leq \sum_{k=1}^n \sum_{l=1}^q \mathbb{SD}(\mathbf{H}_{k,l-1}, \mathbf{H}_{k,l}) \\ &\leq qn \cdot (\nu_c + \nu_v + \nu_u). \end{aligned}$$

□

**Claim 6.**  $\mathbb{P}[\mathbf{B} \text{ wins in } \mathbf{H}] \leq \varepsilon$ .

*Proof.* Towards a contradiction, assume  $\mathbf{B}$  wins in game  $\mathbf{H}$  with probability larger than  $qn \cdot \varepsilon$ . We build an adversary  $\mathbf{B}'$  (using  $\mathbf{B}$ ) that breaks EUF-CMA of  $\mathcal{SS}$ . Adversary  $\mathbf{B}'$  is described below.

Adversary  $B'$ :

- Receive the verification key  $vk$  from the challenger, and return  $vk$  to  $B$ .
- Upon input the  $i$ -th signature query, query the target oracle receiving back a signature  $\sigma_i \leftarrow \text{Sign}(sk, m_i)$  for  $m_i \leftarrow_{\$} \mathcal{M}$ . Return  $\sigma_i$  to  $B$ .
- Upon input a query of the form  $(j, i)$ , query the target oracle receiving back a signature  $\sigma_{i,j} \leftarrow \text{Sign}(sk, \tilde{m}_{i,j})$  for  $\tilde{m}_{i,j} \leftarrow_{\$} \mathcal{M}$ . Return  $\sigma_{i,j}$  to  $B$ .
- Whenever  $B$  outputs  $(m^*, \sigma^*)$ , output  $(m^*, \sigma^*)$ .

For the analysis, note that  $B'$  runs in time similar to that of  $B$  and asks a total of at most  $q + qn$  signing queries. Moreover, the simulation is perfect and thus  $\mathbb{P}[B' \text{ wins}] \geq \varepsilon$ , a contradiction.  $\square$

The proof follows by combining the above two claims.  $\square$

## 5.4 Subversion-Resilient Identification Schemes

We show how to apply EUF-RMA against SAs to the setting of subversion-resilient identification (ID) schemes. Similar applications already appeared in the literature for leakage and tamper resistance [ADW09, FHN<sup>+</sup>12, DFMV13, NVZ14, FNV15].

In a public-key ID scheme a prover with secret key  $sk$  attempts to prove its identity to a verifier holding the corresponding verification key  $vk$ . More formally, an ID scheme  $\mathcal{ID} = (\text{Setup}, \text{KGen}, \text{P}, \text{V})$  consists of four PPT algorithms described as follows: (1) The parameters generation algorithm takes as input the security parameter and outputs public parameters  $\text{params} \leftarrow \text{Setup}(1^\kappa)$ , shared by all users.<sup>16</sup> (2) The key generation algorithm takes as input the security parameter and outputs a verification key/secret key pair  $(vk, sk) \leftarrow \text{KGen}(1^\kappa)$ . (3)  $\text{P}$  and  $\text{V}$  are probabilistic Turing machines interacting in a protocol; at the end of the execution  $\text{V}$  outputs a decision bit  $d \in \{0, 1\}$ , where  $d = 1$  means that the identification was successful. We write  $\langle \text{P}(sk), \text{V}(vk) \rangle$  for the random variable corresponding to the verifier's verdict, and  $\text{P}(sk) \rightleftharpoons \text{V}(vk)$  for the random variable corresponding to transcripts of honest protocol executions.

We now define a variant of passive security, where in a first phase the adversary is allowed to subvert the prover algorithm; in a second phase the adversary has to impersonate the prover. Similarly to the case of signature schemes subversion is modelled by considering a class  $\mathcal{A}$  of SAs, where each  $\tilde{A} \in \mathcal{A}$  is an algorithm replacing the prover algorithm  $\text{P}$  within the ID scheme  $\mathcal{ID}$ .

**Definition 13** (Subversion-Resilient Identification). Let  $\mathcal{ID} = (\text{Setup}, \text{KGen}, \text{Sign}, \text{Vrfy})$  be an ID scheme, and  $\mathcal{A}$  be some class of SAs for  $\mathcal{ID}$ . We say that  $\mathcal{ID}$  is  $(t, n, q, \varepsilon)$ -secure w.r.t. *continuous*  $\mathcal{A}$ -SAs if for all PPT adversaries  $B$  running in time  $t$ , we have  $\mathbb{P}[B \text{ wins}] \leq \varepsilon(\kappa)$  in the following game:

1. The challenger runs  $\text{params} \leftarrow \text{Setup}(1^\kappa)$ ,  $(vk, sk) \leftarrow \text{KGen}(1^\kappa)$ , and forwards  $(\text{params}, vk)$  to  $B$ .
2. The adversary  $B$  can observe  $q$  transcripts  $\text{P}(sk) \rightleftharpoons \text{V}(vk)$  corresponding to honest protocol executions between the prover and the verifier.

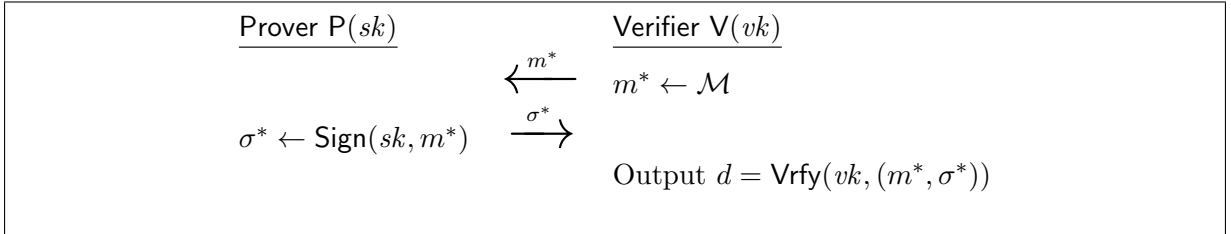
---

<sup>16</sup>In what follows all algorithms take as input  $\text{params}$ , but we omit to explicitly write this for ease of notation.

3. For each  $j \in [n]$ , the adversary  $B$  can adaptively choose an algorithm  $\tilde{A}_j \in \mathcal{A}$ . For each algorithm,  $B$  can observe  $q$  transcripts  $\tilde{A}_j(sk) \rightleftharpoons V(vk)$  corresponding to protocol executions between the subverted prover and the verifier.
4. The adversary  $B$  loses access to all oracles and plays the role of the prover in an execution with an honest verifier  $d \leftarrow \langle B(vk), V(vk) \rangle$ ; we say that  $B$  wins if and only if  $d = 1$ .

Consider the following standard construction (see, e.g., [BFGM01]) of an identification scheme  $\mathcal{ID}$  from a signature scheme  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$ .

- *Parameters generation.* Algorithm  $\text{Setup}$  samples the public parameters  $\text{params}$  for the signature schemes (if any).
- *Key Generation.* Algorithm  $\text{KGen}$  runs the key generation algorithm of the signature scheme, obtaining  $(vk, sk) \leftarrow \text{KGen}(1^\kappa)$ .
- *Identification protocol.* The interaction  $P(sk) \rightleftharpoons V(vk)$  is depicted in Figure 3.



**Figure 3:** Two-round identification using a signature scheme  $\mathcal{SS}$  with message space  $\mathcal{M}$

The theorem below states that the above protocol achieve subversion resilience w.r.t. a given class  $\mathcal{A}$  of SAs, provided that the underlying signature scheme is EUF-RMA w.r.t. the same class  $\mathcal{A}$ .

**Theorem 5.** *Let  $\mathcal{SS}$  be a signature scheme with message space  $\mathcal{M}$ , and let  $\mathcal{A}$  be a class of SAs for  $\mathcal{SS}$ . Assume that  $\mathcal{SS}$  is  $(t, n, q, \varepsilon)$ -EUF-RMA w.r.t. continuous  $\mathcal{A}$ -SAs. Then the ID scheme  $\mathcal{ID}$  from Figure 3 is  $(t', n, q, \varepsilon')$ -secure w.r.t. continuous  $\mathcal{A}$ -SAs where  $t' \approx t$  and  $\varepsilon' \leq \varepsilon + \frac{(n+1)q}{|\mathcal{M}|}$ .*

*Proof.* For the sake of contradiction, assume that there exists an adversary  $B$  breaking security of the identification scheme with probability larger than  $\varepsilon + \frac{(n+1)q}{|\mathcal{M}|}$ . We construct a PPT adversary  $B'$  breaking EUF-RMA of  $\mathcal{SS}$  with probability at least  $\varepsilon$  (a contradiction). Adversary  $B'$  runs in the game of Definition 12 and is described below. The main observation is that the prover's algorithm  $P$  is completely specified by algorithm  $\text{Sign}$ , and thus subverting the ID scheme is equivalent to subverting the signature scheme.

Adversary  $B'$ :

1. Receive the public parameters  $\text{params}$  and the verification key  $vk$  for  $\mathcal{SS}$  and forward  $(\text{params}, vk)$  to  $B$ .
2. Whenever  $B$  wants to observe a honest transcript  $P(sk) \rightleftharpoons V(vk)$ , query the signing oracle obtaining a pair  $(m_i, \sigma_i)$  such that  $\sigma_i \leftarrow \text{Sign}(sk, m_i)$  and  $m_i \leftarrow^* \mathcal{M}$ . Forward  $(m_i, \sigma_i)$  to  $B$ .
3. Whenever  $B$  specifies an algorithm  $\tilde{A}_j \in \mathcal{A}$ , forward  $\tilde{A}_j$  to the challenger. For each query of  $B$  to its own  $j$ -th oracle, query the target  $j$ -th oracle obtaining a pair  $(\tilde{m}_{i,j}, \tilde{\sigma}_{i,j})$  such that  $\tilde{\sigma}_{i,j} \leftarrow \tilde{A}_j(sk, \tilde{m}_{i,j})$  and  $\tilde{m}_{i,j} \leftarrow^* \mathcal{M}$ . Forward  $(\tilde{m}_{i,j}, \tilde{\sigma}_{i,j})$  to  $B$ .

4. Finally, when  $B$  is ready to start the impersonation phase, sample a random message  $m^* \leftarrow_s \mathcal{M}$  and send it to  $B$ . Upon receiving a value  $\sigma^*$  from  $B$  output  $(m^*, \sigma^*)$  as forgery.

It is easy to see that  $B'$ 's simulation of  $B$ 's queries is perfect; moreover, since the message  $m^*$  in the impersonation stage is chosen at random from  $\mathcal{M}$ , also the simulation of this phase has the right distribution and in particular the forgery  $(m^*, \sigma^*)$  will be valid with probability  $\varepsilon$ .

