

Algorithm Substitution Attacks on Public Functions

MIHIR BELLARE¹

DOREEN RIEPEL²

LAURA SHEA³

June 14, 2024

Abstract

We study the possibility of algorithm substitution attacks (ASAs) on functions with no secret-key material, such as hash functions, and verification algorithms of signature schemes and proof systems. We consider big-brother’s goal to be three-fold: It desires utility (it can break the scheme), exclusivity (nobody else can) and undetectability (outsiders can’t detect its presence). We start with a general setting in which big-brother is aiming to subvert an arbitrary public function. We give, in this setting, strong definitions for the three goals. We then present a general construction of an ASA, and prove that it meets these definitions. We use this to derive, as applications, ASAs on hash functions, signature schemes and proof systems. As a further application of the first two, we give an ASA on X.509 certificates. While ASAs were traditionally confined to exfiltrating secret keys, our work shows that they are possible and effective at subverting public functions where there are no keys to exfiltrate. Our constructions serve to help defenders and developers identify potential attacks by illustrating how they might be built.

¹ Department of Computer Science & Engineering, University of California San Diego, 9500 Gilman Drive, La Jolla, California 92093, USA. Email: mbellare@ucsd.edu. URL: <http://cseweb.ucsd.edu/~mihir/>. Supported in part by NSF grant CNS-2154272 and KACST.

² Department of Computer Science & Engineering, University of California San Diego, 9500 Gilman Drive, La Jolla, California 92093, USA. Email: driepel@ucsd.edu. Supported in part by KACST.

³ Department of Computer Science & Engineering, University of California San Diego, 9500 Gilman Drive, La Jolla, California 92093, USA. Email: lmshea@ucsd.edu. Supported by NSF grants CNS-2048563, CNS-1513671 and CNS-2154272.

Contents

1	Introduction	2
2	Preliminaries	6
3	Algorithm substitution attacks on public functions	7
4	ASA construction	10
5	ASAs on hash functions	13
6	ASAs on proof system verification	17
7	ASAs on signature verification	20
8	Application: Forged certificates	23
	References	25
A	Proof of Theorem 6.2	30
B	Proof of Theorem 7.2	32
C	Forged certificate embedding	34

1 Introduction

The Snowden revelations lead researchers to ask how cryptography might be subverted. Algorithm Substitution Attacks (ASAs) [10, 66, 67] emerged as one answer. The first formalism and attacks, by BPR [10], put forth the following template. There is a prescribed cryptographic algorithm A that accesses a secret key K . The adversary (called big-brother in this setting) substitutes A by subverted code \tilde{A} that aims to undetectably exfiltrate K . For this purpose, \tilde{A} and big-brother may share a symmetric key. Such ASAs have now been given or considered for many primitives including symmetric encryption [3, 9, 10, 26, 41], signature schemes [5, 20, 62] and beyond [14, 14, 19, 21, 38, 44, 53–55, 63]. However, the goal has always remained to exfiltrate a *secret key*.

We initiate work in a new and different direction inspired by work on backdooring of machine-learning algorithms [40]. We consider the possibility of ASAs on public functions, such as hash functions, verification algorithms of signature schemes, or verification algorithms of proof systems like SNARGs or SNARKs. In these cases, there is no secret key to exfiltrate, so one has to ask what an ASA would want to do, and what properties big-brother would like it to have. We answer these questions with new definitions. These are given in a general setting where the goal is an ASA on an arbitrary public function. In the same general setting, we then give a construction of an ASA that we prove meets our definitions. We then apply this to obtain ASAs on the specific primitives mentioned above.

THE SETTING. Honest function f is sampled from a prescribed family F , as $f \leftarrow {}^s F$. Big-brother generates a substitution function \tilde{f} together with an associated exploit function (also called a backdoor) e , via $(\tilde{f}, e) \leftarrow {}^s \tilde{F}(f)$, where \tilde{F} is an algorithm of big-brother’s own devising. It now arranges that a user’s code implementing f is substituted with code implementing \tilde{f} . This substitution can take place by a variety of means and is outside our scope. Applications that expected to use f are now (unknowingly) using \tilde{f} instead. In this setting, we consider big-brother to have three goals, as follows:

- **Utility:** Big-brother, through knowledge of the exploit function e , now wants to violate security of the \tilde{f} -using application.
- **Undetectability:** A tester with black-box access to either f or \tilde{f} should not be able to tell which of the two it is.
- **Exclusivity:** Others, meaning entities knowing \tilde{f} but not e , should *not* be able to violate security of the application even when it uses \tilde{f} .

Utility of course represents the main intent of big-brother in creating the ASA. Undetectability is a core ASA requirement [10]. Exclusivity captures that a big-brother government or Intelligence Agency wants to ensure that other parties cannot violate security of the application. Similarly, a big-brother corporation may want to ensure that its competitors and the public cannot violate security of the application.

MOTIVATING APPLICATIONS. Before formalizing the above requirements and giving our general construction, we discuss some motivating applications.

The first is that f is a collision-resistant hash function. The canonical application is hashing the content (data) in a digital certificate before signing it to create the certificate itself. The utility of interest to big-brother would be that, when \tilde{f} is used in place of f , big-brother can forge certificates. Exclusivity means others cannot forge certificates even under \tilde{f} . The ASAs we will give improve the backdoored hash functions of [34] which had lower utility, lower exclusivity and no undetectability requirement.

The second is signatures. Prior work on ASAs [5] aims to substitute the signing algorithm and

exfiltrate the secret signing key, and the ASAs only succeed for schemes which have randomized signing with sufficient entropy. Instead, we let f be the verification algorithm, viewed as embedding a target, public verification key vk . The substitution continues to accept signatures under vk , but utility allows big-brother to have it also accept other signatures that big-brother devises using e . Yet, exclusivity ensures that an outside entity not knowing either e or the secret signing key underlying vk be unable to create new signatures that \tilde{f} accepts.

The third is proof systems. These are extensively used in blockchains and cryptocurrencies. Our ASAs will show how to subvert soundness. We let f be the proof-verification algorithm, possibly embedding a CRS. The substitution \tilde{f} accepts valid proofs under the CRS but also accepts proofs, devised by big-brother using e , of false statements. Exclusivity precludes others from violating soundness even relative to \tilde{f} .

THE GENERAL FRAMEWORK. Rather than treat the above three applications separately, we will derive them via a general framework that we now discuss. Here F is an arbitrary function family from which the honest function $f \leftarrow_s F$ is sampled, and big-brother generates the substitution function \tilde{f} together with an exploit function e via $(\tilde{f}, e) \leftarrow_s \tilde{F}(f)$.

- **Utility:** Utility seems at first glance to be very application dependent, and we have to ask what it may mean in general. We parameterize the requirement by a predicate $P(\cdot, \cdot)$ called the constraint. Then we ask that exploit function e , given a string u called the constraint-parameter, and given a target output y , returns an x such that (1) $f(x) = y$ and (2) $P(x, u) = \text{true}$. Intuitively, big-brother can create an input x , satisfying a desired constraint involving a u of its choice, such that x maps to an output y of its choice under \tilde{f} . Different applications will make different choices of P .
- **Undetectability:** We consider a game that picks a random challenge bit b and gives the detector an oracle that on input x returns $f(x)$ if $b = 1$ and $\tilde{f}(x)$ otherwise. Asking that the detector have negligible advantage in guessing b formalizes black-box undetectability, a strong requirement in line with prior work [10, 40].
- **Exclusivity:** We consider an adversary that is given the descriptions of f and \tilde{f} and oracle access to the exploit function e , and ask that it cannot come up with an input x at which f and \tilde{f} differ, except trivially, meaning through use of its oracle.

Our first result is Theorem 3.1, saying that exclusivity, in our strong formulation, actually implies undetectability. We define the latter separately because it is a core ASA requirement insisted on by prior work [10], but for our ASAs, we will prove exclusivity and then conclude undetectability via Theorem 3.1.

OUR GENERAL ASA CONSTRUCTION. Is it possible to build an ASA \tilde{F} on an arbitrary function family F that meets the three conditions above? We show, through construction, that the answer is “yes.”

To expand on this, first note that one cannot hope to achieve this for all constraint predicates P . It is, for example, impossible for the predicate $P(x, u)$ that returns true iff $x = u$; intuitively, one needs some “room” in x to exploit. We show how to build an ASA \tilde{F} for any constraint predicate satisfying a certain condition, that we define and call *embeddability*. The class of predicates meeting this condition is large and includes in particular ones allowing our above-discussed applications.

Our construction is a transform **ASA** that takes (1) the target function family F (2) a signature scheme S and (3) an *embedding function* Emb , for the constraint predicate P , that is compatible with S, F , as we will define in Section 4. (Embeddability, for now, just asks that such an embedding exists.) It returns an ASA $\tilde{F} = \mathbf{ASA}[F, S, \text{Emb}]$ built from these three components. Proposition 4.1 establishes utility with respect to P , assuming correctness of the embedding and signature scheme.

Theorem 4.2 establishes exclusivity assuming strong unforgeability of the signature scheme S . Undetectability follows from Theorem 3.1. The design of **ASA** extends ideas of [40], who used a signature scheme to backdoor a machine learning model.

ASAS ON HASH FUNCTIONS. Towards our first application, to collision-resistant hash functions and certificate forgeries, we start with some background. While ASAs on hash functions have not to our knowledge been explicitly discussed before, Fischlin, Janson and Mazaheri (FJM) [34] considered backdooring hash functions, which is similar. Let H be the target set of hash functions from which, via $h \leftarrow_s H$, one generates an honest hash function $h : \{0, 1\}^* \rightarrow \{0, 1\}^\ell$ that is assumed collision-resistant. A backdooring of H can be seen, like an ASA, as specified by an algorithm \tilde{H} that takes an instance h of H and via $(\tilde{h}, e) \leftarrow_s \tilde{H}(h)$ generates a substitution function $\tilde{h} : \{0, 1\}^* \rightarrow \{0, 1\}^\ell$ as well as a backdoor (or exploitation) function e . FJM [34] have two requirements that in our language represent utility and exclusivity. FJM-utility asks that knowledge of e allows big-brother to violate collision resistance (cr) of \tilde{h} , meaning find distinct x_1, x_2 such that $\tilde{h}(x_1) = \tilde{h}(x_2)$. FJM-exclusivity asks that an adversary given \tilde{h} but not e cannot violate cr of \tilde{h} . They give no undetectability condition.

FJM [34] build a backdooring \tilde{H} of H satisfying their two conditions. Roughly, \tilde{h} embeds the image \bar{y} of a random point \bar{x} under a one-way function g , and \tilde{h} behaves anomalously if (and only if) its input maps to \bar{y} under g . Other works also satisfy this basic notion of a backdoored hash function: Albertini, Aumasson, Eichlseder, Mendel and Schl affer (AAEMS) [2] give a backdoored version of SHA1. The VSH (Very Smooth Hash) algorithm of Contini, Lenstra and Steinfeld [22] achieves FJM-utility and FJM-exclusivity when viewed as a backdoored hash function.

This notion of utility is however limited; in its quest for subversion, big-brother wants significantly more. We illustrate with two examples:

Forgery of TLS certificates. Forgery of TLS certificates, to allow impersonation of a target website, is a commonly considered big-brother goal [60, 61, 69]. Hashing is used here by the CA to compress the certificate data before it signs to create the certificate. An arbitrary hash collision will not help towards a forgery; big-brother needs not only to find a preimage x of a target point y (the hash of a legitimate certificate) but to be able to embed in x information of its choice, such as a public key for which it knows the secret key.

Breaking password-based authentication. Big-brother may be interested in subverting password-based authentication, where the server holds a hash of the user password and salt. Big-brother needs to find a preimage of this same hash that contains a password that it knows and the given salt as a suffix.

We note that FJM [34] and AAEMS [2] suggest that utility be the ability to violate preimage resistance but their formulations still do not suffice for the above tasks because they do not allow big-brother to embed information of its choice in the preimage. Meanwhile, FJM-exclusivity too is limited, as collisions created by big-brother using e can reveal e , making it easy for outside observers to find further collisions.

In Section 5, we give an ASA on hash functions that satisfies our strong utility and exclusivity requirements. In particular, the ASA permits the two exploits discussed above, meaning TLS certificate forgery and breaking of password-based authentication. We obtain our hash-function ASA by applying our above-discussed general **ASA** and results about it. To obtain the two exploits, we make particular choices of the constraint predicates that parameterize our definition. Section 8 describes the X.509 forgery in detail. We note that this application relies crucially on our strong, predicate-parameterized definition of utility.

ASAS ON PROOF SYSTEMS. As a second example, we illustrate ASAs on proof systems, targeting

the substitution of the verification algorithm. We use the notation $(p, v) \leftarrow^s \text{PS}$ to denote the generation of a proving and verification algorithm which hardcode the CRS. Our ASA then targets the family of verification functions \mathbf{V} , where our goal is to replace $v \in \mathbf{V}$ with a substitution function \tilde{v} , and to produce an exploit function e via an ASA generation algorithm $\tilde{\mathbf{V}}$.

The property we presume an ASA on verification to target is soundness. Using e , we want to be able to create an accepting proof for any statement ϕ which is not in the language described by the proof system. This means that our utility definition asks to effectively violate soundness. Further, we ask that soundness is maintained for any other party not having access to e , even if it has observed forged proofs for statements of its choice. This is captured by our exclusivity definition which provides oracle access to e . In order to apply our general **ASA** transform, we view a proof forgery as finding a preimage of 1. (Recall that a verification function outputs a bit where 1 indicates success.) In line with the general setting, we define a corresponding predicate P_{PS} along with an embedding function. We include these results in Section 6.

ASAS ON SIGNATURE SCHEMES. Finally, we turn to signature schemes. Similar to proof systems, we want to substitute the verification algorithm. However in this setting, different users create their own signing-verification key pair, hence it is desired that an exploitation algorithm works for not only one but multiple users. Again we consider a family of functions \mathbf{V} , where a function $v_i \in \mathbf{V}$ hardcodes the verification key for user i , and we want to construct an ASA $\tilde{\mathbf{V}}$ which takes v_1, \dots, v_n and generates substitution functions $\tilde{v}_1, \dots, \tilde{v}_n$ and a single exploit function e for all n users.

Utility might ask that possession of e allows creation of some forgery, meaning a message m and a valid signature for it under the verification function \tilde{v}_i of some user. We ask for a stronger notion of utility, namely that given any user i , and *any target message* m , the exploit algorithm allows creation of a signature for m relative to \tilde{v}_i . Our exclusivity definition is similarly strong, asking that unforgeability of signatures still holds even with access to an oracle for producing signatures using the exploit algorithm for any user. We observe that our general ASA transform needs to be adapted so that signatures produced by the exploit algorithm are tied to the user. In Section 7, we therefore give a modified multi-user ASA transform \mathbf{ASA}_n and show that it indeed achieves our specified goals above.

As a concrete application, we consider again the certificate forgery example. Using our signature ASA, big-brother can use the exploitation function to create certificates with validating signatures. This is discussed in detail in Section 8.

RELATED WORK. “Subversion” is a recurring concern in cryptography. Attacks have been observed in a variety of settings, including parameter generation in Dual EC [15, 57], malicious code changes in Linux [32], and governmental exceptional access [1]. To situate our investigation of ASAs on public functions, it is useful to consider different categories in this area of subverted cryptography.

In a first category, code can be maliciously modified from its algorithmic specification. This has been studied as algorithm substitution attacks [10] and as kleptography [66–68]. Here, an adversary’s goal is to both modify an algorithm such that it exfiltrates secret information, and to keep this modification undetected. We call our considered attacks ASAs because they involve a code substitution step, and a second exploitation step. However, we define different goals than exfiltration and undetectability. ASAs and kleptography have been studied for symmetric encryption [3, 9, 10, 26], KEMs [19, 44, 53], signatures [5, 62], and protocols [14, 14, 21, 38, 63]. Defenses have been studied from the perspective of preventing exfiltration [30, 50] or other “subversion-resistant” notions [4, 13, 54, 55]. These include more fine-grained online/offline detectability notions and other modes of computation.

In another category of subversion, violations occur through authoritarian means. That is, an authority can overcome usual security guarantees, but this extra power is well known to the

public. Anamorphic cryptography [7, 47, 52] and earlier work on subliminal channels [43, 58, 59] have proposed defenses to this type of subversion, which is different than the ASA model.

A final category, somewhat more related to ASAs, is maliciously designed algorithms or parameters. These have been studied for PRGs [27, 29], NIZKs [8], PKE [6], and hash functions [2, 34]. Unlike an ASA, algorithms are assumed to be implemented honestly, and code can be inspected. Nonetheless, some of the goals and techniques are similar. In relevant sections, we will compare related work on particular primitives in more detail.

The 2022 work of Goldwasser, Kim, Vaikuntanathan and Zamir (GKVZ) studies the possibility of inserting undetectable backdoors in a machine learning model [40]. We leverage their techniques in Section 4. In GKVZ, a strongly unforgeable signature triggers alternate execution in a model, modifying classifier output when a signature is correctly parsed and verified. The non-replicability condition of GKVZ offers an oracle providing backdoored model inputs, which is similar to an additional oracle that our exclusivity game provides. We note that the predicates and embeddings we additionally formalize in Section 4 generalize the parsing of GKVZ, and that our techniques for “general public functions” seem to be quite applicable beyond cryptographic functions.

As a real-world motivation to study ASAs, the **xz** backdoor was discovered on March 29, 2024 [35]. Current understanding of its cryptographic portion [64] shows interesting similarities with our ASAs, such as the embedding of a signature and attacker-chosen data in a certificate which triggers alternate execution during certificate validation. The discovery of the **xz** backdoor shows that ASAs targeting high levels of utility are a realistic possibility, and motivates research, such as ours, on this topic.

2 Preliminaries

NOTATION AND TERMINOLOGY. By ε we denote the empty string. By $|Z|$ we denote the length of a string Z . By $x \parallel y$ we denote the concatenation of strings x, y . If Z is a string, we let $Z[a..b]$ be the substring of Z between indices a and b , inclusive, or ε if $b < a$. If S is a finite set, then $|S|$ denotes its size. We say that a set S is *length-closed* if, for any $x \in S$ it is the case that $\{0, 1\}^{|x|} \subseteq S$. (This will be a requirement for certain spaces.)

If \mathcal{X} is a finite set, we let $x \leftarrow^* \mathcal{X}$ denote picking an element of \mathcal{X} uniformly at random and assigning it to x . If A is an algorithm, we let $y \leftarrow A[\mathcal{O}_1, \dots](x_1, \dots; r)$ denote running A on inputs x_1, \dots and coins r with oracle access to \mathcal{O}_1, \dots , and assigning the output to y . We let $y \leftarrow^* A[\mathcal{O}_1, \dots](x_1, \dots)$ be the result of picking r at random and computing $y \leftarrow A[\mathcal{O}_1, \dots](x_1, \dots; r)$. We let $\text{OUT}(A[\mathcal{O}_1, \dots](x_1, \dots))$ denote the set of all possible outputs of A when invoked with inputs x_1, \dots and oracles \mathcal{O}_1, \dots . Algorithms are randomized unless otherwise indicated. Running time is worst case, which for an algorithm with access to oracles means across all possible replies from the oracles. The abbreviation “p.p.t.” denotes “probabilistic polynomial time.”

An adversary is an algorithm. We use \perp (bot) as a special symbol to denote rejection, and it is assumed to not be in $\{0, 1\}^*$. The image of a function $f : \mathcal{D} \rightarrow \mathcal{R}$ is the set $\text{Im}(f) = \{f(x) : x \in \mathcal{D}\} \subseteq \mathcal{R}$. We may interchangeably refer to the boolean `false` and integer 0, or to the boolean `true` and integer 1.

GAMES. We use the code-based game-playing framework of BR [11]. A game G starts with an optional `INIT` procedure, followed by a non-negative number of additional procedures called oracles, and ends with a `FIN` procedure. Execution of adversary A with game G begins by running `INIT` (if present) to produce `input` $\leftarrow^* \text{INIT}$. A is then given `input` and is run with query access to the game oracles. When A terminates with some `output`, execution of game G ends by returning `FIN(output)`.

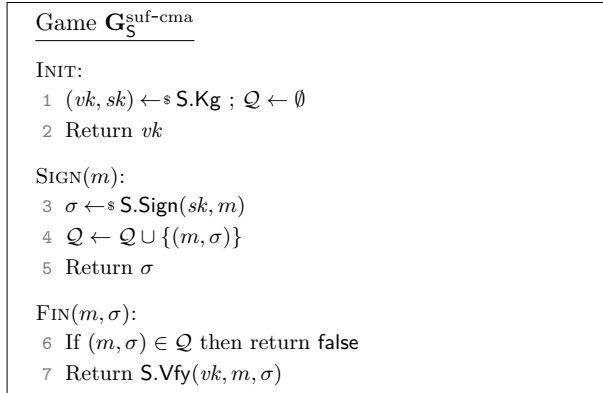


Figure 1: Strong unforgeability for a signature scheme \mathcal{S} .