It remains to compute the probability that  $B'$  is successful. Observe that  $B'$  is successful whenever  $(m^*, \sigma^*)$  is valid and  $m^* \notin \mathcal{Q} \cup \tilde{\mathcal{Q}}$ . Also, note that  $m^*$  is independent from  $\tilde{\mathcal{Q}}$ , so in particular

$$\mathbb{P}[m^* \in \mathcal{Q} \cup \tilde{\mathcal{Q}}] \leq \frac{|\mathcal{Q}| + |\tilde{\mathcal{Q}}|}{|\mathcal{M}|} = \frac{(n+1)q}{|\mathcal{M}|}.$$

Let  $E$  be the event that  $m^* \notin \mathcal{Q} \cup \tilde{\mathcal{Q}}$ . We have,

$$\begin{aligned} \mathbb{P}[B' \text{ wins}] &\geq \mathbb{P}[B \text{ wins} \wedge E] \geq \mathbb{P}[B \text{ wins}] - \mathbb{P}[\neg E] \\ &\geq \mathbb{P}[B \text{ wins}] - \frac{(n+1)q}{|\mathcal{M}|} \\ &> \varepsilon, \end{aligned}$$

where the last inequality follows by our initial assumption on  $B$ 's advantage. This concludes the proof.  $\square$

## 6 Reverse Firewalls for Signatures

In Section 5 we have shown that unique signatures are secure against a restricted class of SAs, namely all SAs that meet the so-called verifiability condition. As discussed in Section 3, by removing the latter requirement (i.e., allowing for arbitrary classes of SAs in Definition 5 and 6) would require that a signature scheme  $\mathcal{SS}$  remains unforgeable even against an adversary allowed *arbitrary tampering with the computation* performed by the signing algorithm. This is impossible without making further assumptions.

In this section we explore to what extent one can model signature schemes secure against arbitrary tampering with the computation, by making the extra assumption of an un-tamperable cryptographic reverse firewall (RF) [MS15]. Roughly, a RF for a signature scheme is a (possibly stateful) algorithm that takes as input a message/signature pair and outputs an updated signature; importantly the firewall has to do so using only public information (in particular, without knowing the signing key). A formal definition follows.

**Definition 14** (RF for signatures). Let  $\mathcal{SS}$  be a signature scheme. A RF for  $\mathcal{SS}$  is a pair of algorithms  $\mathcal{FW} = (\text{Setup}, \text{Patch})$  specified as follows: (i) **Setup** takes as input the security parameter and a verification key  $vk \in \mathcal{VK}$ , and outputs some initial (public) state  $\delta \in \{0, 1\}^*$ ; (ii) **Patch** takes as input the current (public) state  $\delta$ , and a message/signature pair  $(m, \sigma)$  and outputs a possibly modified signature or a special symbol  $\perp$  and an updated (public) state  $\delta'$ . We write this as  $\sigma' \leftarrow \text{Patch}_\delta(m, \sigma)$  (and omit to denote the updated state  $\delta'$  as an explicit output).

We will typically assume that the current state  $\delta_{\text{cur}}$  of the RF, can be computed efficiently given just the verification key  $vk$ , the initial state  $\delta$  and the entire history of all inputs to the RF.

## 6.1 Properties

Below, we discuss the correctness and security requirements of cryptographic RF  $\mathcal{FW}$  for a signature scheme  $\mathcal{SS}$ .

**Maintaining functionality.** The first basic property of a RF is that it should preserve the functionality of the underlying signature scheme, i.e. if a signature  $\sigma$  on a message  $m$  is computed using signing key  $sk$ , and the firewall is initialized with the corresponding verification key  $vk$ , the patched signatures  $\sigma'$  should (almost always) be a valid signatures for  $m$  under  $vk$ . More precisely, we say that  $\mathcal{FW}$  is *functionality maintaining* for  $\mathcal{SS}$ , if for any polynomial  $p(\kappa)$  and any vector of inputs  $(m_1, \dots, m_p) \in \mathcal{M}$ , there exists a negligible function  $\nu : \mathbb{N} \rightarrow [0, 1]$  such that

$$\mathbb{P} \left[ \exists i \in [p] \text{ s.t. } \text{Vrfy}(vk, (m_i, \sigma'_i)) = 0 : \begin{array}{l} (vk, sk) \leftarrow \text{KGen}(1^\kappa), \delta \leftarrow \text{Setup}(vk, 1^\kappa) \\ \sigma_1 \leftarrow \text{Sign}(sk, m_1), \dots, \sigma_p \leftarrow \text{Sign}(sk, m_p) \\ \sigma'_1 \leftarrow \text{Patch}_\delta(m_1, \sigma_1), \dots, \sigma'_p \leftarrow \text{Patch}_\delta(m_p, \sigma_p) \end{array} \right] \leq \nu(\kappa),$$

where the probability is taken over the coin tosses of all involved algorithms. Recall that each invocation of algorithm  $\text{Patch}$  updates the (public) state  $\delta$  of the RF.

**Preserving Unforgeability.** The second property of a RF is a security requirement. Note that a firewall can never “create” security (as it does not know the signing key). Below we define what it means for a RF to *preserve* unforgeability of a signature scheme against *arbitrary* tampering attacks.

**Definition 15** (Unforgeability preserving RF). Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a signature scheme with RF  $\mathcal{FW} = (\text{Setup}, \text{Patch})$ . We say that  $\mathcal{FW}$   $(t, n, q, \varepsilon)$ -preserves unforgeability for  $\mathcal{SS}$  against continuous SAs if for all adversaries  $\mathbf{B}$  running in time  $t$  we have that  $\mathbb{P}[\mathbf{B} \text{ wins}] \leq \varepsilon$  in the following game:

1. The challenger runs  $(vk, sk) \leftarrow \text{KGen}(1^\kappa)$ ,  $\delta \leftarrow \text{Setup}(vk, 1^\kappa)$ , and gives  $(vk, \delta)$  to  $\mathbf{B}$ .
2. The adversary  $\mathbf{B}$  is given oracle access to  $\text{Sign}(sk, \cdot)$ . Upon input the  $i$ -th query  $m_i$ , this oracle returns  $\sigma_i \leftarrow \text{Sign}(sk, m_i)$ . Let  $\mathcal{Q} = \{m_1, \dots, m_q\}$  be the set of all signature queries.
3. The adversary  $\mathbf{B}$  can adaptively choose an arbitrary algorithm  $\tilde{\mathbf{A}}_j$ , and correspondingly obtain oracle access to  $\text{Patch}_\delta(\cdot, \tilde{\mathbf{A}}_j(sk, \cdot))$ :
  - Upon input the  $i$ -th query  $\tilde{m}_{i,j}$ , for  $i \in [q]$  and  $j \in [n]$ , the oracle returns  $\tilde{\sigma}_{i,j} \leftarrow \text{Patch}_\delta(\tilde{m}_{i,j}, \tilde{\mathbf{A}}_j(sk, \tilde{m}_{i,j}))$  and updates the public state  $\delta$ ;
  - Whenever  $\tilde{\sigma}_{i,j} = \perp$  the oracle enters a special self-destructs mode, in which the answer to all future queries is by default set to  $\perp$ .

Let  $\tilde{\mathcal{Q}}_j = \{\tilde{m}_{1,j}, \dots, \tilde{m}_{q,j}\}$  be the set of all queries for each  $\tilde{\mathbf{A}}_j$ .

4. Finally,  $\mathbf{B}$  outputs a pair  $(m^*, \sigma^*)$ ; we say that  $\mathbf{B}$  wins iff  $\text{Vrfy}(vk, (m^*, \sigma^*)) = 1$  and  $m^* \notin \mathcal{Q} \cup \tilde{\mathcal{Q}}$ , where  $\tilde{\mathcal{Q}} := \bigcup_{j=1}^n \tilde{\mathcal{Q}}_j$ .

Whenever  $t = \text{poly}(\kappa)$ ,  $q = \text{poly}(\kappa)$  and  $\varepsilon = \text{negl}(\kappa)$  we simply say that  $\mathcal{FW}$  preserves unforgeability for  $\mathcal{SS}$ . Furthermore, in case  $\mathbf{A}$  specifies all of its queries  $\{\tilde{\mathbf{A}}_j, \tilde{m}_{i,j}\}_{j \in [n], i \in [q]}$  at the same time we say that  $\mathcal{FW}$  *non-adaptively* preserves unforgeability.



We observe that Definition 15 is very similar to Definition 6, except for a few crucial differences. First, note that the above definition considers arbitrary classes of SAs instead of SAs within a given class  $\mathcal{A}$ ; this is possible because the output of each invocation of the subverted signing algorithm is patched using the firewall (which is assumed to be un-tamperable).

Second, observe that the above definition relies on the so-called self-destruct capability: Whenever the firewall returns  $\perp$ , all further queries to any of the oracles results in  $\perp$ ; as we show in Section 6.2 this is necessary as without such a capability there exists simple generic attacks that allow for complete security breaches. We stress, however, that the assumption of the self-destruct capability does not make the problem of designing an unforgeability preserving reverse firewall trivial. In fact, the biased-randomness attacks of Section 4 allow to break all randomized scheme *without* ever provoking a self-destruct. On the positive side, in Section 6.3, we show how to design an unforgeability preserving RF for any *re-randomizable* signature scheme.

**Exfiltration resistance.** More in general, one might require a stronger security property from a RF. Namely, we could ask that patched signatures are indistinguishable from real signatures to the eyes of an attacker. This property, which is called exfiltration resistance in [MS15], would be similar in spirit to our definition of indistinguishability w.r.t. continuous SAs (see Definition 5).

It is not hard to see that exfiltration resistance against arbitrary SAs is impossible to achieve in the case of signature schemes; this is because the attacker could simply set the subverted signing algorithm to always output the all-zero string, in which case the RF has no way to patch its input to a valid signature (and thus the adversary can easily distinguish subverted patched signatures from real signatures).<sup>17</sup>

## 6.2 Necessity of Self-Destruct

We show that no RF can preserve both functionality and unforgeability, without assuming the self-destruct capability. This is achieved via a generic (non-adaptive) attack that allows to extract the secret key in case the RF does not self-destruct. The attack itself is a generalization of a similar attack by Gennaro *et al.* [GLM<sup>+</sup>04] in the context of memory tampering.

**Theorem 6.** *Let  $\mathcal{SS}$  be an EUF-CMA signature scheme. No RF  $\mathcal{FW}$  can at the same time be functionality maintaining and non-adaptively  $(\text{poly}(\kappa), 1, \text{poly}(\kappa), \text{negl}(\kappa))$ -preserve unforgeability for  $\mathcal{SS}$ , without assuming the self-destruct capability.*

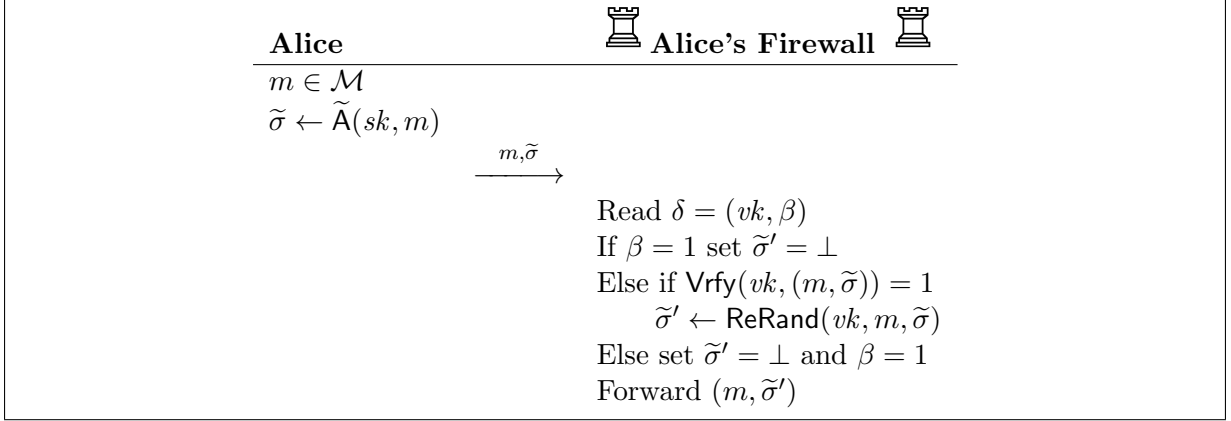
*Proof sketch.* Consider the following adversary  $\mathcal{B}$  playing the game of Definition 15 (omitting the self-destruct capability).