By $\Pr[G(A)]$ we denote the probability that the execution of game G with adversary A results in $\text{FIN}(\text{output})$ being the boolean true.

Different games may have procedures (oracles) with the same names. If we need to disambiguate, we may write $G.O$ to refer to oracle O of game G . In games, integer variables, set variables, boolean variables and string variables are assumed initialized, respectively, to 0, the empty set \emptyset , the boolean false and \perp . Tables are initialized with all entries being \perp . Games may occasionally **Require**: some condition, which means that all adversaries must obey this condition. This is used to rule out trivial wins.

UNFORGEABILITY OF SIGNATURES. A signature scheme \mathcal{S} specifies algorithms $\mathcal{S}.\text{Kg}$, $\mathcal{S}.\text{Sign}$, $\mathcal{S}.\text{Vfy}$, key spaces $\mathcal{S}.\text{VK}$, $\mathcal{S}.\text{SK}$, and signature length $\mathcal{S}.\text{sl}$. Key generation $\mathcal{S}.\text{Kg}$ produces a verification key $vk \in \mathcal{S}.\text{VK}$ and signing key $sk \in \mathcal{S}.\text{SK}$ via $(vk, sk) \leftarrow \mathcal{S}.\text{Kg}$. Signing takes as input a signing key $sk \in \mathcal{S}.\text{SK}$ and message $m \in \{0, 1\}^*$ to return a signature $\sigma \in \{0, 1\}^{\mathcal{S}.\text{sl}}$ via $\sigma \leftarrow \mathcal{S}.\text{Sign}(sk, m)$, where $\mathcal{S}.\text{sl} \in \mathbb{N}$ is a constant signature length. Deterministic algorithm $\mathcal{S}.\text{Vfy}$ takes as input a verification key $vk \in \mathcal{S}.\text{VK}$, message $m \in \{0, 1\}^*$, and signature $\sigma \in \{0, 1\}^{\mathcal{S}.\text{sl}}$ to return a bit d via $d \leftarrow \mathcal{S}.\text{Vfy}(vk, m, \sigma)$.

Correctness of scheme \mathcal{S} asks that for all $(vk, sk) \in \text{OUT}(\mathcal{S}.\text{Kg})$, for all $m \in \{0, 1\}^*$, it holds that $\mathcal{S}.\text{Vfy}(vk, m, \mathcal{S}.\text{Sign}(sk, m)) = 1$.

The security notion that we will make use of is strong unforgeability. This is captured by game $\mathbf{G}_S^{\text{suf-cma}}$ of Figure 1. If A is an adversary, we let $\mathbf{Adv}_S^{\text{suf-cma}}(A) = \Pr[\mathbf{G}_S^{\text{suf-cma}}(A)]$ be its suf-cma advantage. Strongly unforgeable signatures have been constructed based on bilinear CDH [17], strong RSA [24, 37], and generally from one-way functions [39, Section 6.5].

3 Algorithm substitution attacks on public functions

We begin with new definitions for algorithm substitution attacks on arbitrary functions. In Section 4 we will provide a construction satisfying our notions, and in the remainder of the paper we extend both our definitions and construction to more specific settings. As our first task, we introduce a generic ASA syntax.

SYNTAX. Let F be a function family; a member function $f : \{0, 1\}^* \rightarrow \{0, 1\}^{F.\text{ol}}$ is selected via $f \leftarrow \mathcal{S}.\text{F}$. An ASA aims to substitute f with \tilde{f} which has certain malicious behavior on particular inputs, while ensuring that the ability to find such inputs is exclusive to the attacker. Concretely,

Game $\mathbf{G}_{F,\tilde{F},P}^{\text{exc}}$	Game $\mathbf{G}_{F,\tilde{F},P}^{\text{det}}$
INIT: 1 $f \leftarrow \mathfrak{s} F ; (\tilde{f}, e) \leftarrow \mathfrak{s} \tilde{F}(f)$ 2 $\mathcal{X} \leftarrow \emptyset$ 3 Return (f, \tilde{f}) GETPMG(u, y): 4 $x \leftarrow \mathfrak{s} e(u, y) ; \mathcal{X} \leftarrow \mathcal{X} \cup \{x\}$ 5 Return x FIN(x^*): 6 Return $(x^* \notin \mathcal{X}) \wedge (f(x^*) \neq \tilde{f}(x^*))$	INIT: 1 $f \leftarrow \mathfrak{s} F ; (\tilde{f}, e) \leftarrow \mathfrak{s} \tilde{F}(f)$ 2 $b \leftarrow \mathfrak{s} \{0, 1\}$ 3 Return ε EVAL(x): 4 $y_0 \leftarrow f(x) ; y_1 \leftarrow \tilde{f}(x)$ 5 Return y_b FIN(b'): 6 Return $(b = b')$

Figure 2: Exclusivity (left) and detectability (right) of an ASA \tilde{F} on public function f . While detectability is black-box, exclusivity returns the functions (f, \tilde{f}) to the adversary.

we have an algorithm \tilde{F} which generates $(\tilde{f}, e) \leftarrow \mathfrak{s} \tilde{F}(f)$, where \tilde{f} is the *substitution* and e is the *exploitation* algorithm. Before elaborating on this ability, let us clarify the ASA model.

ASA MODEL FOR PUBLIC FUNCTIONS. An attacker who mounts an ASA proceeds in two steps, substitution and exploitation. First, they generate algorithms $(\tilde{f}, e) \leftarrow \mathfrak{s} \tilde{F}(f)$ and replace a user’s implementation of f with one of \tilde{f} . The exploitation algorithm e remains secret and is retained by the attacker. Second, in the exploitation step, the attacker may use e to find a preimage x of target point y . The user, who has \tilde{f} on their device, now computes $\tilde{f}(x) = y$. The interface of this exploitation algorithm is broad, applying to settings where (1) it is reasonable that the attacker chooses inputs x to send to the user, and (2) it is useful to an attacker to find preimages of target points (where the target point could even be 1 or true).

UTILITY. The *utility* capabilities of \tilde{F} are captured with respect to a constraint predicate P . Let $P : \{0, 1\}^* \times \{0, 1\}^* \rightarrow \{\text{true}, \text{false}\}$ be a predicate. We say that \tilde{F} achieves utility relative to P , if for every every $f \in F$, every constraint-parameter $u \in \{0, 1\}^*$ and every $y \in \{0, 1\}^{F, \text{ol}}$, if $(\tilde{f}, e) \leftarrow \mathfrak{s} \tilde{F}(f)$ and if $x \leftarrow \mathfrak{s} e(u, y)$, then we have (1) $\tilde{f}(x) = y$ and (2) $P(x, u) = 1$. In other words, e allows one to compute a preimage of any target output, where the preimage also satisfies the constraint predicate. In the following section we discuss predicates in more detail. At a high level, this utility definition relative to a predicate allows one to specify fine-grained notions of successful attacks, beyond only finding *some* preimage.

EXCLUSIVITY. Effectiveness of an ASA calls for exclusivity as well as utility. *Exclusivity* requires that, even after seeing inputs x such that $f(x) \neq \tilde{f}(x)$, it remains hard to find another nontrivial input x^* on which f and \tilde{f} differ. This is captured by the exclusivity game $\mathbf{G}_{F,\tilde{F},P}^{\text{exc}}$ of Figure 2. If A

is an adversary, we let $\mathbf{Adv}_{F,\tilde{F},P}^{\text{exc}}(A) = \Pr \left[\mathbf{G}_{F,\tilde{F},P}^{\text{exc}}(A) \right]$ be its exc advantage.

We will see exclusivity applied to more specific settings in later sections, but for now, we remark that an ASA is only useful if the substituted algorithm \tilde{f} differs from f on some inputs (else f could be attacked directly). Exclusivity asks that is hard for anyone other than the ASA to find these inputs which produce different behavior. The game $\mathbf{G}_{F,\tilde{F},P}^{\text{exc}}$ provides an additional oracle GETPMG for seeing these “trigger inputs” and an adversary A wins if it produces another such input x^* which was not generated by GETPMG directly.

UNDETECTABILITY. We define (black-box) undetectability as a distinguishing game which is de-

<p>Games $G_0, \boxed{G_1}$</p> <p>INIT:</p> <ol style="list-style-type: none"> 1 $f \leftarrow \mathfrak{F} ; (\tilde{f}, e) \leftarrow \mathfrak{F}(f)$ 2 $b \leftarrow \mathfrak{s} \{0, 1\}$ 3 Return ε <p>EVAL(x):</p> <ol style="list-style-type: none"> 4 $y_0 \leftarrow f(x) ; y_1 \leftarrow \tilde{f}(x)$ 5 If $y_0 \neq y_1$ then bad \leftarrow true ; Return \perp 6 Return y_b <p>FIN(b'):</p> <ol style="list-style-type: none"> 7 Return $(b = b')$ 	<p>Adversary $A'(f, \tilde{f})$:</p> <ol style="list-style-type: none"> 1 $b \leftarrow \mathfrak{s} \{0, 1\}$ 2 $b' \leftarrow A[\text{EVAL}]()$ <p>Oracle EVAL(x):</p> <ol style="list-style-type: none"> 3 $y_0 \leftarrow f(x) ; y_1 \leftarrow \tilde{f}(x)$ 4 If $y_0 \neq y_1$ then FIN(x) 5 Return y_b
--	---

Figure 3: Games G_0, G_1 (left) and adversary A' (right) for the proof of Theorem 3.1. G_1 contains the boxed code and G_0 does not.

scribed on the right side of Figure 2. If A is an adversary, we let $\mathbf{Adv}_{\mathfrak{F}, \tilde{\mathfrak{F}}, \mathfrak{P}}^{\det}(A) = 2 \cdot \Pr \left[\mathbf{G}_{\mathfrak{F}, \tilde{\mathfrak{F}}, \mathfrak{P}}^{\det}(A) \right] - 1$ be its det advantage. The following statement shows that exclusivity implies undetectability. For this reason, we will continue our focus on exclusivity, as it is stronger. Nonetheless, black-box undetectability remains a standard notion in works on ASAs.

Theorem 3.1 *Let \mathfrak{F} be a family of functions and $\tilde{\mathfrak{F}}$ an ASA on \mathfrak{F} relative to a predicate \mathfrak{P} . Given an adversary A against the undetectability of $\tilde{\mathfrak{F}}$ we can build an adversary A' such that*

$$\mathbf{Adv}_{\mathfrak{F}, \tilde{\mathfrak{F}}, \mathfrak{P}}^{\det}(A) \leq 2 \cdot \mathbf{Adv}_{\mathfrak{F}, \tilde{\mathfrak{F}}, \mathfrak{P}}^{\text{exc}}(A'). \quad (1)$$

A' makes no GETPMG queries. The running time of A' is close to that of A .

Proof of Theorem 3.1: Consider game G_0 of Figure 3 which is exactly the det game, except that it additionally checks whether y_0 equals y_1 for queries to EVAL and sets flag **bad** if they are not the same. The output, however, is not modified, so

$$\mathbf{Adv}_{\mathfrak{F}, \tilde{\mathfrak{F}}, \mathfrak{P}}^{\det}(A) = 2 \cdot \Pr[G_0(A)] - 1.$$

We now turn to G_1 . Games G_0, G_1 are identical-until-**bad**, so by the Fundamental Lemma of Game Playing [11] we have

$$\begin{aligned} \Pr[G_0(A)] &= \Pr[G_1(A)] + (\Pr[G_0(A)] - \Pr[G_1(A)]) \\ &\leq \Pr[G_1(A)] + \Pr[G_1(A) \text{ sets bad}]. \end{aligned}$$

We construct adversary A' on the right side of Figure 3. A' gets as input functions f and \tilde{f} , and can simulate the detectability game in a straightforward way by picking its own challenge bit. Since A' does not ask any GETPMG queries, then if **bad** is set in G_1 , A' has found a winning output in the exclusivity game. We have

$$\Pr[G_1(A) \text{ sets bad}] \leq \mathbf{Adv}_{\mathfrak{F}, \tilde{\mathfrak{F}}, \mathfrak{P}}^{\text{exc}}(A').$$

Finally note that $\Pr[G_1(A)] = 1/2$ since the outputs of EVAL are independent of the challenge bit. Collecting the probabilities proves Eq. (1). ■

<p><u>$P_{\text{pfx}}^n(x, u)$:</u></p> <ol style="list-style-type: none"> 1 If $(x \neq u + n)$ then return false 2 Return $(x[1.. u] = u)$ <p><u>$\text{Emb}_{\text{pfx}}^n(z, u)$:</u></p> <ol style="list-style-type: none"> 3 Return $u \parallel z$ <p><u>$\text{Emblnv}_{\text{pfx}}^n(x)$:</u></p> <ol style="list-style-type: none"> 4 If $(x < n)$ then return \perp 5 $z \leftarrow x[(x - n + 1).. x]$ 6 $u \leftarrow x[1..(x - n)]$ 7 Return (z, u) 	<p><u>$P_{\text{sfx}}^n(x, u)$:</u></p> <ol style="list-style-type: none"> 1 If $(x \neq u + n)$ then return false 2 Return $(x[(x - u + 1).. x] = u)$ <p><u>$\text{Emb}_{\text{sfx}}^n(z, u)$:</u></p> <ol style="list-style-type: none"> 3 Return $z \parallel u$ <p><u>$\text{Emblnv}_{\text{sfx}}^n(x)$:</u></p> <ol style="list-style-type: none"> 4 If $(x < n)$ then return \perp 5 $z \leftarrow x[1..n]$ 6 $u \leftarrow x[(n + 1).. x]$ 7 Return (z, u)
--	--

Figure 4: Practical examples of predicates and embedding functions: a prefix embedding (left) and suffix embedding (right). The embedding space is $\text{Emb}_{\text{pfx}}^n \cdot \text{ES} = \text{Emb}_{\text{sfx}}^n \cdot \text{ES} = \{0, 1\}^n$ for a fixed n .

4 ASA construction

PREDICATES AND EMBEDDINGS. In order to proceed to our construction, we first turn to some details of realizing constraint predicates in constructions. For this, we introduce a message embedding function Emb for predicate P . This is a map $\text{Emb} : \text{Emb} \cdot \text{ES} \times \{0, 1\}^* \rightarrow \{0, 1\}^*$, where $\text{Emb} \cdot \text{ES}$ is a set (the “embedding space”) that must be specified and depends on the intended predicate. There is also an inverse $\text{Emb}^{-1} : \{0, 1\}^* \rightarrow (\text{Emb} \cdot \text{ES} \times \{0, 1\}^*) \cup \{\perp\}$ such that (1) for all $(z, u) \in \text{Emb} \cdot \text{ES} \times \{0, 1\}^*$, if $x \leftarrow \text{Emb}(z, u)$ then $P(x, u) = 1$ and $\text{Emb}^{-1}(x) = (z, u)$, and (2) $\text{Emb}^{-1}(x) = \perp$ for all $x \notin \text{Im}(\text{Emb})$. We say that Emb is a *correct embedding function* for predicate P if these two properties are satisfied.

In other words, Emb is a bijection from $(\text{Emb} \cdot \text{ES} \times \{0, 1\}^*)$ to $\text{Im}(\text{Emb})$ with inverse Emb^{-1} .

Two illustrative examples are prefix and suffix embeddings. Suppose one wants to find a preimage x of function output y , with the constraint that x begins with prefix u , or ends with suffix u . For a prefix embedding, let $n \in \mathbb{N}$ and $\text{Emb}_{\text{pfx}}^n \cdot \text{ES} = \{0, 1\}^n$. Then the predicate P_{pfx}^n and embedding function $\text{Emb}_{\text{pfx}}^n$ are given on the left side of Figure 4. Similarly, for suffixes, let $\text{Emb}_{\text{sfx}}^n \cdot \text{ES} = \{0, 1\}^n$, with the predicate P_{sfx}^n and embedding function $\text{Emb}_{\text{sfx}}^n$ on the right side of Figure 4. The choice of n here is important: it will matter for the feasibility of designing an effective ASA in some settings.

A variety of other predicates and embedding functions may be desirable in practice. Predicate $P(x, u)$ could capture whether x is a valid X.509 certificate containing information u ; this is considered in more detail in Section 8. A predicate could capture whether x can be parsed as human-readable text or otherwise does not “look suspicious.” Increasingly useful predicates will come with implementation challenges beyond mounting an ASA, but the ones described above are already potent.

WARMUP CONSTRUCTION. We begin with a basic construction of an ASA. This is essentially a generalization of the construction of FJM [34, Section 7.1] which targets backdoored hash functions. More specifically, we build an ASA algorithm \tilde{F}_0 for a function instance f using a one-way function $g : \{0, 1\}^k \rightarrow \{0, 1\}^\ell$. \tilde{F}_0 is specified on the left of Figure 5 and it achieves some degree of effectiveness. To an entity with e_{bd} , finding a preimage of target y is simple: $bd \parallel y$ is a preimage because the trigger $(g(bd') = t)$ passes in f_t , which then returns y . In terms of exclusivity, t is public but finding bd remains hard to find assuming g is one-way. However, this does not hold if a preimage produced by e_{bd} is observed. Thus our notion of exclusivity is not satisfied, but black-box undetectability is.

<p><u>$\tilde{F}_0(f)$:</u></p> <ol style="list-style-type: none"> 1 $bd \leftarrow_s \{0, 1\}^k$; $t \leftarrow g(bd)$ 2 Define as below: <ul style="list-style-type: none"> $\tilde{f}_t : \{0, 1\}^* \rightarrow \{0, 1\}^{\text{F.ol}}$ $e_{bd} : \{0, 1\}^* \times \{0, 1\}^{\text{F.ol}} \rightarrow \{0, 1\}^*$ 3 Return (\tilde{f}_t, e_{bd}) <p><u>$\tilde{f}_t(x)$:</u></p> <ol style="list-style-type: none"> 4 If $(x = k + \text{F.ol})$ then 5 $bd' \leftarrow x[1..k]$ 6 If $(g(bd') = t)$ then return $x[(k+1).. x]$ 7 Else return $f(x)$ <p><u>$e_{bd}(u, y)$:</u></p> <ol style="list-style-type: none"> 8 // In the warmup, u is ignored 9 Require: $y \in \{0, 1\}^{\text{F.ol}}$ 10 $x \leftarrow bd \parallel y$ 11 Return x 	<p><u>$\tilde{F}(f)$:</u></p> <ol style="list-style-type: none"> 1 $(vk, sk) \leftarrow_s \text{S.Kg}$ 2 Define as below: <ul style="list-style-type: none"> $\tilde{f}_{vk} : \{0, 1\}^* \rightarrow \{0, 1\}^{\text{F.ol}}$ $e_{sk} : \{0, 1\}^* \times \{0, 1\}^{\text{F.ol}} \rightarrow \{0, 1\}^*$ 3 Return (\tilde{f}_{vk}, e_{sk}) <p><u>$\tilde{f}_{vk}(x)$:</u></p> <ol style="list-style-type: none"> 4 $w \leftarrow \text{Emb}^{-1}(x)$ 5 If $(w = \perp)$ then return $f(x)$ 6 $((y \parallel \sigma), u) \leftarrow w$ 7 If $\text{S.Vfy}(vk, (y, u), \sigma)$ then return y 8 Else return $f(x)$ <p><u>$e_{sk}(u, y)$:</u></p> <ol style="list-style-type: none"> 9 Require: $y \in \{0, 1\}^{\text{F.ol}}$ 10 $\sigma \leftarrow_s \text{S.Sign}(sk, (y, u))$ 11 $x \leftarrow \text{Emb}((y \parallel \sigma), u)$ 12 Return x
--	---

Figure 5: **Left:** Warmup construction of an ASA \tilde{F}_0 using a one-way function g . **Right:** Construction of an ASA \tilde{F} , relative to predicate P with associated embedding function Emb , using a signature scheme S .

We omit a formal analysis as we will next consider a construction which does meet all of our notions. In particular, we would like to find preimages satisfying a predicate P , and achieve exclusivity in addition to undetectability.

CONSTRUCTION. Let us now turn to meeting these requirements. Let S be an suf-cma signature scheme and let F be a family of functions. We say that message embedding function $\text{Emb} : \text{Emb.ES} \times \{0, 1\}^* \rightarrow \{0, 1\}^*$ is *compatible* with S, F if $\text{Emb.ES} = \{0, 1\}^{\text{F.ol} + \text{S.sl}}$. That is, the embedding information consists of an output of a function from F and a signature. Our transform **ASA** associates to S, F , and an Emb compatible to S, F an ASA algorithm $\tilde{F} = \text{ASA}[F, S, \text{Emb}]$ which is defined on the right side of Figure 5.

Note that we assume a correct embedding function Emb for predicate P ; these have been given for common predicates in Figure 4, but it may not be the case that every predicate has a correct embedding function, or that every embedding function is compatible with S, F .