- Upon input the verification key  $vk$ , and the initial state  $\delta$ , initialize  $\tau := 1$ .
- Forward  $\tilde{\mathcal{A}}_\tau$  to the challenger, where algorithm  $\tilde{\mathcal{A}}_\tau$  is defined as follows: Upon input a message  $\tilde{m}_i$ , set  $j = \tau \bmod \ell$  (where  $\ell := |sk|$ ) and
  - If  $sk[j] = 1$ , output  $\tilde{\sigma}_i \leftarrow \text{Sign}(sk, \tilde{m}_i)$ .
  - Else, output  $0^{|\sigma|}$ .

Update  $\tau \leftarrow \tau + 1$ .

---

<sup>17</sup>We note, however, that our techniques from Section 5 can be extended to design a RF that is *weakly* exfiltration resistant, namely it is exfiltration resistant against restricted SAs that satisfy the verifiability condition.



**Figure 4:** A cryptographic reverse firewall preserving unforgeability of any re-randomizable signature scheme against arbitrary SAs.

- Let  $(\bar{m}, \tilde{\sigma}'_1), \dots, (\bar{m}, \tilde{\sigma}'_\ell)$  be the set of tampered signature queries (and answers to these queries) asked by  $\mathbf{B}$ , where  $\tilde{\sigma}'_i \leftarrow \mathbf{Patch}_\delta(\bar{m}, \tilde{\mathbf{A}}_\tau(sk, \bar{m}))$ . Define  $sk'[i] = \mathbf{Vrfy}(vk, (\bar{m}, \tilde{\sigma}'_i))$  and return  $sk' := (sk'[1], \dots, sk'[\ell])$ .

Notice that  $\mathbf{B}$  specifies its queries non-adaptively, and moreover it only uses one subversion which is queried upon a fixed message  $\bar{m} \in \mathcal{M}$ . We will show that the extracted key  $sk'$  is equal to the original secret key  $sk$  with overwhelming probability, which clearly implies the statement. The proof is by induction; assume that the statement is true up to some index  $i \geq 1$ . We claim that  $sk'[i+1] = sk[i+1]$  with all but negligible probability. To see this, define the event  $E_{i+1}$  that  $sk[i+1] = 0$  and  $\mathbf{Vrfy}(vk, (\bar{m}, \tilde{\sigma}'_{i+1})) = 1$  or  $sk[i+1] = 1$  and  $\mathbf{Vrfy}(vk, (\bar{m}, \tilde{\sigma}'_{i+1})) = 0$ . By the assumption that the RF does not self-destruct and is functionality maintaining, we get that the latter sub-case happens only with negligible probability. On the other hand, if the former sub-case happens we get that the RF forged a signature on  $\bar{m}$ , which contradicts EUF-CMA security of  $\mathcal{SS}$ . By a union bound, we get that  $\mathbb{P}[E_{i+1}]$  is negligible as desired.  $\square$

### 6.3 Patching Re-Randomizable Signatures

We design a RF preserving unforgeability of so-called *re-randomizable* signature schemes (that include unique signatures as a special case).

**Definition 16** (Re-randomizable signatures). A signature scheme  $\mathcal{SS} = (\mathbf{KGen}, \mathbf{Sign}, \mathbf{Vrfy})$  is efficiently  $\nu_r$ -re-randomizable, if there exists a PPT algorithm  $\mathbf{ReRand}$  such that for all messages  $m \in \mathcal{M}$  and for all  $(vk, sk) \leftarrow \mathbf{KGen}(1^\kappa)$  and  $\sigma \leftarrow \mathbf{Sign}(sk, m)$ , we have that  $\mathbb{SD}(\mathbf{ReRand}(vk, m, \sigma); \mathbf{Sign}(sk, m)) \leq \nu_r$ .

Note that unique signatures are efficiently re-randomizable, for  $\mathbf{ReRand}(vk, m, \sigma) = \sigma$  and  $\nu_r = 0$ ; Waters' signature scheme [Wat05], and its variant by Hofheinz *et al.* [HJK12], are also efficiently re-randomizable.

Our firewall, which is formally described in Fig. 4, first checks if  $\sigma$  is a valid signature on message  $m$  under key  $vk$  (provided that a self-destruct was not provoked yet). If not, it self-destructs and returns  $\perp$ ; otherwise it re-randomizes  $\sigma$  and outputs the result. The self-destruct capability is implemented using a one-time writable bit  $\beta$  (which is included in the public state).

**Theorem 7.** Let  $\mathcal{SS}$  be a  $(t, (q+1)n, \varepsilon)$ -EUF-CMA signature scheme that is efficiently  $\nu_r$ -re-randomizable and that satisfies  $\nu_c$ -correctness. Then, the RF of Fig. 4 maintains functionality and  $(t', q, \varepsilon')$ -preserves unforgeability for  $\mathcal{SS}$ , where  $t' \approx t$  and  $\varepsilon' \leq qn \cdot (\nu_c + \nu_r + \varepsilon)$ .

*Proof.* The fact that the firewall maintains functionality follows directly by  $\nu_c$ -correctness of  $\mathcal{SS}$ . We now show the firewall preserves unforgeability. Let  $\mathbf{G}$  be the game of Definition 15; we write  $(i^*, j^*) \in [q] \times [n]$  for the pair of indexes in which the firewall self-destructs (if any). Consider the modified game  $\mathbf{H}$  that is identical to  $\mathbf{G}$  except that tampered signature queries are answered as described below:

- For all  $j < j^*$ , upon input  $(j, \tilde{m}_{i,j})$  return  $\sigma_{i,j} \leftarrow \text{Sign}(sk, \tilde{m}_{i,j})$  for all  $i \in [q]$ .
- For  $j = j^*$ , upon input  $(j, \tilde{m}_{i,j})$  if  $i < i^*$  return  $\sigma_{i,j} \leftarrow \text{Sign}(sk, \tilde{m}_{i,j})$ ; else return  $\perp$ .
- For all  $j > j^*$ , upon input message  $\tilde{m}_{i,j}$  return  $\perp$  for all  $i \in [q]$ .

**Claim 7.**  $|\mathbb{P}[\mathbf{B} \text{ wins in } \mathbf{G}] - \mathbb{P}[\mathbf{B} \text{ wins in } \mathbf{H}]| \leq qn \cdot (\nu_c + \nu_r)$ .

*Proof.* For an index  $k \in [0, n]$ , consider the hybrid game  $\mathbf{H}_k$  that answers each query  $(j, \tilde{m}_{i,j})$  such that  $j \leq k$  as in game  $\mathbf{G}$ , while all queries  $(j, \tilde{m}_{i,j})$  such that  $j > k$  are answered as in  $\mathbf{H}$ . We note that  $\mathbf{H}_0 \equiv \mathbf{H}$  and  $\mathbf{H}_n \equiv \mathbf{G}$ . Abusing notation, let us write  $\mathbf{H}_k$  for the distribution of the random variable corresponding to  $\mathbf{B}$ 's view in game  $\mathbf{H}_k$ .

We will show that  $\mathbb{SD}(\mathbf{H}_{k-1}, \mathbf{H}_k) \leq q \cdot (\nu_c + \nu_r)$  for all  $k$ . Fix a particular  $k \in [0, n]$ , and for an index  $l \in [0, q]$  consider the hybrid game  $\mathbf{H}_{k,l}$  that is identical to  $\mathbf{H}_k$  except that it answers queries  $(k, \tilde{m}_{i,k})$  with  $i \leq l$  as in game  $\mathbf{G}$ , while all queries  $(k, \tilde{m}_{i,k})$  with  $i > l$  are treated as in  $\mathbf{H}$ . Observe that  $\mathbf{H}_{k,0} \equiv \mathbf{H}_{k-1}$ , and  $\mathbf{H}_{k,q} \equiv \mathbf{H}_k$ .

We now argue that for each  $l \in [q]$ , one has that  $\mathbb{SD}(\mathbf{H}_{k,l-1}, \mathbf{H}_{k,l}) \leq \nu_c + \nu_r$ . Observe that, since for  $k > j^*$  both games always return  $\perp$ , we can assume without loss of generality that  $k \leq j^*$ . Note that the only difference between  $\mathbf{H}_{k,l-1}$  and  $\mathbf{H}_{k,l}$  is how the two games answer the query  $(k, \tilde{m}_{l,k})$ :  $\mathbf{H}_{k,l-1}$  returns  $\sigma_{l,k} \leftarrow \text{Sign}(sk, \tilde{m}_{l,k})$  whereas  $\mathbf{H}_{k,l}$  returns  $\tilde{\sigma}'_{l,k} \leftarrow \text{Patch}_\delta(\tilde{m}_{l,k}, \tilde{\sigma}_{l,k})$  where  $\tilde{\sigma}_{l,k} \leftarrow \tilde{\mathbf{A}}_k(sk, \tilde{m}_{l,k})$ . Let  $E_{l,k}$  be the event that  $\text{Vrfy}(vk, (\tilde{m}_{l,k}, \sigma_{l,k})) = 0$ . We have

$$\mathbb{SD}(\mathbf{H}_{k,l-1}; \mathbf{H}_{k,l}) \leq \mathbb{SD}(\mathbf{H}_{k,l-1}; \mathbf{H}_{k,l} | \neg E_{l,k}) + \mathbb{P}[E_{l,k}] \quad (8)$$

$$\leq \nu_r + \nu_c. \quad (9)$$

Eq. (8) follows by Lemma 1 and Eq. (9) by the fact that  $\mathbf{H}_{k,l-1}$  and  $\mathbf{H}_{k,l}$  are statistically close (up to distance  $\nu_r$ ) conditioned on  $E_{l,k}$  not happening, and moreover  $\mathbb{P}[E_{l,k}] \leq \nu_c$ . The former is because signatures are re-randomizable, and thus (as long as the firewall did not self-destruct) the output of  $\text{ReRand}$  is statistically close (up to distance  $\nu_r$ ) to the output of the original signing algorithm; the latter follows by  $\nu_c$ -correctness of the signature scheme.

The statement now follows by the above argument and by the triangle inequality, as

$$\begin{aligned} \mathbb{SD}(\mathbf{G}, \mathbf{H}) &\leq \sum_{k=1}^n \mathbb{SD}(\mathbf{H}_{k-1}, \mathbf{H}_k) \\ &\leq \sum_{k=1}^n \sum_{l=1}^q \mathbb{SD}(\mathbf{H}_{k,l-1}, \mathbf{H}_{k,l}) \\ &\leq qn \cdot (\nu_c + \nu_r). \end{aligned}$$

□

**Claim 8.**  $\mathbb{P}[\mathbf{B} \text{ wins in } \mathbf{H}] \leq qn \cdot \varepsilon$ .

*Proof.* Towards a contradiction, assume  $\mathbf{B}$  wins in game  $\mathbf{H}$  with probability larger than  $qn \cdot \varepsilon$ . Wlog. we assume that  $\mathbf{B}$  always outputs its forgery after provoking a self-destruct.<sup>18</sup> We build an adversary  $\mathbf{B}'$  (using  $\mathbf{B}$ ) that breaks EUF-CMA of  $\mathcal{SS}$ . Adversary  $\mathbf{B}'$  is described below.