In the remainder of this section, we show that \tilde{F} produced by transform $\text{ASA}[F, S, \text{Emb}]$ is in fact effective, achieving utility as long as S, Emb are correct (Proposition 4.1) and achieving exclusivity as long as S is suf-cma (Theorem 4.2). The bottom line of this result is that effective ASAs *are* possible to construct from standard building blocks, for useful constraints.

Proposition 4.1 *Let S be a signature scheme, F a family of functions, and Emb an embedding function for predicate P which is compatible with S, F . Let $\tilde{F} = \text{ASA}[F, S, \text{Emb}]$. If S is a correct signature scheme and Emb is a correct embedding function for predicate P , then \tilde{F} achieves utility for P .*

Proof of Proposition 4.1: Consider any $f \in F$, $(\tilde{f}_{vk}, e_{sk}) \in \text{OUT}(\tilde{F}(f))$, $u \in \{0, 1\}^*$, and $y \in \{0, 1\}^{\text{F.ol}}$. The function e_{sk} , on inputs u and y , returns $x \leftarrow \text{Emb}((y \parallel \sigma), u)$ where $\sigma \leftarrow_s \text{S.Sign}(sk,$

<p>Games $G_0, \boxed{G_1}$</p> <p>INIT:</p> <ol style="list-style-type: none"> 1 $f \leftarrow \text{\\$ } F ; (vk, sk) \leftarrow \text{\\$ } S.Kg$ 2 Define \tilde{f}_{vk} and e_{sk} as in Fig. 5 (using vk, sk chosen above) 3 $\mathcal{X} \leftarrow \emptyset$ 4 Return (f, \tilde{f}_{vk}) <p>GETPMG(u, y):</p> <ol style="list-style-type: none"> 5 $\sigma \leftarrow \text{\\$ } S.Sign(sk, (y, u)) ; x \leftarrow Emb((y \parallel \sigma), u)$ 6 $\mathcal{X} \leftarrow \mathcal{X} \cup \{x\}$ 7 Return x <p>FIN(x^*):</p> <ol style="list-style-type: none"> 8 If $x^* \in \mathcal{X}$ then return false 9 $y' \leftarrow f(x^*) ; w \leftarrow Emb^{-1}(x^*)$ 10 If $(w \neq \perp)$ then 11 $((y \parallel \sigma), u) \leftarrow w$ 12 If $S.Vfy(vk, (y, u), \sigma)$ then 13 $bad \leftarrow true ; \boxed{\text{Return false}}$ 14 Else: $y \leftarrow y'$ 15 Else: $y \leftarrow y'$ 16 Return $(y' \neq y)$

Figure 6: Games G_0, G_1 for the proof of Theorem 4.2. G_1 contains the boxed code and G_0 does not.

(y, u)). Property (1) of utility requires that $\tilde{f}(x) = y$. Let us consider $\tilde{f}_{vk}(x)$ of our construction. On lines 4,5, since $x \in \text{Im}(Emb)$, $w \neq \perp$. If Emb, Emb^{-1} satisfy our notion of a correct embedding, w is recovered as $((y \parallel \sigma), u)$. That is, $Emb^{-1}(Emb((y \parallel \sigma), u)) = ((y \parallel \sigma), u)$. Next, the signature verification on line 7 runs $S.Vfy(vk, (y, u), \sigma)$ where $\sigma \leftarrow \text{\$ } S.Sign(sk, (y, u))$. This passes as long as S is a correct signature scheme, and $\tilde{f}_{vk}(x)$ thus returns y on line 7.

Property (2) of utility asks that $P(x, u) = 1$. This is proven by line 11, where $x \leftarrow Emb((y \parallel \sigma), u)$. Correctness of the embedding function Emb for predicate P implies that $P(x, u) = 1$. ■

Theorem 4.2 *Let S be a signature scheme, F a family of functions, and Emb a correct embedding function for predicate P which is compatible with S, F . Let $\tilde{F} = \mathbf{ASA}[F, S, Emb]$. Given an adversary A against the exclusivity of \tilde{F} we can build an adversary A_S such that*

$$\mathbf{Adv}_{F, \tilde{F}, P}^{\text{exc}}(A) \leq \mathbf{Adv}_S^{\text{suf-cma}}(A_S). \quad (2)$$

If A makes q GETPMG queries, then A_S makes q SIGN queries. The running time of A_S is close to that of A .

Proof of Theorem 4.2: Consider game G_0 of Figure 6. We claim that

$$\mathbf{Adv}_{F, \tilde{F}, P}^{\text{exc}}(A) = \Pr[G_0(A)]. \quad (3)$$

To justify Eq. (3), we claim that the $\text{FIN}(x^*)$ return value is the same in G_0 as it is in $\mathbf{G}_{F, \tilde{F}, P}^{\text{exc}}$. (The INIT and GETPMG oracles are identical, instantiated with scheme \tilde{F} .) To begin with, the check made in line 8 is identical. Now consider y and y' of G_0 . G_0 sets $y' = f(x^*)$. Further, lines 9-15 of G_0 correspond to lines 4-7 of \tilde{f}_{vk} in Figure 5. That is, $y = y' = f(x^*)$ if there is no valid parsing of

<p>Adversary $A_S(vk)$:</p> <ol style="list-style-type: none"> 1 $f \leftarrow \mathfrak{s} \mathbf{F}$; $\mathcal{X} \leftarrow \emptyset$ 2 Define \tilde{f}_{vk} as in Fig. 5 (using vk provided as input) 3 $x^* \leftarrow A[\text{GETPMG}_S](f, \tilde{f}_{vk})$ 4 $w \leftarrow \text{Emb}^{-1}(x^*)$ 5 If $(w \neq \perp)$ then 6 $((y \parallel \sigma), u) \leftarrow w$ 7 If $\mathbf{S.Vfy}(vk, (y, u), \sigma)$ then 8 Return $((y, u), \sigma)$ <p>Oracle $\text{GETPMG}_S(u, y)$:</p> <ol style="list-style-type: none"> 9 $\sigma \leftarrow \text{SIGN}((y, u))$ 10 $x \leftarrow \text{Emb}((y \parallel \sigma), u)$; $\mathcal{X} \leftarrow \mathcal{X} \cup \{x\}$ 11 Return x

Figure 7: Adversary A_S for the proof of Theorem 4.2.

w nor verified signature. Otherwise, y is as output by Emb^{-1} . Hence, the final check in line 16 of G_0 is identical to that of the exclusivity game. This proves Eq. (3).

We now turn to G_1 . Games G_0, G_1 are identical-until-bad, so by the Fundamental Lemma of Game Playing [11] we have

$$\begin{aligned} \Pr[G_0(A)] &= \Pr[G_1(A)] + (\Pr[G_0(A)] - \Pr[G_1(A)]) \\ &\leq \Pr[G_1(A)] + \Pr[G_1(A) \text{ sets bad}]. \end{aligned}$$

It is easy to see that $\Pr[G_1(A)] = 0$. This is because either the game will return **false** in line 13, or it will set y to y' in which case line 16 will return **false**.

It remains to bound the difference between G_0 and G_1 . For this, we construct an adversary A_S and claim that

$$\Pr[G_1(A) \text{ sets bad}] \leq \mathbf{Adv}_S^{\text{suf-cma}}(A_S). \quad (4)$$

This will complete the proof of Eq. (2) and the theorem statement.

We now explain adversary A_S , which is in game $\mathbf{G}_S^{\text{suf-cma}}$ and runs A as specified in Figure 7. Note that A_S can define \tilde{f}_{vk} using vk which it receives as input and f which it chooses by itself. Further, it simulates A 's GETPMG oracle using its own SIGN oracle.

If **bad** is set in G_1 , then the message-signature pair $((y, u), \sigma)$ that was parsed from $w \leftarrow \text{Emb}^{-1}(x^*)$ has passed verification. We also know that **bad** is only set when $x^* \notin \mathcal{X}$ due to the check in line 8 of G_1 . Since the embedding is correct and deterministic, $((y, u), \sigma)$ was not used in the simulation of the signing oracle and is a winning output in game $\mathbf{G}_S^{\text{suf-cma}}(A_S)$. This completes the proof of Eq. (4).

Note that A_S makes one SIGN query for each of A 's GETPMG_S queries, proving the running time in the theorem statement. ■

5 ASAs on hash functions

We now turn to applying our general definitions and transform to hash functions, but we begin with some related work on hash function subversion. For one, “backdoored hash functions” consist

of malicious hash function designs or parameter selection, and have been considered by [2, 34]. An explicit construction of maliciously designed SHA1 parameters was given by [2]. Implementations of hash functions have been studied from a kleptographic perspective as subverted random oracles [56] and from a proof-techniques perspective as programmable hash functions [42]. From the kleptographic perspective, [12] studied defense techniques for subverted hash functions, for which the model excluded domains of $\{0, 1\}^*$ (which we consider) among other definitional differences. “Subverted hash functions” also brings to mind asymmetric notions including chameleon hash functions [28, 46], trapdoor hash functions [31], and provably secure hash functions (with a trapdoor) [22]. However, these works build hash functions with trapdoor properties for constructive applications, and the interface differs from standard hash functions.

SYNTAX. We now proceed to ASAs on hash functions. A hash function family H is a set of hash functions; one generates hash function $h : \{0, 1\}^* \rightarrow \{0, 1\}^\ell$ via $h \leftarrow_s H$. Following the syntax from Section 3, an ASA on H is specified by an algorithm \tilde{H} that takes an instance h of H and via $(\tilde{h}, e) \leftarrow_s \tilde{H}(h)$ generates a substitution function $\tilde{h} : \{0, 1\}^* \rightarrow \{0, 1\}^\ell$ as well as an exploitation function e .

The following definitions of utility and exclusivity strengthen comparable notions from prior work including [34], as discussed in the Introduction. In particular, we consider collision resistance in the presence of an exploit-finding oracle, and we define utility as the ability to find highly structured preimages, not only *some* preimage.

UTILITY. Let $P : \{0, 1\}^* \times \{0, 1\}^* \rightarrow \{\text{true}, \text{false}\}$ be a constraint predicate. We say that \tilde{H} achieves *utility* relative to P if for every constraint-parameter $u \in \{0, 1\}^*$ and every $y \in \{0, 1\}^\ell$, if $(\tilde{h}, e) \in \text{OUT}(\tilde{H})$ and if $x \leftarrow_s e(u, y)$, then we have (1) $\tilde{h}(x) = y$ and (2) $P(x, u) = 1$. In other words, e allows one to compute a preimage of any target hash, where the preimage also satisfies the constraint predicate.