Adversary  $\mathbf{B}'$ :

- Receive the verification key  $vk$  from the challenger, sample a random pair  $(j^*, i^*) \leftarrow_{\$} [n] \times [q]$ , and return  $vk$  to  $\mathbf{B}$ .
- Upon input the  $i$ -th signature query  $m_i$ , forward this value to the signing oracle receiving back a signature  $\sigma_i \leftarrow \text{Sign}(sk, m_i)$ . Return  $\sigma_i$  to  $\mathbf{B}$ .
- Upon input a query of the form  $(j, \tilde{m}_{i,j})$  answer as follows:
  - In case  $j < j^*$ , forward  $\tilde{m}_{i,j}$  to the signing oracle, obtaining  $\tilde{\sigma}_{i,j} \leftarrow \text{Sign}(sk, \tilde{m}_{i,j})$ , and return  $\tilde{\sigma}_{i,j}$  to  $\mathbf{B}$ .
  - In case  $j = j^*$ , if  $i < i^*$  forward  $\tilde{m}_{i,j}$  to the signing oracle, obtaining  $\tilde{\sigma}_{i,j} \leftarrow \text{Sign}(sk, \tilde{m}_{i,j})$ , and return  $\tilde{\sigma}_{i,j}$  to  $\mathbf{B}$ . Else, return  $\perp$ .
  - In case  $j > j^*$  answer with  $\perp$ .
- Whenever  $\mathbf{B}$  outputs  $(m^*, \sigma^*)$ , output  $(m^*, \sigma^*)$ .

For the analysis, note that  $\mathbf{B}'$  runs in time similar to that of  $\mathbf{B}$  and asks a total of at most  $q + qn$  signing queries. Moreover, define the event  $E$  that  $\mathbf{B}'$  guesses correctly the query  $(j^*, i^*)$  where  $\mathbf{B}$  provokes a self-destruct. Clearly, in case  $E$  happens we have that  $\mathbf{B}'$  perfectly simulates the distribution of game  $\mathbf{H}$ . Hence  $\mathbb{P}[\mathbf{B}' \text{ wins}] \geq (qn \cdot \varepsilon)/(qn) = \varepsilon$ , a contradiction.  $\square$

The proof follows by combining the above two claims.  $\square$

## 7 The Multi-User Setting

In this section we consider the multi-user setting for all definitions that appear in this paper. We also provide a complete picture of relationships between all definitions, as shown in Fig. 5 and Fig. 6.

### 7.1 Multi-User Impersonation

Analogous to the single-user setting, we consider two security definitions corresponding to different adversarial goals.

In the indistinguishability definition for the multi-user setting adversary  $\mathbf{B}$  now receives  $u \geq 1$  key pairs from the challenger and can continuously subvert each user independently. A formal definition follows.

**Definition 17** (Indistinguishability against SAs—Multi-User). Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a signature scheme, and  $\mathcal{A}$  be some class of SAs for  $\mathcal{SS}$ . We say that  $\mathcal{SS}$  is  $u$ -users indistinguishable w.r.t *continuous*  $\mathcal{A}$ -SAs if for all PPT adversaries  $\mathbf{B}$  there exists a negligible function  $\nu : \mathbb{N} \rightarrow [0, 1]$ , such that  $|\mathbb{P}[\mathbf{B} \text{ wins}] - \frac{1}{2}| \leq \nu(\kappa)$  in the following game:

1. The challenger samples  $b \leftarrow_{\$} \{0, 1\}$ , generates  $(vk_i, sk_i) \leftarrow \text{KGen}(1^\kappa)$  for  $i \in [u]$  and gives  $vk_1, \dots, vk_u$  to  $\mathbf{B}$ .

---

<sup>18</sup>If not we can always modify  $\mathbf{B}$  in such a way that it asks one additional query provoking a self-destruct; this clearly does not decrease  $\mathbf{B}$ 's advantage.

2. The adversary  $\mathbf{B}$  can specify polynomially many queries (adaptively and in an arbitrary order) of the form  $(i, \tilde{\mathbf{A}})$  for  $i \in [u]$ .
  - (a) For each such query,  $\mathbf{B}$  is given access to an oracle that can be queried polynomially many times on inputs  $m \in \mathcal{M}$ .
  - (b) The answer to each query  $m$  depends on the value of the secret bit  $b$ . In particular, if  $b = 1$ , the output is  $\sigma \leftarrow \text{Sign}(sk_i, m)$ ; if  $b = 0$ , the output is  $\tilde{\sigma} \leftarrow \tilde{\mathbf{A}}(sk_i, m)$ .
3. Finally,  $\mathbf{B}$  outputs a value  $b' \in \{0, 1\}$ ; we say that  $\mathbf{B}$  wins iff  $b' = b$ .

In the impersonation definition for the multi-user setting adversary  $\mathbf{B}$  now receives  $u \geq 1$  key pairs from the challenger and can continuously subvert each user independently; adversary  $\mathbf{B}$  is successful if it can impersonate *any* of the users. A formal definition follows.

**Definition 18** (EUF-CMA against SAs—Multi-User). Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a signature scheme, and  $\mathcal{A}$  be some class of SAs for  $\mathcal{SS}$ . We say that  $\mathcal{SS}$  is  $u$ -users EUF-CMA w.r.t. *continuous*  $\mathcal{A}$ -SAs if for all PPT adversaries  $\mathbf{B}$  there exists a negligible function  $\nu : \mathbb{N} \rightarrow [0, 1]$ , such that  $\mathbb{P}[\mathbf{B} \text{ wins}] \leq \nu(\kappa)$  in the following game:

1. The challenger generates  $(vk_i, sk_i) \leftarrow \text{KGen}(1^\kappa)$  for  $i \in [u]$  and gives  $vk_1, \dots, vk_u$  to  $\mathbf{B}$ .
2. Adversary  $\mathbf{B}$  is given oracle access to  $\text{Sign}(sk_i, \cdot)$ . Upon input query  $m \in \mathcal{M}$ , this oracle returns  $\sigma \leftarrow \text{Sign}(sk_i, m)$ ; let  $\mathcal{Q}$  be the set of all messages queried to this oracle.
3. The adversary  $\mathbf{B}$  can specify polynomially many queries (adaptively and in an arbitrary order) of the form  $(i, \tilde{\mathbf{A}})$  for  $i \in [u]$ .
  - (a) For each such query,  $\mathbf{B}$  is given access to an oracle that can be queried polynomially many times upon inputs  $m \in \mathcal{M}$ .
  - (b) The answer to each query  $m$  is  $\tilde{\sigma} \leftarrow \tilde{\mathbf{A}}(sk_i, m)$ ; let  $\tilde{\mathcal{Q}}$  be the set containing all queried messages to all oracles.
4. Finally,  $\mathbf{B}$  outputs a tuple  $(m^*, \sigma^*, i^*)$ ; we say that  $\mathbf{B}$  wins iff  $\text{Vrfy}(vk_{i^*}, (m^*, \sigma^*)) = 1$  and  $m^* \notin \mathcal{Q} \cup \tilde{\mathcal{Q}}$ .

## 7.2 Multi-User Public/Secret Undetectability

In the undetectability definition for the multi-user setting user  $\mathbf{U}$  now receives  $u \geq 1$  key pairs from the challenger (only the verification keys for public undetectability) and is allowed to make polynomially many signature queries for all users (key pairs). The answer to these queries are either computed using the real signature algorithm or a subverted algorithm previously chosen by the challenger possibly depending on the verification keys of the users. A formal definition follows.

**Definition 19** (Public/Secret Undetectability—Multi-User). Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a signature scheme, and  $\mathcal{A}$  be some class of SAs for  $\mathcal{SS}$ . We say that  $\mathcal{A}$  is  $u$ -users *secretly* undetectable w.r.t.  $\mathcal{SS}$  if for all PPT users  $\mathbf{U}$ , there exists a negligible function  $\nu : \mathbb{N} \rightarrow [0, 1]$  and an efficient challenger such that  $|\mathbb{P}[\mathbf{U} \text{ wins}] - \frac{1}{2}| \leq \nu(\kappa)$  in the following game:

1. The challenger samples  $b \leftarrow_{\$} \{0, 1\}$ , generates  $(vk_i, sk_i) \leftarrow \text{KGen}(1^\kappa)$  for  $i \in [u]$ , chooses  $\tilde{\mathbf{A}} \leftarrow \mathcal{A}$  (possibly depending on  $vk_1, \dots, vk_u$ ), and gives  $(vk_1, sk_1, \dots, vk_u, sk_u)$  to  $\mathbf{B}$ .

2. The user  $U$  can ask polynomially many queries of the form  $(i, m)$ , where  $i \in [u]$  and  $m \in \mathcal{M}$ . The answer to each query depends on the secret bit  $b$ . In particular, if  $b = 1$ , the challenger returns  $\sigma \leftarrow \text{Sign}(sk_i, m)$ ; if  $b = 0$ , the challenger returns  $\tilde{\sigma} \leftarrow \tilde{A}(sk_i, m)$ .
3. Finally,  $U$  outputs a value  $b' \in \{0, 1\}$ ; we say that  $U$  wins iff  $b' = b$ .

We say that  $\mathcal{A}$  is  $u$ -users *publicly* undetectable w.r.t.  $\mathcal{SS}$  if in step 1. of the above game,  $U$  is only given the verification keys of the  $u$ -users.

### 7.3 Impersonation Relations

Theorem 8 formalizes the relations depicted in Fig. 5.

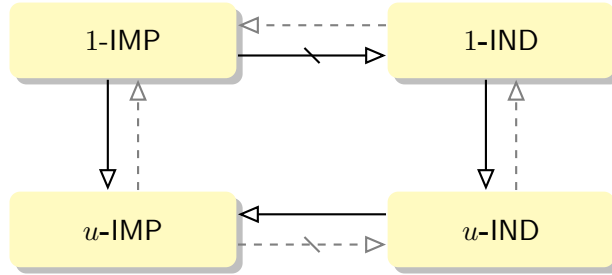


Figure 5: Diagram of the relationships between the subversion notions considered in this paper.  $X \rightarrow Y$  means that  $X$  implies  $Y$ ;  $X \nrightarrow Y$  indicates a separation between  $X$  and  $Y$ . The lighter arrows indicates trivial implications (or implications that follow from Theorem 8). Definition 17 is represented by  $u$ -IND and Definition 18 is represented by  $u$ -IMP.

**Theorem 8.** *The following relations hold: (i)  $1\text{-IND} \Rightarrow u\text{-IND}$ , (ii)  $1\text{-IMP} \nrightarrow 1\text{-IND}$ , (iii)  $u\text{-IND} \Rightarrow u\text{-IMP}$  for any EUF-CMA signature scheme, and (iv)  $1\text{-IMP} \Rightarrow u\text{-IMP}$ .*

*Proof.* (i)  $1\text{-IND} \Rightarrow u\text{-IND}$ . Towards contradiction, consider an adversary  $B$  that wins the game described in Definition 17. We build an adversary  $B'$  that (using  $B$ ) wins the game described in Definition 5.

Let  $\mathbf{G}$  be the game described in Definition 17. Consider the game  $\mathbf{G}_0$ , an identical copy of game  $\mathbf{G}$  when  $b = 0$ , and consider the game  $\mathbf{G}_1$  an identical copy of game  $\mathbf{G}$  when  $b = 1$ .

For an index  $k \in [0, u]$ , consider the hybrid game  $\mathbf{H}_k$  where each oracle corresponding to query  $(i, \tilde{A})$  such that  $i \leq k$  behaves as  $\tilde{A}(sk_i, \cdot)$  (i.e., as in game  $\mathbf{G}_0$ ), while all oracles corresponding to queries  $(i, \tilde{A})$  such that  $i > k$  behave as  $\text{Sign}(sk_i, \cdot)$  (i.e., as in game  $\mathbf{G}_1$ ). We note that  $\mathbf{H}_0 \equiv \mathbf{G}_1$  and  $\mathbf{H}_u \equiv \mathbf{G}_0$ . By assumption, we know that  $B$  can distinguish between the extreme hybrid games  $\mathbf{H}_0$  and  $\mathbf{H}_u$ . So there must exist a pair of hybrids  $\mathbf{H}_i, \mathbf{H}_{i-1}$  that  $B$  can distinguish with a non-negligible advantage. We can construct  $B'$  as follows.

Adversary  $B'$ :

1. Receive  $vk^*$  from the challenger and sample  $(vk_j, sk_j) \leftarrow \text{KGen}(1^\kappa)$  for all  $j \in [u] \setminus \{i\}$ . Define  $vk_i = vk^*$  and forward  $(vk_1, \dots, vk_u)$  to adversary  $B$ .
2. Upon input a query  $(j, \tilde{A})$  from  $B$ , behave as follows.
  - If  $j \leq i - 1$  answer all queries  $m \in \mathcal{M}$  as  $\tilde{\sigma} \leftarrow \tilde{A}(sk_j, m)$ ;
  - if  $j = i$  forward all queries  $m \in \mathcal{M}$  to the challenger;
  - if  $j \geq i$  answer all queries  $m \in \mathcal{M}$  as  $\sigma \leftarrow \text{Sign}(sk_j, m)$ .
3. Output whatever  $B$  outputs.

We observe that adversary  $B'$  simulates perfectly the distribution of the games  $\mathbf{H}_{i-1}$  (when  $b = 0$ ) and  $\mathbf{H}_i$  (when  $b = 1$ ). Since adversary  $B$  can distinguish this pair of hybrids with non-negligible probability it follows that adversary  $B'$  wins in the single-user game with the same probability.

(ii)  $1\text{-IMP} \not\Rightarrow 1\text{-IND}$ . We sketch a separation for the definitions. Consider  $\mathcal{SS}$  to be an EUF-CMA signature scheme with signature size  $\ell$  bits, and let  $\mathcal{A}$  be the class of SAs for  $\mathcal{SS}$  such that for all  $\tilde{A} \in \mathcal{A}$  the output of  $\tilde{A}$  is  $0^\ell$ . By  $\mathcal{SS}$  being EUF-CMA, adversary  $B$  has only a negligible probability of winning at the game described in Definition 18. However adversary  $B$  can clearly win the game described in Definition 17, because it is easy for  $B$  to distinguish between real signatures and subverted signatures.

(iii)  $u\text{-IND} \Rightarrow u\text{-IMP}$ . Consider  $\mathcal{SS}$  to be an EUF-CMA signature scheme and let  $\mathcal{A}$  be a class of SAs for  $\mathcal{SS}$ . The objective here is to show that if  $\mathcal{A}$  is  $u$ -users indistinguishable w.r.t continuous SAs (Definition 17) then  $\mathcal{A}$  is also  $u$ -users EUF-CMA w.r.t continuous SA (Definition 18). We sketch a proof by considering a modified game for Definition 18, where all oracles behave like the real signing oracle (one oracle for each signing key). Since the class  $\mathcal{A}$  is  $u$ -users indistinguishable we get that the advantage of any adversary in winning the game of Definition 18 is negligibly close to the advantage of winning in the modified game. However, since  $\mathcal{SS}$  is EUF-CMA no PPT adversary can win the modified game with non-negligible advantage, and so  $\mathcal{SS}$  satisfies  $u\text{-IMP}$ .

(iv)  $1\text{-IMP} \Rightarrow u\text{-IMP}$ . Towards contradiction, consider an adversary  $B$  that wins the game described in Definition 18. We build an adversary  $B'$  that (using  $B$ ) wins the game described in Definition 6.

Adversary  $B'$ :

1. Receive  $vk^*$  from the challenger, sample  $i^* \leftarrow_s \{1, \dots, u\}$  and  $(vk_i, sk_i) \leftarrow \text{KGen}(1^\kappa)$  for all  $i \in [u] \setminus \{i^*\}$ . Set  $vk_{i^*} := vk^*$  and forward  $(vk_1, \dots, vk_u)$  to adversary  $B$ .
2. Upon each query  $(i, m)$ , for  $i \in [u]$  and  $m \in \mathcal{M}$ : If  $i \neq i^*$  reply with  $\sigma = \text{Sign}(sk_i, m)$ , else forward the query to the challenger.
3. Upon each query  $(i, \tilde{A})$ , with  $i \in [u]$ , behave as follows.
  - For each  $m \in \mathcal{M}$  chosen by the adversary  $B$ , if  $i \neq i^*$  answer with  $\tilde{\sigma} = \tilde{A}(sk_j, m)$ , else forward the query to the challenger.
4. Eventually  $B$  outputs a forgery  $(i', m', \sigma')$ ; adversary  $B'$  outputs  $(m', \sigma')$  as its own forgery.

Adversary  $B'$  is successful if adversary  $B$  outputs a forgery for user  $i^*$ . Define  $E$ , to be the event that  $B'$  guesses correctly the index  $i' = i^*$ ; note that  $\mathbb{P}[E] = 1/u$ . Therefore adversary  $B'$  has a non-negligible probability of winning at the game described in Definition 6.  $\square$

## 7.4 Undetectability Relations

The following theorem (9) formalizes the relations depicted in Fig. 6.

**Theorem 9.** *The following relations hold. (i)  $u\text{-sUND} \Rightarrow u\text{-pUND}$ , (ii)  $1\text{-pUND} \not\Rightarrow 1\text{-sUND}$ , (iii)  $1\text{-pUND} \not\Rightarrow u\text{-pUND}$ , and (iv)  $1\text{-sUND} \not\Rightarrow u\text{-sUND}$ .*

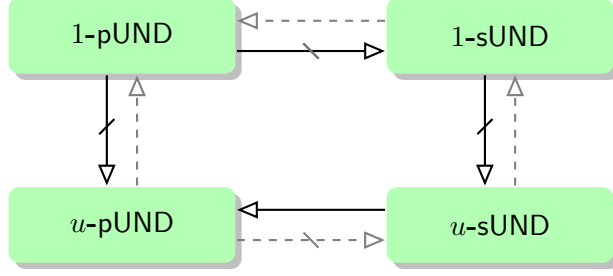


Figure 6: Diagram of the relationships between the undetectability notions considered in this paper.  $X \rightarrow Y$  means that  $X$  implies  $Y$ ;  $X \not\rightarrow Y$  indicates a separation between  $X$  and  $Y$ . The lighter arrows indicates trivial implications (or implications that follow from Theorem 9). Public undetectability (Definition 19) is represented by  $u$ -pUND and the secret undetectability (Definition 18) is represented by  $u$ -sUND.

*Proof.* (i)  $u$ -sUND  $\Rightarrow$   $u$ -pUND. Fix any class  $\mathcal{A}$  of SAs for  $\mathcal{SS}$ , and let  $C^*$  be the (efficient) challenger that exists by the assumption that  $\mathcal{A}$  is secretly undetectable. We claim that  $\mathcal{A}$  is also publicly undetectable for the same choice of the challenger  $C^*$ . Towards contradiction, consider a user  $U$  that wins the public undetectability game described in Definition 19 against  $C^*$ . We build a user  $U'$  (using  $U$ ) that wins the secret undetectability game described in Definition 19 against  $C^*$ .

User  $U'$ :

1. Receive  $(vk_i, sk_i) \leftarrow \text{KGen}(1^\kappa)$ , for  $i \in [u]$ , from the challenger  $C^*$  and forward it to user  $U$ .
2. User  $U$  asks polynomially many queries of the type  $(i, m)$  which are forwarded to the challenger  $C^*$ .
3. Output whatever  $U$  outputs.

We note that the simulation performed by user  $U'$  is perfect, therefore  $U'$  wins the secret undetectability game with the same probability that user  $U$  wins the public undetectability game.

(ii)  $1$ -pUND  $\not\Rightarrow$   $1$ -sUND. We sketch a separation between the definitions. Let  $\mathcal{SS}$  be a randomized signature scheme, and let  $\mathcal{SS}'$  be its derandomized implementation s.t.  $sk' := (sk, s)$ ,  $vk' := vk$ , and  $\sigma' := \text{Sign}(sk, m; r')$  with  $r' := F'_s(m)$  (for a PRF  $F$ ). We note that the only difference between  $\mathcal{SS}$  and  $\mathcal{SS}'$  is how the randomness  $r$  is computed for the signing algorithm. Let  $\mathcal{A}_{\text{bias}}^F$  be the class of SAs described in Fig. 1. By security of the PRF, and by Theorem 1, a user  $U$  has only a negligible probability of winning at the public undetectability game described in Definition 7. However, a user  $U$  playing the secret undetectability game knows  $sk' = (sk, s)$  and thus  $U$  can easily distinguish subverted and real signatures by simply re-computing  $r'$  and re-signing the input message; if both signatures match then with a high probability the target oracle is the real signing oracle. Notice that the last statement holds no matter how the subversion algorithm  $\tilde{A} \in \mathcal{A}_{\text{bias}}^F$  is selected by the challenger in the secret undetectability game.

(iii)  $1$ -pUND  $\not\Rightarrow$   $u$ -pUND. We sketch a separation for the definitions. Consider  $\mathcal{SS}$  to be a contrived signature scheme such that the signature of a message  $m \in \mathcal{M}$  is  $\sigma = \text{Sign}(sk, m; r) || r$ , where  $r \leftarrow_{\$} \{0, 1\}^\kappa$ . Let  $\mathcal{A} = \{\tilde{A}_{\tau, \bar{r}}\}_{\tau=0, \bar{r} \in \{0, 1\}^\kappa}$  to be class of SAs for  $\mathcal{SS}$  described next.

$\tilde{A}_{\tau, \bar{r}}(sk, m)$ :



1. If  $\tau = 0$  then let  $r := \bar{r}$ , else let  $r \leftarrow_{\$} \{0, 1\}^\kappa$ .
2. Output  $\sigma \leftarrow \text{Sign}(sk, m; r) || r$  and update  $\tau = \tau + 1$ .