Notably, this is the same as our notion of utility for general public functions; our first application is hash functions because it most immediately matches the general case.

EXCLUSIVITY. When it comes to exclusivity, hash functions do introduce an additional property: collision resistance is now desired. We define exclusivity for \tilde{H} via game $\mathbf{G}_{H, \tilde{H}, P}^{\text{cfe}}$ of Figure 8, where “cfe” denotes collision-finding exclusivity. If A is an adversary, we let $\text{Adv}_{H, \tilde{H}, P}^{\text{cfe}}(A) = \Pr \left[\mathbf{G}_{H, \tilde{H}, P}^{\text{cfe}}(A) \right]$ be its cfe advantage. This cfe notion requires that \tilde{h} remains collision-resistant, but the difference from standard cr is the addition of the GETPMG oracle, which allows an adversary to view preimages that have been produced by the exploit algorithm. These are subject to adversary-chosen constraint-parameters u . An adversary A wins game $\mathbf{G}_{H, \tilde{H}, P}^{\text{cfe}}$ if it produces any nontrivial collision, meaning, it cannot have asked for a preimage. The addition of this oracle can be viewed as a formalization of “backdoor key exposure” as raised by [34].

CONSTRUCTION. We construct an ASA on hash functions using our transform described in the previous section. In the following, we show that it achieves utility and cfe exclusivity if the target hash function is collision-resistant.

The difference from the general result is that when the function under consideration is a hash function, an effective ASA would expect that no other party can find collisions, as this is the standard security notion of a hash function. The cfe game asks that the presence of the ASA, and of an oracle to the preimage-finding exploit algorithm e , does not help anyone else find collisions in the hash function \tilde{h} . It is natural, then, that cfe additionally relies on cr of h , which matches \tilde{h} on most inputs.

Game $\mathbf{G}_{\tilde{H}}^{\text{cr}}$	Game $\mathbf{G}_{\tilde{H}, \tilde{H}, \tilde{P}}^{\text{cfe}}$
INIT: 1 $h \leftarrow \mathfrak{s} \mathbf{H}$; Return h	INIT: 1 $h \leftarrow \mathfrak{s} \mathbf{H}$; $(\tilde{h}, e) \leftarrow \mathfrak{s} \tilde{\mathbf{H}}(h)$
FIN(x_1, x_2): 2 If $(x_1 = x_2)$ then return false 3 Return $(h(x_1) = h(x_2))$	2 $\mathcal{X} \leftarrow \emptyset$ 3 Return (h, \tilde{h})
	GETPMG(u, y): 4 $x \leftarrow \mathfrak{s} e(u, y)$; $\mathcal{X} \leftarrow \mathcal{X} \cup \{x\}$ 5 Return x
	FIN(x_1, x_2): 6 Return $(x_1 \notin \mathcal{X}) \wedge (x_2 \notin \mathcal{X})$ $\quad \wedge (x_1 \neq x_2) \wedge (\tilde{h}(x_1) = \tilde{h}(x_2))$

Figure 8: **Left:** Collision resistance (cr) for a family of hash functions. **Right:** Collision-finding exclusivity (cfe) of an ASA on hash functions.

Proposition 5.1 *Let \mathbf{S} be a signature scheme, \mathbf{H} a family of hash functions, and Emb an embedding function for predicate \mathbf{P} which is compatible with \mathbf{S}, \mathbf{H} . Let $\tilde{\mathbf{H}} = \mathbf{ASA}[\mathbf{H}, \mathbf{S}, \text{Emb}]$. If \mathbf{S} and Emb are correct, then $\tilde{\mathbf{H}}$ achieves utility for \mathbf{P} .*

Utility follows directly from Proposition 4.1 by observing that utility for ASAs on hash functions is defined exactly as for general functions. Compatible embedding functions are, for example, the prefix and suffix embedding from Figure 4 or the certificate embedding in Section 8.

We now turn to cfe exclusivity; the proof of the below theorem will take advantage of Theorem 4.2 in the prior section.

Theorem 5.2 *Let \mathbf{S} be a signature scheme, \mathbf{H} a family of hash functions, and Emb a correct embedding function for predicate \mathbf{P} which is compatible with \mathbf{S}, \mathbf{H} . Let $\tilde{\mathbf{H}} = \mathbf{ASA}[\mathbf{H}, \mathbf{S}, \text{Emb}]$. Given an adversary A against the cfe exclusivity of $\tilde{\mathbf{H}}$ we can build adversaries $A_{\mathbf{S}}, A_{\mathbf{H}}$ such that*

$$\mathbf{Adv}_{\tilde{\mathbf{H}}, \tilde{\mathbf{H}}, \tilde{\mathbf{P}}}^{\text{cfe}}(A) \leq \mathbf{Adv}_{\mathbf{S}}^{\text{suf-cma}}(A_{\mathbf{S}}) + \mathbf{Adv}_{\mathbf{H}}^{\text{cr}}(A_{\mathbf{H}}). \quad (5)$$

If A makes q GETPMG queries, then $A_{\mathbf{S}}$ makes q SIGN queries. The running times of $A_{\mathbf{S}}, A_{\mathbf{H}}$ are close to that of A .

Proof of Theorem 5.2: Consider game \mathbf{G}_0 of Figure 9. It is easy to see that this game is exactly the cfe game. We already include line 7 which sets a flag `bad`, but in \mathbf{G}_0 this has no effect on the final output. Hence,

$$\mathbf{Adv}_{\tilde{\mathbf{H}}, \tilde{\mathbf{H}}, \tilde{\mathbf{P}}}^{\text{cfe}}(A) = \Pr[\mathbf{G}_0(A)].$$

We next turn to game \mathbf{G}_1 which outputs `false` whenever `bad` is set. Therefore, games $\mathbf{G}_0, \mathbf{G}_1$ are identical-until-`bad` and by the Fundamental Lemma of Game Playing [11] we have

$$\begin{aligned} \Pr[\mathbf{G}_0(A)] &= \Pr[\mathbf{G}_1(A)] + (\Pr[\mathbf{G}_0(A)] - \Pr[\mathbf{G}_1(A)]) \\ &\leq \Pr[\mathbf{G}_1(A)] + \Pr[\mathbf{G}_1(A) \text{ sets bad}]. \end{aligned}$$

Let us now take a closer look at the event defined in line 7. The game outputs `false` when the output of \tilde{h} and h differ on either x_1 or x_2 which is exactly captured by the exclusivity of $\tilde{\mathbf{H}}$. That

<p>Games $G_0, \boxed{G_1}, \boxed{G_2}$</p> <p>INIT:</p> <ol style="list-style-type: none"> 1 $h \leftarrow s H ; (\tilde{h}, e) \leftarrow s \tilde{H}(h) ; \mathcal{X} \leftarrow \emptyset$ 2 Return (h, \tilde{h}) <p>GETPMG(u, y):</p> <ol style="list-style-type: none"> 3 $x \leftarrow s e(u, y) ; \mathcal{X} \leftarrow \mathcal{X} \cup \{x\}$ 4 Return x <p>FIN(x_1, x_2):</p> <ol style="list-style-type: none"> 5 If $(x_1 = x_2)$ then return false 6 If $(x_1 \in \mathcal{X}) \vee (x_2 \in \mathcal{X})$ then return false 7 If $\tilde{h}(x_1) \neq h(x_1) \vee \tilde{h}(x_2) \neq h(x_2)$ then bad \leftarrow true ; return false 8 $y_1 \leftarrow \tilde{h}(x_1) ; y_2 \leftarrow \tilde{h}(x_2)$ 9 $y_1 \leftarrow h(x_1) ; y_2 \leftarrow h(x_2)$ // Game G_2 10 Return $(y_1 = y_2)$
--

Figure 9: Games G_0, G_1, G_2 for the proof of Theorem 5.2. G_1, G_2 contain the boxed code and G_0 does not. Line 9 is only present in G_2 .

<p>Adversary $A_{\tilde{H}}(\tilde{h}, \tilde{h})$:</p> <ol style="list-style-type: none"> 1 $\mathcal{X} \leftarrow \emptyset$ 2 $(x_1, x_2) \leftarrow A[\text{GETPMG}_{\tilde{H}}](\tilde{h}, \tilde{h})$ 3 If $(x_1 = x_2)$ then return \perp 4 If $(x_1 \in \mathcal{X}) \vee (x_2 \in \mathcal{X})$ then return \perp 5 If $\tilde{h}(x_1) \neq h(x_1)$ then return x_1 6 If $\tilde{h}(x_2) \neq h(x_2)$ then return x_2 <p>Oracle $\text{GETPMG}_{\tilde{H}}(u, y)$:</p> <ol style="list-style-type: none"> 7 $x \leftarrow \text{GETPMG}(y, u)$ 8 $\mathcal{X} \leftarrow \mathcal{X} \cup \{x\}$; Return x 	<p>Adversary $A_H(h)$:</p> <ol style="list-style-type: none"> 1 $(\tilde{h}, e) \leftarrow s \tilde{H}(h) ; \mathcal{X} \leftarrow \emptyset$ 2 $(x_1, x_2) \leftarrow A[\text{GETPMG}_H](h, \tilde{h})$ 3 Return (x_1, x_2) <p>Oracle $\text{GETPMG}_H(u, y)$:</p> <ol style="list-style-type: none"> 4 $x \leftarrow s e(u, y) ; \mathcal{X} \leftarrow \mathcal{X} \cup \{x\}$ 5 Return x
--	--

Figure 10: Adversaries $A_{\tilde{H}}$ (left) and A_H (right) for the proof of Theorem 5.2.

is, we can construct an adversary $A_{\tilde{H}}$ for which

$$\Pr[G_1(A) \text{ sets bad}] \leq \mathbf{Adv}_{H, \tilde{H}, P}^{\text{exc}}(A_{\tilde{H}}). \quad (6)$$

(We will discuss exclusivity and then move to suf-cma via Theorem 4.2.) Adversary $A_{\tilde{H}}$ is in game $\mathbf{G}_{H, \tilde{H}, P}^{\text{exc}}$ and runs A as specified on the left side of Figure 10. $A_{\tilde{H}}$ simulates A 's preimage oracle using its own preimage oracle. Suppose now that **bad** \leftarrow **true** on line 7 of G_1 . Then A has output x_1, x_2 such that $\tilde{h}(x_i) \neq h(x_i)$ for at least one $i \in \{1, 2\}$, while both are not in \mathcal{X} , meaning have not been the output of a query to GETPMG. Hence, $A_{\tilde{H}}$ has found a winning output $x^* \in \{x_1, x_2\}$ in the exclusivity game. This proves Eq. (6).