Clearly, the class  $\mathcal{A}$  is publicly undetectable for a single user because the output of the subverted signature algorithm is indistinguishable from that of the real signing algorithm, even for the first query (when  $\tau = 0$ ). However, the class  $\mathcal{A}$  is clearly 2-users publicly *detectable* since (no matter the strategy of the challenger) it suffices to ask one query for each user and compare the last  $\kappa$  bits of the signatures to distinguish between real and subverted signatures.

(iv) 1-sUND  $\not\Rightarrow$   $u$ -sUND. We show a separation between the definitions. Consider  $\mathcal{SS}$  to be a randomized, coin-extractable signature scheme, with randomness size of  $\ell$ -bits, where  $\ell = |sk|$ , and  $\mathcal{A}_{\text{cext}}$  to be the class of SAs for  $\mathcal{SS}$  described in Fig. 2. We already showed in Theorem 2 that (for the challenger  $\mathbf{C}^*$  that chooses  $\tilde{\mathbf{A}}$  at random from  $\mathcal{A}_{\text{cext}}$ ) any PPT user  $\mathbf{U}$  playing the secret undetectability game described in Definition 7 has a negligible advantage. Now consider the same user  $\mathbf{U}$  playing the 2-users secret undetectability game described in Definition 19; user  $\mathbf{U}$  now has 2 key pairs that can be used to detect the attack in the following way.

User  $\mathbf{U}$ :

1. Receive  $(vk_i, sk_i) \leftarrow \text{KGen}(1^\kappa)$ , for  $i = 1, 2$ .
2. Fix a message  $\bar{m} \in \mathcal{M}$  and query  $(1, \bar{m})$  and  $(2, \bar{m})$  to the challenger, that replies with  $\sigma_1$  and  $\sigma_2$ .
3. Use  $\text{CExt}$  to extract the randomness from  $\sigma_1$  and  $\sigma_2$  to get  $r_1 \leftarrow \text{CExt}(vk_1, \bar{m}, \sigma_1)$  and  $r_2 \leftarrow \text{CExt}(vk_2, \bar{m}, \sigma_2)$ .
4. Compute  $sk_1 \oplus sk_2$  and return 0 iff the result equals  $r_1 \oplus r_2$ .

We note that the above detection strategy works regardless what strategy the challenger uses to select an algorithm from the class  $\mathcal{A}_{\text{cext}}$ . We conclude that user  $\mathbf{U}$  has an overwhelming probability of distinguishing between real and subverted signatures. □

## 8 Mounting Multi-User SAs

In this section we extend the attacks of Fig. 1 and Fig. 2 to the multi-user setting.

### 8.1 Attacking Coin-Injective Schemes (Multi-User Version)

The attack described in Fig. 1 can be extended to the multi-user setting with minor modifications. We create an SA class  $\mathcal{A}_{\text{bias}}^{F,u}$  from the class  $\mathcal{A}_{\text{bias}}^F$  of Fig. 1 by just appending the index  $j$ , that represents each user, to the function  $g(\cdot) = \text{Sign}(sk_j, m) || \tau || j$ , so that each application of the random function  $f(g(\cdot))$  remains independent.

The two lemmas below (Lemma 3 and Lemma 4) are needed for the proof of undetectability in the multi-user setting. The two lemmas combined roughly state that the statistical distance of a joint distribution of  $u$  random variables is *at most*  $u$  times the statistical distance of each pair of the random variables.

**Lemma 3** ([Rey11]). *Let  $\mathbf{X}$  and  $\mathbf{Y}$  be two random variables over some finite domain, and let  $G$  be a randomized function. Then  $\mathbb{SD}(G(\mathbf{X}), G(\mathbf{Y})) \leq \mathbb{SD}(\mathbf{X}, \mathbf{Y})$ .*

**Lemma 4.** Let  $\mathbf{X}$  and  $\mathbf{Y}$  be two random variables over some finite domain, and let  $(\mathbf{X}_1, \dots, \mathbf{X}_u)$  and  $(\mathbf{Y}_1, \dots, \mathbf{Y}_u)$  be  $u$  independent copies of the random variables  $\mathbf{X}$  and  $\mathbf{Y}$ . Then

$$\mathbb{SD}((\mathbf{X}_1, \dots, \mathbf{X}_u), (\mathbf{Y}_1, \dots, \mathbf{Y}_u)) \leq u \cdot \mathbb{SD}(\mathbf{X}, \mathbf{Y}).$$

*Proof.* We prove this lemma by induction. First we consider the basis case where  $i = 1$ , which trivially holds as  $\mathbb{SD}(\mathbf{X}_1, \mathbf{Y}_1) \leq \mathbb{SD}(\mathbf{X}, \mathbf{Y})$ .

For the induction step we define the random functions  $G_1(\cdot) = (\mathbf{X}_1, \dots, \mathbf{X}_{i-1}, \cdot)$  and  $G_2(\cdot) = (\cdot, \mathbf{Y}_i)$ . We assume that the statement holds up to  $i - 1$  random variables and then we proceed to show that it also holds for  $i$  random variables.

$$\begin{aligned} & \mathbb{SD}((\mathbf{X}_1, \dots, \mathbf{X}_i), (\mathbf{Y}_1, \dots, \mathbf{Y}_i)) \\ & \leq \mathbb{SD}((\mathbf{X}_1, \dots, \mathbf{X}_i), (\mathbf{X}_1, \dots, \mathbf{X}_{i-1}, \mathbf{Y}_i)) + \mathbb{SD}((\mathbf{X}_1, \dots, \mathbf{X}_{i-1}, \mathbf{Y}_i), (\mathbf{Y}_1, \dots, \mathbf{Y}_i)) \\ & = \mathbb{SD}(G_1(\mathbf{X}_i), G_1(\mathbf{Y}_i)) + \mathbb{SD}(G_2(\mathbf{X}_1, \dots, \mathbf{X}_{i-1}), G_2(\mathbf{Y}_1, \dots, \mathbf{Y}_{i-1})) \\ & \leq \mathbb{SD}(\mathbf{X}, \mathbf{Y}) + \mathbb{SD}((\mathbf{X}_1, \dots, \mathbf{X}_{i-1}), (\mathbf{Y}_1, \dots, \mathbf{Y}_{i-1})) \\ & \leq i \cdot \mathbb{SD}(\mathbf{X}, \mathbf{Y}), \end{aligned}$$

where the first inequality follows by the triangle inequality, the second inequality follows by Lemma 3, and the third inequality follows by the induction hypothesis.  $\square$

The theorem below quantifies the effectiveness of the attack of Fig. 1 in the multi-user setting.

**Theorem 10.** Let  $F : \{0, 1\}^\kappa \times \{0, 1\}^* \rightarrow \{0, 1\}$  be a secure PRF. For a randomized, coin-injective signature scheme  $\mathcal{SS}$  with randomness space of size  $\rho = |\mathcal{R}|$ , consider the class of SAs  $\mathcal{A}_{\text{bias}}^{F,u}$  described above. Then,

- (i)  $\mathcal{A}_{\text{bias}}^{F,u}$  is  $u$ -users secretly undetectable.
- (ii) Each  $\tilde{\mathbf{A}} \in \mathcal{A}_{\text{bias}}^{F,u}$  recovers the signing key of any of the users with probability at least  $(1 - (1/2 + \varepsilon_{\text{prf}})^\rho)^\ell$ , where  $\ell$  is the size of the signing key.

*Proof.* (i) Let  $\mathbf{G}$  be the game described in Definition 19. Consider the game  $\mathbf{G}_0$ , an identical copy of game  $\mathbf{G}$  when  $b = 0$ , and consider the game  $\mathbf{G}_1$ , an identical copy of game  $\mathbf{G}$  when  $b = 1$ . For the first part of the proof the objective is to show that  $\mathbf{G}_0 \approx \mathbf{G}_1$ .

Now consider game  $\mathbf{G}'_0$  an identical copy of game  $\mathbf{G}_0$  except that  $\mathbf{G}'_0$  utilizes the distribution from the random function  $f$  (analogous to Eq. (1) in the single user attack) instead of the distribution from the PRF  $F$  (analogous to Eq. (2) in the single user attack).

**Claim 9.**  $|\mathbb{P}[\text{U wins in } \mathbf{G}_0] - \mathbb{P}[\text{U wins in } \mathbf{G}'_0]| \leq \text{negl}(\kappa)$ .

The above claim follows by a standard reduction argument to the hardness of the PRF  $F$  to distinguishing games  $\mathbf{G}_0$  and  $\mathbf{G}'_0$ . The proof is similar to the one in Theorem 1 and is therefore omitted.

**Claim 10.**  $|\mathbb{P}[\text{U wins in } \mathbf{G}'_0] - \mathbb{P}[\text{U wins in } \mathbf{G}_1]| \leq \text{negl}(\kappa)$ .

*Proof.* Abusing notation, let us write  $\mathbf{G}'_0$  and  $\mathbf{G}_1$  for the distribution of the random variables corresponding to  $\text{U}$ 's view in games  $\mathbf{G}'_0$  and  $\mathbf{G}_1$  respectively. For an index  $i \in [0, q]$  consider the hybrid game  $\mathbf{H}_i$  that answers the first  $i$  signature queries as in game  $\mathbf{G}'_0$  while all the subsequent queries are answered as in  $\mathbf{G}_1$ . We note that  $\mathbf{H}_0 = \mathbf{G}_1$  and  $\mathbf{H}_q = \mathbf{G}'_0$ .

We claim that for all  $i \in [q]$ , we have  $\mathbb{SD}(\mathbf{H}_{i-1}, \mathbf{H}_i) \leq 2^{-(\rho+1)}$ . To see this, fix some  $i \in [q]$  and denote with  $\mathbf{R}_1, \dots, \mathbf{R}_u$  (resp.  $\tilde{\mathbf{R}}_1, \dots, \tilde{\mathbf{R}}_u$ ) the random variables defined by sampling an element from  $\mathcal{R}$  (resp.  $\tilde{\mathcal{R}}$ ) uniformly at random. Clearly,

$$\begin{aligned} \mathbb{SD}(\mathbf{H}_{i-1}, \mathbf{H}_i) &\leq \mathbb{SD}((\mathbf{R}_1, \dots, \mathbf{R}_u), (\tilde{\mathbf{R}}_1, \dots, \tilde{\mathbf{R}}_u)) \\ &\leq u \cdot \mathbb{SD}(\mathbf{R}, \tilde{\mathbf{R}}) \end{aligned} \tag{10}$$

$$= u \cdot 2^{-(\rho+1)}, \tag{11}$$

where Eq. (10) follows by Lemma 4 and Eq. (11) follows by Eq. (3). The claim now follows by the triangle inequality, as

$$\mathbb{SD}(\mathbf{G}_1, \mathbf{G}'_0) \leq \sum_{i=1}^q \mathbb{SD}(\mathbf{H}_{i-1}, \mathbf{H}_i) \leq qu \cdot 2^{-(\rho+1)}$$

and the last term becomes negligible for  $u, q = \text{poly}(\kappa)$  and for  $\rho$  large enough.  $\square$

The two claims above finish the proof of statement (i).