Now given this adversary $A_{\tilde{H}}$ in the exc game, we have adversary A_S in the suf-cma game by applying Theorem 4.2. Thus

$$\Pr[G_1(A) \text{ sets bad}] \leq \mathbf{Adv}_{H, \tilde{H}, P}^{\text{exc}}(A_{\tilde{H}}) \leq \mathbf{Adv}_S^{\text{suf-cma}}(A_S). \quad (7)$$

Note that if A makes q GETPMG queries then $A_{\tilde{H}}$ makes q GETPMG queries, and thus A_S from

Theorem 4.2 makes q SIGN queries.

We next turn to game G_2 , where the assignment in line 8 is replaced by the one in line 9. More specifically, y_1 and y_2 are now computed using h instead of using \tilde{h} . We claim that

$$\Pr[G_2(A)] = \Pr[G_1(A)]. \quad (8)$$

To justify Eq. (8), we observe that whenever the outputs of h and \tilde{h} for inputs x_1 or x_2 differ, the game has already returned `false`. Hence, we have $h(x_i) = \tilde{h}(x_i)$ for both $i \in \{0, 1\}$.

Finally, we claim that

$$\Pr[G_2(A)] \leq \mathbf{Adv}_{\mathbb{H}}^{\text{cr}}(A_{\mathbb{H}}). \quad (9)$$

We construct adversary $A_{\mathbb{H}}$ in game $\mathbf{G}_{\mathbb{H}}^{\text{cr}}$ as specified on the right side of Figure 10. A 's view is that of game G_2 ; initialization and `GETPMGH` return the same responses as in G_2 . Now, if $G_2(A)$ returns `true`, and since the boxed code is executed in G_2 , then it must be that $y_1 = y_2$ and thus $h(x_1) = h(x_2)$. This is precisely the winning condition of $A_{\mathbb{H}}$'s game $\mathbf{G}_{\mathbb{H}}^{\text{cr}}$ and proves Eq. (9). $A_{\mathbb{H}}$ maintains running time close to that of A . ■

HASHED PASSWORDS. We conclude our discussion of hash functions with an application to password-based authentication. Suppose a server stores a user password pwd along with a random salt s as $y = h(pwd \parallel s)$, as specified by PBKDF1 in PKCS#5 [45]. When someone tries to log in with password pwd' the server checks whether $h(pwd' \parallel s) = y$. Now suppose an ASA is mounted with respect to the suffix predicate, so that $(\tilde{h}, e) \leftarrow_{\mathbb{S}} \mathbb{H}(h)$ and \tilde{h} replaces h . Using e , the attacker can find $x \leftarrow_{\mathbb{S}} e(s, y)$ such that $x = pwd' \parallel s$ and $\tilde{h}(x) = y$. Thus someone in possession of e can effectively log in as any user. Note that this example requires a notion of a predicate and constraint parameter (the salt).

6 ASAs on proof system verification

Our second application of a public function is verification, for which we specifically look at proof systems. These have seen widespread usage in blockchains and cryptocurrencies, where attackers may have significant financial motivation to attack verification. However, we remark that verification is a common part of many schemes, to which our generic ASA of Section 4 applies.

In a (non-interactive) proof system, a prover who holds a statement ϕ and a witness ω for a given polynomial-time-decidable relation R wants to convince a verifier that ϕ is true; that is, $(\phi, \omega) \in R$. Usually, both parties have access to a common reference string (CRS). Instead of studying security in the presence of an untrusted CRS, as for example done in [8, 36], we target the verification procedure of proof systems, which can be attacked by an ASA. In related veins, ASAs have been considered on proof-of-work components by [63]; defenses against exfiltration have also been considered by [18]. Defenses specifically for verification programs were studied by [33]; they however target only non-malicious implementation errors.

SYNTAX. We now turn to ASAs on proof system verification. Following our general syntax, let \mathbf{PS} be a proof system. We write $(p, v) \leftarrow_{\mathbb{S}} \mathbf{PS}$ to denote the generation of corresponding proving and verification algorithms for \mathbf{PS} . Let \mathbf{V} be the function generator which runs as $\mathbf{V} : (p, v) \leftarrow_{\mathbb{S}} \mathbf{PS}$; return v . A verification function $v \in \text{OUT}(\mathbf{V})$ is usually described by a CRS and a relation R . We consider those to be hardcoded in v . Hence, verification takes as input a statement $\phi \in \{0, 1\}^*$ and proof $\pi \in \{0, 1\}^{\mathbf{PS}, \text{pl}}$ to produce a bit $d \in \{0, 1\}$; we write $d \leftarrow v(\phi, \pi)$. Further, we denote by L_R the language consisting of statements for which there exist matching witnesses in R .

<p>$\mathbf{P}_{\text{PS}}(x, u = \phi')$:</p> <ol style="list-style-type: none"> 1 $(\phi, \pi) \leftarrow x$ 2 Return $(\phi = \phi')$ <p>$\mathbf{Emb}_{\text{PS}}^{\ell, k}(z = (y \parallel \sigma), u = \phi')$:</p> <ol style="list-style-type: none"> 3 Require: $y = 1$ 4 $\pi \leftarrow \sigma \parallel 0^{k-\ell}$ 5 $x \leftarrow (\phi', \pi)$ 6 Return x <p>$\mathbf{Emblnv}_{\text{PS}}^{\ell, k}(x)$:</p> <ol style="list-style-type: none"> 7 $(\phi, \pi) \leftarrow x$ 8 $\sigma \leftarrow \pi[1..\ell]$ 9 Return $((1 \parallel \sigma), \phi)$ 	<p>$\tilde{\mathbf{V}}(v)$:</p> <ol style="list-style-type: none"> 1 $(vk, sk) \leftarrow \text{s S.Kg}$ 2 Define as below: <li style="padding-left: 20px;">$\tilde{v} : \{0, 1\}^* \rightarrow \{0, 1\}$ <li style="padding-left: 20px;">$e : \{0, 1\}^* \times \{0, 1\} \rightarrow \{0, 1\}^*$ 3 Return (\tilde{v}, e) <p>$\tilde{v}(x)$:</p> <ol style="list-style-type: none"> 4 $((1 \parallel \sigma), \phi) \leftarrow \mathbf{Emblnv}_{\text{PS}}(x)$ 5 If $\text{S.Vfy}(vk, (1, \phi), \sigma)$ then return 1 6 Else return $v(x)$ <p>$e(u = \phi', y)$:</p> <ol style="list-style-type: none"> 7 Require: $y = 1$ 8 $\sigma \leftarrow \text{s S.Sign}(sk, (1, \phi'))$ 9 $x \leftarrow \mathbf{Emb}_{\text{PS}}((1 \parallel \sigma), \phi')$ 10 Return x
---	---

Figure 11: **Left:** Fixed predicate \mathbf{P}_{PS} and a correct embedding for proof systems. We assume $\ell \leq k$. The embedding space is $\mathbf{Emb}_{\text{PS}}^{\ell, k}.\text{ES} = \{0, 1\}^{\ell+1}$; we have $z \in \{0, 1\}^{\ell+1}$. In our construction we would have $\ell = \text{S.sl}$ and $k = \text{PS.pl}$. **Right:** Our ASA construction on public functions, applied to proof system verification.

On the left side of Figure 12 we define soundness for PS. If A is an adversary, we let $\mathbf{Adv}_{\text{PS}}^{\text{snd}}(A) = \Pr \left[\mathbf{G}_{\text{PS}}^{\text{snd}}(A) \right]$ be its snd advantage.

An ASA on \mathbf{V} is specified by an algorithm $\tilde{\mathbf{V}}$ which produces $(\tilde{v}, e) \leftarrow \text{s } \tilde{\mathbf{V}}(v)$.

UTILITY. The canonical goal of a malicious prover is a forgery of a false statement; we use this as the utility goal of an ASA, viewing forgery as finding preimages of 1. To capture the ability of an attacker to forge arbitrary statements, we fix the predicate \mathbf{P}_{PS} as in Figure 11. There we also give an embedding function that is correct for \mathbf{P}_{PS} . Utility now asks that for any ϕ' , if $x \leftarrow \text{s } e(u = \phi', y = 1)$ then $x = (\phi, \pi)$ where both $\tilde{v}(x) = 1$ and $\phi = \phi'$. (The latter requirement is captured by $\mathbf{P}_{\text{PS}}(x, u)$.) In short, the exploit algorithm e allows one to forge a proof for any statement.

EXCLUSIVITY. We define exclusivity via the pfe (“proof-finding exclusivity”) game in Figure 12. Note that the pfe game is essentially the soundness game with the addition of the exploit-finding SIM oracle. If A is an adversary, we let $\mathbf{Adv}_{\text{PS}, \tilde{\mathbf{V}}}^{\text{pfe}}(A) = \Pr \left[\mathbf{G}_{\text{PS}, \tilde{\mathbf{V}}}^{\text{pfe}}(A) \right]$ be its pfe advantage. Recall that for general functions and hash functions, exclusivity is also parameterized by the predicate \mathbf{P} ; however for verification we assume the fixed predicate \mathbf{P}_{PS} above.

CONSTRUCTION. We consider the ASA construction using our transform from Section 4. We depict $\tilde{\mathbf{V}} = \mathbf{ASA}[\mathbf{V}, \mathbf{S}, \mathbf{Emb}_{\text{PS}}]$ for a family of verification functions \mathbf{V} , a signature scheme \mathbf{S} and compatible embedding \mathbf{Emb}_{PS} in Figure 11. In the following, we show effectiveness of this ASA, proving utility and pfe exclusivity.

Proposition 6.1 *Let \mathbf{S} be a signature scheme, PS a target proof system with verification family \mathbf{V} , and \mathbf{Emb}_{PS} an embedding function for predicate \mathbf{P}_{PS} which is compatible with \mathbf{S}, \mathbf{V} . Assume $\text{S.sl} \leq \text{PS.pl}$. Let $\tilde{\mathbf{V}} = \mathbf{ASA}[\mathbf{V}, \mathbf{S}, \mathbf{Emb}_{\text{PS}}]$. If \mathbf{S} and \mathbf{Emb}_{PS} are correct, then $\tilde{\mathbf{V}}$ achieves utility for*

Game $\mathbf{G}_{\text{PS}}^{\text{snd}}$	Game $\mathbf{G}_{\text{PS}, \tilde{\mathbf{V}}}^{\text{pfe}}$
INIT: 1 $(p, v) \leftarrow \text{PS}$ 2 Return (p, v) FIN(ϕ^*, π^*): 3 Return $((\phi^* \notin L_R)$ $\wedge (v(\phi^*, \pi^*) = 1))$	INIT: 1 $(p, v) \leftarrow \text{PS}$ 2 $(\tilde{v}, e) \leftarrow \tilde{\mathbf{V}}(v)$ 3 $\mathcal{X} \leftarrow \emptyset$ 4 Return (p, v, \tilde{v}) SIM(ϕ'): 5 $(\phi, \pi) \leftarrow e(\phi', 1)$; $\mathcal{X} \leftarrow \mathcal{X} \cup \{\phi\}$ 6 Return π FIN(ϕ^*, π^*): 7 Return $((\phi^* \notin \mathcal{X}) \wedge (\phi^* \notin L_R)$ $\wedge (\tilde{v}(\phi^*, \pi^*) = 1))$

Figure 12: **Left:** Soundness (snd) for a proof system. **Right:** Proof-finding exclusivity (pfe) of an ASA on proof system verification.