(ii) For the second part of the proof we proceed as in Theorem 1. We note that the specified class of SAs  $\mathcal{A}_{\text{bias}}^{F,u}$  maintains each application of the random function  $f$  independent by appending the index  $j$ , that represents each user, to the function  $g$ , obtaining  $g(\cdot) = \text{Sign}(sk_j, m) \parallel \tau \parallel j$ . The statement follows.  $\square$

## 8.2 Attacking Coin-Extractable Schemes (Multi-User Version)

The attack against coin-extractable schemes described in Fig. 2 becomes easily detectable in the presence of two or more users (see the proof of Theorem 9, item (iv)). An easy solution is to modify the SA class such that each algorithm in the class uses a different one-time pad key for each target user. We describe this class of SAs in Fig 7.

**Theorem 11.** *For a randomized,  $\nu_{\text{ext}}$ -coin-extractable, signature scheme  $\mathcal{SS}$  with randomness space  $\mathcal{R}$  of size  $\rho = 2^d$ , consider the class of SAs  $\mathcal{A}_{\text{cext}}^u$  described in Fig. 7. Then,*

(i)  $\mathcal{A}_{\text{cext}}^u$  is  $u$ -users secretly undetectable.

(ii) Each  $\tilde{\mathbf{A}} \in \mathcal{A}_{\text{cext}}^u$  recovers the signing key of any of the users with probability at least  $(1 - \nu_{\text{ext}})^{\ell/d}$ .

*Proof.* (i) Let  $\mathbf{G}$  be the game described in Definition 19, where the challenger first generates all key pairs  $(vk_i, sk_i)$  (for  $i \in [u]$ ) and afterwards selects the algorithm  $\tilde{\mathbf{A}} \leftarrow \mathcal{A}_{\text{cext}}^u$  such that  $\vec{vk} := (vk_1, \dots, vk_u)$ . Consider the game  $\mathbf{G}_0$ , an identical copy of game  $\mathbf{G}$  when  $b = 0$ , and consider the game  $\mathbf{G}_1$ , an identical copy of game  $\mathbf{G}$  when  $b = 1$ . For the first part of the proof the objective is to show that  $\mathbf{G}_0 \approx \mathbf{G}_1$ .

**Claim 11.**  $\mathbf{G}_0 \equiv \mathbf{G}_1$ .

*Proof.* Abusing notation, let us write  $\mathbf{G}_0$  and  $\mathbf{G}_1$  for the distribution of the random variables corresponding to  $\mathbf{U}$ 's view in games  $\mathbf{G}_0$  and  $\mathbf{G}_1$  respectively. For an index  $i \in [0, q]$  consider the hybrid game  $\mathbf{H}_i$  that answers the first  $i$  signature queries as in game  $\mathbf{G}_0$  while all the subsequent queries are answered as in  $\mathbf{G}_1$ . We note that  $\mathbf{H}_0 \equiv \mathbf{G}_1$  and  $\mathbf{H}_q \equiv \mathbf{G}_0$ .

### SA class $\mathcal{A}_{\text{cext}}^u$

Let  $\mathcal{SS} = (\text{KGen}, \text{Sign}, \text{Vrfy})$  be a coin-extractable randomized signature scheme with randomness space  $\mathcal{R}$  of size  $\rho = 2^d$ . The class  $\mathcal{A}_{\text{cext}}^u$  consists of a set of algorithms  $\{\tilde{\mathbf{A}}_{\vec{s}, \vec{vk}, \vec{\tau}}\}_{\vec{s} \in \{0,1\}^{\ell \cdot u}, \vec{vk} \in \mathcal{VK}^u, \vec{\tau} = 0^u}$ , where  $\ell = |sk|$ , and where each algorithm in the class behaves as follows:

$\tilde{\mathbf{A}}_{\vec{s}, \vec{vk}, \vec{\tau}}(sk_i, m)$ :

- Parse  $\vec{s}$  as  $(s_1, \dots, s_u)$ ,  $\vec{vk}$  as  $(vk_1, \dots, vk_u)$ , and  $\vec{\tau}$  as  $(\tau_1, \dots, \tau_u)$ .
- Find the index  $i$  such that  $\text{Vrfy}(vk_i, m, \text{Sign}(sk_i, m)) = 1$ .
- If  $\tau_i \geq \ell$  output a real signature  $\sigma \leftarrow \text{Sign}(sk_i, m)$ .
- Else,
  - For each value  $j \in [d]$  compute the biased random bit  $\tilde{r}[j] = s_i[\tau_i + j] \oplus sk_i[\tau_i + j]$ .
  - Return the signature  $\sigma := \text{Sign}(sk_i, m; \tilde{r})$ , and update the state  $\tau_i \leftarrow \tau_i + d$ .

**Extracting the signing key.** Given as input a vector of signatures  $\vec{\sigma} = (\sigma_1, \dots, \sigma_{\ell/d})$  of user  $i$ , represent the trapdoor  $s_i$  as  $\ell/d$  chunks of  $d$  bits  $s_i = \{s_{i,1}, \dots, s_{i,\ell/d}\}$ . For each signature  $\sigma_k \in \vec{\sigma}$  try to extract the  $d$ -bit chunk  $sk'_{i,k}$  of the signing key as follows.

- Extract the randomness from the  $k$ -th signature  $\tilde{r} \leftarrow \text{CExt}(vk_i, m_k, \sigma_k)$ .
- For each value  $j \in [d]$  compute the secret key bit  $sk'_{i,k}[j] = \tilde{r}[j] \oplus s_{i,k}[j]$ .

Return the signing key  $sk'_i := (sk'_{i,k}, \dots, sk'_{i,\ell/d})$ .

**Figure 7:** Attacking coin-extractable schemes in the multi-user setting

We claim that for all  $i \in [q]$ , we have  $\mathbf{H}_{i-1} \equiv \mathbf{H}_i$ . To see this, fix some  $i \in [q]$  and denote with  $\mathbf{R}_1, \dots, \mathbf{R}_u$  the random variables defined by sampling an element from  $\mathcal{R}$  uniformly at random and with  $\tilde{\mathbf{R}}_1, \dots, \tilde{\mathbf{R}}_u$  the random variables defined by sampling an element from the biased distribution  $\tilde{\mathcal{R}}$  also uniformly at random. It is easy to see that  $\mathbf{R}_i$  and  $\tilde{\mathbf{R}}_i$ , for  $i \in [q]$ , are identically distributed, as the biased distribution consists of a one-time pad encryption of (part of) the signing key with a uniform key (a different key for each user). The claim follows.  $\square$

(ii) For the second part of the proof we note that the attack of Fig. 7 successfully recovers the biased randomness  $\tilde{r}$  of each  $\sigma_i \in \{\sigma_1, \dots, \sigma_{\ell/d}\}$  and computes the chunk  $sk_{j,i}$  of the signing key of a user  $j$  with probability at least  $1 - \nu_{\text{ext}}$ . This gives a total probability of recovering an entire signing key of at least  $(1 - \nu_{\text{ext}})^{\ell/d}$ .  $\square$

## Acknowledgements

This work was supported by the European Commission H2020 program under the SUNFISH project, grant N.644666, and by the European Commission Directorate General Home Affairs, under the GAINS project, HOME/2013/CIPS/AG/4000005057.

We are grateful to Abhishek Jain for his insightful comments, suggestions, and contributions during the early stages of this work.

## References

- [ACF14] Michel Abdalla, Dario Catalano, and Dario Fiore. Verifiable random functions: Relations to identity-based key encapsulation and new constructions. *J. Cryptology*, 27(3):544–593, 2014.
- [ACM<sup>+</sup>14] Per Austrin, Kai-Min Chung, Mohammad Mahmoody, Rafael Pass, and Karn Seth. On the impossibility of cryptography with tamperable randomness. In *CRYPTO*, pages 462–479, 2014.
- [ADL14] Divesh Aggarwal, Yevgeniy Dodis, and Shachar Lovett. Non-malleable codes from additive combinatorics. In *STOC*, pages 774–783, 2014.
- [ADW09] Joël Alwen, Yevgeniy Dodis, and Daniel Wichs. Leakage-resilient public-key cryptography in the bounded-retrieval model. In *CRYPTO*, pages 36–54, 2009.
- [AFPW11] Martin R. Albrecht, Pooya Farshim, Kenneth G. Paterson, and Gaven J. Watson. On cipher-dependent related-key attacks in the ideal-cipher model. In *FSE*, pages 128–145, 2011.
- [AGM<sup>+</sup>15] Shashank Agrawal, Divya Gupta, Hemanta K. Maji, Omkant Pandey, and Manoj Prabhakaran. A rate-optimizing compiler for non-malleable codes against bit-wise tampering and permutations. In *TCC*, pages 375–397, 2015.
- [AHI11] Benny Applebaum, Danny Harnik, and Yuval Ishai. Semantic security under related-key attacks and applications. In *Innovations in Computer Science*, pages 45–60, 2011.
- [AMV15] Giuseppe Ateniese, Bernardo Magri, and Daniele Venturi. Subversion-resilient signature schemes. In *CCS*, pages 364–375, 2015.
- [BB08] Dan Boneh and Xavier Boyen. Short signatures without random oracles and the SDH assumption in bilinear groups. *J. Cryptology*, 21(2):149–177, 2008.
- [BBG13] James Ball, Julian Borger, and Glenn Greenwald. Revealed: how US and UK spy agencies defeat internet privacy and security. *Guardian Weekly*, September 2013.
- [BC10] Mihir Bellare and David Cash. Pseudorandom functions and permutations provably secure against related-key attacks. In *CRYPTO*, pages 666–684, 2010.
- [BCM11] Mihir Bellare, David Cash, and Rachel Miller. Cryptography secure against related-key attacks and tampering. In *ASIACRYPT*, pages 486–503, 2011.
- [BDI<sup>+</sup>99] Mike Burmester, Yvo Desmedt, Toshiya Itoh, Kouichi Sakurai, and Hiroki Shizuya. Divertible and subliminal-free zero-knowledge proofs for languages. *J. Cryptology*, 12(3):197–223, 1999.
- [Ber08] Daniel J. Bernstein. Proving tight security for Rabin-Williams signatures. In *EUROCRYPT*, pages 70–87, 2008.
- [BFGM01] Mihir Bellare, Marc Fischlin, Shafi Goldwasser, and Silvio Micali. Identification protocols secure against reset attacks. In *EUROCRYPT*, pages 495–511, 2001.
- [BJK15] Mihir Bellare, Joseph Jaeger, and Daniel Kane. Mass-surveillance without the state: Strongly undetectable algorithm-substitution attacks. In *CCS*, pages 1431–1440, 2015.

- [BK03] Mihir Bellare and Tadayoshi Kohno. A theoretical treatment of related-key attacks: RKA-PRPs, RKA-PRFs, and applications. In *EUROCRYPT*, pages 491–506, 2003.
- [BPR14] Mihir Bellare, Kenneth G. Paterson, and Phillip Rogaway. Security of symmetric encryption against mass surveillance. In *CRYPTO*, pages 1–19, 2014.
- [BR96] Mihir Bellare and Phillip Rogaway. The exact security of digital signatures—How to sign with RSA and Rabin. In *EUROCRYPT*, pages 399–416, 1996.
- [CL02] Jan Camenisch and Anna Lysyanskaya. A signature scheme with efficient protocols. In *SCN*, pages 268–289, 2002.
- [CL04] Jan Camenisch and Anna Lysyanskaya. Signature schemes and anonymous credentials from bilinear maps. In *CRYPTO*, pages 56–72, 2004.
- [Cor02] Jean-Sébastien Coron. Optimal security proofs for PSS and other signature schemes. In *EUROCRYPT*, pages 272–287, 2002.
- [CS00] Ronald Cramer and Victor Shoup. Signature schemes based on the strong RSA assumption. *ACM Trans. Inf. Syst. Secur.*, 3(3):161–185, 2000.
- [Des88a] Yvo Desmedt. Abuses in cryptography and how to fight them. In *CRYPTO*, pages 375–389, 1988.
- [Des88b] Yvo Desmedt. Subliminal-free authentication and signature (extended abstract). In *EUROCRYPT*, pages 23–33, 1988.
- [DFMV13] Ivan Damgård, Sebastian Faust, Pratyay Mukherjee, and Daniele Venturi. Bounded tamper resilience: How to go beyond the algebraic barrier. In *ASIACRYPT*, pages 140–160, 2013.
- [DFMV15] Ivan Damgård, Sebastian Faust, Pratyay Mukherjee, and Daniele Venturi. The chaining lemma and its application. In *ICITS*, pages 181–196, 2015.
- [DFP15] Jean Paul Degabriele, Pooya Farshim, and Bertram Poettering. A more cautious approach to security against mass surveillance. In *FSE*, 2015. To appear.
- [DGG<sup>+</sup>15] Yevgeniy Dodis, Chaya Ganesh, Alexander Golovnev, Ari Juels, and Thomas Ristenpart. A formal treatment of backdoored pseudorandom generators. In *EUROCRYPT*, pages 101–126, 2015.
- [DK12] Dana Dachman-Soled and Yael Tauman Kalai. Securing circuits against constant-rate tampering. In *CRYPTO*, pages 533–551, 2012.
- [DK14] Dana Dachman-Soled and Yael Tauman Kalai. Securing circuits and protocols against  $1/\text{poly}(k)$  tampering rate. In *TCC*, pages 540–565, 2014.
- [DLSZ15] Dana Dachman-Soled, Feng-Hao Liu, Elaine Shi, and Hong-Sheng Zhou. Locally decodable and updatable non-malleable codes and their applications. In *TCC*, pages 427–450, 2015.
- [DMS15] Yevgeniy Dodis, Ilya Mironov, and Noah Stephens-Davidowitz. Message transmission with reverse firewalls - secure communication on corrupted machines. *IACR Cryptology ePrint Archive*, 2015:548, 2015.

- [Dod03] Yevgeniy Dodis. Efficient construction of (distributed) verifiable random functions. In *PKC*, pages 1–17, 2003.
- [DPW10] Stefan Dziembowski, Krzysztof Pietrzak, and Daniel Wichs. Non-malleable codes. In *Innovations in Computer Science*, pages 434–452, 2010.
- [DY05] Yevgeniy Dodis and Aleksandr Yampolskiy. A verifiable random function with short proofs and keys. In *PKC*, pages 416–431, 2005.
- [FHN<sup>+</sup>12] Sebastian Faust, Carmit Hazay, Jesper Buus Nielsen, Peter Sebastian Nordholt, and Angela Zottarel. Signature schemes secure against hard-to-invert leakage. In *ASIACRYPT*, pages 98–115, 2012.
- [Fis03] Marc Fischlin. The Cramer-Shoup strong-RSA signature scheme revisited. In *PKC*, pages 116–129, 2003.
- [FMNV14] Sebastian Faust, Pratyay Mukherjee, Jesper Buus Nielsen, and Daniele Venturi. Continuous non-malleable codes. In *TCC*, pages 465–488, 2014.
- [FMNV15] Sebastian Faust, Pratyay Mukherjee, Jesper Buus Nielsen, and Daniele Venturi. A tamper and leakage resilient von Neumann architecture. In *PKC*, pages 579–603, 2015.
- [FMVW14] Sebastian Faust, Pratyay Mukherjee, Daniele Venturi, and Daniel Wichs. Efficient non-malleable codes and key-derivation for poly-size tampering circuits. In *EUROCRYPT*, pages 111–128, 2014.
- [FNV15] Antonio Faonio, Jesper Buus Nielsen, and Daniele Venturi. Mind your coins: Fully leakage-resilient signatures with graceful degradation. In *ICALP*, pages 456–468, 2015.
- [FPV11] Sebastian Faust, Krzysztof Pietrzak, and Daniele Venturi. Tamper-proof circuits: How to trade leakage for tamper-resilience. In *ICALP*, pages 391–402, 2011.
- [Fry00] Niklas Frykholm. Countermeasures against buffer overflow attacks. Technical report, RSA Data Security, Inc., November 2000.
- [GHR99] Rosario Gennaro, Shai Halevi, and Tal Rabin. Secure hash-and-sign signatures without the random oracle. In *EUROCRYPT*, pages 123–139, 1999.
- [GIP<sup>+</sup>14] Daniel Genkin, Yuval Ishai, Manoj Prabhakaran, Amit Sahai, and Eran Tromer. Circuits resilient to additive attacks with applications to secure computation. In *STOC*, pages 495–504, 2014.
- [GL10] David Goldenberg and Moses Liskov. On related-secret pseudorandomness. In *TCC*, pages 255–272, 2010.
- [GLM<sup>+</sup>04] Rosario Gennaro, Anna Lysyanskaya, Tal Malkin, Silvio Micali, and Tal Rabin. Algorithmic tamper-proof (ATP) security: Theoretical foundations for security against hardware tampering. In *TCC*, pages 258–277, 2004.
- [GOR11] Vipul Goyal, Adam O’Neill, and Vanishree Rao. Correlated-input secure hash functions. In *TCC*, pages 182–200, 2011.

- [Gre14] Glenn Greenwald. No place to hide: Edward Snowden, the NSA, and the U.S. surveillance state. *Metropolitan Books*, May 2014.
- [HJK12] Dennis Hofheinz, Tibor Jager, and Edward Knapp. Waters signatures with optimal security reduction. In *PKC*, pages 66–83, 2012.
- [HK12] Dennis Hofheinz and Eike Kiltz. Programmable hash functions and their applications. *J. Cryptology*, 25(3):484–527, 2012.
- [HW09a] Susan Hohenberger and Brent Waters. Realizing hash-and-sign signatures under standard assumptions. In *EUROCRYPT*, pages 333–350, 2009.
- [HW09b] Susan Hohenberger and Brent Waters. Short and stateless signatures from the RSA assumption. In *CRYPTO*, pages 654–670, 2009.
- [IPSW06] Yuval Ishai, Manoj Prabhakaran, Amit Sahai, and David Wagner. Private circuits II: keeping secrets in tamperable circuits. In *EUROCRYPT*, pages 308–327, 2006.
- [Jag15] Tibor Jager. Verifiable random functions from weaker assumptions. In *TCC*, pages 121–143, 2015.
- [JW15] Zahra Jafargholi and Daniel Wichs. Tamper detection and continuous non-malleable codes. In *TCC*, pages 451–480, 2015.
- [KKS11] Yael Tauman Kalai, Bhavana Kanukurthi, and Amit Sahai. Cryptography with tamperable and leaky memory. In *CRYPTO*, pages 373–390, 2011.
- [KT13] Aggelos Kiayias and Yiannis Tselekounis. Tamper resilient circuits: The adversary at the gates. In *ASIACRYPT*, pages 161–180, 2013.
- [KW03] Jonathan Katz and Nan Wang. Efficiency improvements for signature schemes with tight security reductions. In *ACM CCS*, pages 155–164, 2003.
- [LL12] Feng-Hao Liu and Anna Lysyanskaya. Tamper and leakage resilience in the split-state model. In *CRYPTO*, pages 517–532, 2012.
- [Luc04] Stefan Lucks. Ciphers secure against related-key attacks. In *FSE*, pages 359–370, 2004.
- [Lys02] Anna Lysyanskaya. Unique signatures and verifiable random functions from the DH-DDH separation. In *CRYPTO*, pages 597–612, 2002.
- [MRV99] Silvio Micali, Michael O. Rabin, and Salil P. Vadhan. Verifiable random functions. In *FOCS*, pages 120–130, 1999.
- [MS15] Ilya Mironov and Noah Stephens-Davidowitz. Cryptographic reverse firewalls. In *EUROCRYPT*, pages 657–686, 2015.
- [NIS07] NIST (National Institute of Standards and Technology). Special Publication 800-90: Recommendation for random number generation using deterministic random bit generators, March 2007.
- [NPS01] David Naccache, David Pointcheval, and Jacques Stern. Twin signatures: an alternative to the hash-and-sign paradigm. In *ACM CCS*, pages 20–27, 2001.



- [NVZ14] Jesper Buus Nielsen, Daniele Venturi, and Angela Zottarel. Leakage-resilient signatures with graceful degradation. In *PKC*, pages 362–379, 2014.
- [One96] Aleph One. Smashing the stack for fun and profit. *Phrack Magazine*, 7(49):File 14, 1996.
- [PB04] Jonathan D. Pincus and Brandon Baker. Beyond stack smashing: Recent advances in exploiting buffer overruns. *IEEE Security & Privacy*, 2(4):20–27, 2004.
- [PLS13] Nicole Perlroth, Jeff Larson, and Scott Shane. N.S.A. able to foil basic safeguards of privacy on web. *The New York Times*, September 2013.
- [Rey11] Leo Reyzin. Lecture notes: Extractors and the leftover hash lemma, March 2011.
- [RTYZ15] Alexander Russell, Qiang Tang, Moti Yung, and Hong-Sheng Zhou. Cliptography: Clipping the power of kleptographic attacks. *IACR Cryptology ePrint Archive*, 2015:695, 2015.
- [Sch12] Sven Schäge. Strong security from probabilistic signature schemes. In *PKC*, pages 84–101, 2012.
- [SFKR15] Bruce Schneier, Matthew Fredrikson, Tadayoshi Kohno, and Thomas Ristenpart. Surreptitiously weakening cryptographic systems. *IACR Cryptology ePrint Archive*, 2015:97, 2015.
- [Sim83] Gustavus J. Simmons. The prisoners’ problem and the subliminal channel. In *CRYPTO*, pages 51–67, 1983.
- [Sim84] Gustavus J. Simmons. The subliminal channel and digital signature. In *EUROCRYPT*, pages 364–378, 1984.
- [VV83] Umesh V. Vazirani and Vijay V. Vazirani. Trapdoor pseudo-random number generators, with applications to protocol design. In *FOCS*, pages 23–30, 1983.
- [Wat05] Brent Waters. Efficient identity-based encryption without random oracles. In *EUROCRYPT*, pages 114–127, 2005.
- [Wee12] Hoeteck Wee. Public key encryption against related key attacks. In *PKC*, pages 262–279, 2012.
- [YY96] Adam L. Young and Moti Yung. The dark side of “black-box” cryptography, or: Should we trust Capstone? In *CRYPTO*, pages 89–103, 1996.
- [YY97] Adam L. Young and Moti Yung. Kleptography: Using cryptography against cryptography. In *EUROCRYPT*, pages 62–74, 1997.
- [YY04] Adam L. Young and Moti Yung. *Malicious Cryptography: Exposing Cryptovirology*. John Wiley & Sons, Inc., first edition, 2004.