P_{PS} .

Proof of Proposition 6.1: We need to show that for every $v \in \text{OUT}(\mathbf{V})$ and every constraint-parameter (here, statement) $\phi \in \{0, 1\}^*$, if $(\tilde{v}, e) \leftarrow \tilde{\mathbf{V}}(v)$ and if $x \leftarrow e(\phi, 1)$, then we have (1) $\tilde{v}(x = (\phi, \pi)) = 1$ and (2) $\text{P}_{\text{PS}}(x, \phi) = 1$.

As defined in Figure 11, e produces $x \leftarrow e(u = \phi', y = 1)$ where $x = \text{Emb}_{\text{PS}}((1 \parallel \sigma), \phi')$ and $\sigma \leftarrow \text{S.Sign}(sk, (1, \phi'))$. That is, e outputs some verification input x which embeds both ϕ' and σ . If Emb_{PS} is a correct embedding function then $\text{P}_{\text{PS}}(x, u = \phi') = 1$. (An example of a correct embedding is in Figure 11. Recall that P_{PS} ensures that ϕ' is in fact the statement in x .)

The second requirement of utility is that $\tilde{v}(x) = 1$; that is, that a correct proof for ϕ' is produced. This is ensured by lines 4,5 of \tilde{v} in Figure 11. On this input x , which is a correctly produced embedding, line 4 parses $((1 \parallel \sigma), \phi) \leftarrow \text{EmbIn}_{\text{PS}}(x)$, where $\phi = \phi'$ by correctness of the embedding. Now line 5 computes $\text{S.Vfy}(vk, (1, \phi'), \sigma)$ which passes because, as in the paragraph above, σ is a signature on $(1, \phi')$ using scheme S . Verification thus returns 1 on line 5, proving that utility is achieved. ■

Theorem 6.2 *Let S be a signature scheme, PS a target proof system with verification family \mathbf{V} , and Emb_{PS} a correct embedding function for predicate P_{PS} which is compatible with S, \mathbf{V} . Assume $\text{S.sl} \leq \text{PS.pl}$. Let $\tilde{\mathbf{V}} = \text{ASA}[\mathbf{V}, \text{S}, \text{Emb}_{\text{PS}}]$. Given an adversary A against the pfe exclusivity of $\tilde{\mathbf{V}}$ we can build adversaries $A_{\text{S}}, A_{\text{PS}}$ such that*

$$\mathbf{Adv}_{\text{PS}, \tilde{\mathbf{V}}}^{\text{pfe}}(A) \leq \mathbf{Adv}_{\text{S}}^{\text{suf-cma}}(A_{\text{S}}) + \mathbf{Adv}_{\text{PS}}^{\text{snd}}(A_{\text{PS}}). \quad (10)$$

If A makes q queries to SIM, then A_{S} makes q SIGN queries. The running times of $A_{\text{S}}, A_{\text{PS}}$ are close to that of A .

We present the proof of Theorem 6.2 in Appendix A. We note that Emb_{PS} could be *any* correct and compatible embedding for P_{PS} . One example is $\text{Emb}_{\text{PS}}^{\ell, k}$ of Figure 11 but more clever embeddings could be constructed.

7 ASAs on signature verification

The prior section exhibits a verification function as a useful target of ASAs; there are examples beyond proof systems including commitment schemes, protocols and more. We turn to signatures as a slightly different setting than proof systems, as there are now per-user keys. Prior work has studied ASAs on signature schemes when key-generation [25] or when signing [5,62] are substituted. In our study of ASAs on public functions, we turn to consider signature verification. In particular, existing work looks at attacks on randomized signatures with the goal of exfiltrating the signing key. In this section, we consider even schemes that are deterministic or have unique signatures, in contrast to [5].

SYNTAX. Following our general ASA syntax, let TS be a (target) signature scheme. We write $(s, v) \leftarrow^s \text{TS}$ to denote the generation of corresponding signing and verification algorithms for signature scheme TS . Note that compared to usual notation, the signing key sk is now hardcoded in s and vk is hardcoded in v . Let \mathbf{V} be the function generator which runs as $\mathbf{V} : (s, v) \leftarrow^s \text{TS} ; \text{return } v$. A verification function $v \in \text{OUT}(\mathbf{V})$ takes as input a message $m \in \{0, 1\}^*$ and signature $\sigma \in \{0, 1\}^{\text{TS.sl}}$ to produce a bit $d \in \{0, 1\}$; we write $d \leftarrow v(m, \sigma)$. We introduced more standard signature syntax in Section 2 and in our constructions, but we use this syntax for TS to make clear that it is the target of the ASA.

An ASA on \mathbf{V} is specified by an algorithm $\tilde{\mathbf{V}}$ which produces $(\tilde{v}, e) \leftarrow^s \tilde{\mathbf{V}}(v)$. In the multi-user (mu) setting with n users, we write $((\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_n), e) \leftarrow^s \tilde{\mathbf{V}}(v_1, v_2, \dots, v_n)$. That is, there are n substitutions and only one exploitation algorithm. This is a more realistic attack to consider than simply mounting per-user exploit algorithms; in particular an attacker could modify library code to substitute the algorithms of n users while needing only one e . Once again we ask whether the ASA achieves both utility and exclusivity. (Effectiveness, which calls for utility, exclusivity and undetectability, is implied by the former two.)

UTILITY. The canonical goal of a signature attacker is forgery; we use this as the utility goal of an ASA, viewing forgery as finding preimages of 1. One could continue to ask, as in the general case, that the ASA be able to find inputs $x = (m, \sigma)$ such that $\tilde{v}(x) = 1$ and $\mathbf{P}(x, u) = 1$ for some predicate \mathbf{P} and constraint-parameter u . In this way, an ASA on signature verification is a special case of an ASA on general public functions. However, for utility of an ASA on signature verification, we ask for a quite strong notion: the attacker, via the exploit algorithm e , should be able to forge a signature *for any* message, not only messages satisfying a potentially restrictive predicate. We studied a similar notion of utility for proof system verification in the prior section.

In the mu setting, we let the exploit algorithm $e : \{0, 1\}^{\lceil \log n \rceil} \times \{0, 1\}^* \times \{0, 1\} \rightarrow \{0, 1\}^*$ take one additional input, indicating the user i for which to produce a forgery. We write $x \leftarrow^s e(i, u, y = 1)$ which constrains the target output to be 1 and allows for i as input in addition to the usual constraint-parameter u . Utility asks that for any m' and any i , if $x \leftarrow^s e(i, u = m', y = 1)$ then $x = (m', \sigma_{\text{TS}})$ for σ_{TS} such that $\tilde{v}_i(m', \sigma_{\text{TS}}) = 1$. That is, e allows one to forge a signature for any m' and i . This is captured by the predicate and embedding described in Figure 13; utility equivalently asks that both $\mathbf{P}_{\text{sig}}(x, u) = 1$ and that $\tilde{v}_i(m', \sigma_{\text{TS}}) = 1$.

The embedding $\text{Emb}_{\text{sig}}^{\ell, k}$ given for predicate \mathbf{P}_{sig} in Figure 13 is a correct embedding function, which is assumed to exist in the coming construction. Other embedding functions for \mathbf{P}_{sig} exist as well, and care could be taken to make them look more innocuous. In particular, embedding $\text{Emb}_{\text{sig}}^{\ell, k}$ is simplest when one uses the same signature scheme $\mathbf{S} = \text{TS}$ where $\ell = k$.

EXCLUSIVITY. We define exclusivity via the ffe (“forgery-finding exclusivity”) game in Figure 14, still in the mu setting. Note that the ffe game is essentially the multi-user uf-cma game with the

<p><u>$P_{\text{sig}}(x, u = m')$:</u></p> <ol style="list-style-type: none"> 1 $(m, \sigma_{\text{TS}}) \leftarrow x$ 2 Return $(m = m')$ 	<p><u>$\text{Emb}_{\text{sig}}^{\ell, k}(z = (y \parallel \sigma_{\text{S}}), u = m')$:</u></p> <ol style="list-style-type: none"> 1 Require: $y = 1$ 2 $\sigma_{\text{TS}} \leftarrow \sigma_{\text{S}} \parallel 0^{k-\ell}$ 3 $x \leftarrow (m', \sigma_{\text{TS}})$ 4 Return x <p><u>$\text{Emblnv}_{\text{sig}}^{\ell, k}(x)$:</u></p> <ol style="list-style-type: none"> 5 $(m, \sigma_{\text{TS}}) \leftarrow x$ 6 $\sigma_{\text{S}} \leftarrow \sigma_{\text{TS}}[1..\ell]$ 7 Return $((1 \parallel \sigma_{\text{S}}), m)$
---	---

Figure 13: Fixed predicate P_{sig} and a correct embedding for signatures. We assume $\ell \leq k$. The embedding space is $\text{Emb}_{\text{sig}}^{\ell, k} \cdot \text{ES} = \{0, 1\}^{\ell+1}$; we have $z \in \{0, 1\}^{\ell+1}$. In our construction we would have $\ell = \text{S.sl}$ and $k = \text{TS.sl}$.

<p><u>Game $\mathbf{G}_{\text{TS}, n}^{\text{mu-ufcma}}$</u></p> <p>INIT:</p> <ol style="list-style-type: none"> 1 For $i = 1, 2, \dots, n$ do $(s_i, v_i) \leftarrow_{\text{S}} \text{TS}$ 2 $\mathcal{Q} \leftarrow \emptyset$ 3 Return (v_1, v_2, \dots, v_n) <p>SIGN(i, m):</p> <ol style="list-style-type: none"> 4 $\sigma \leftarrow_{\text{S}} s_i(m)$; $\mathcal{Q} \leftarrow \mathcal{Q} \cup \{(i, m)\}$ 5 Return σ <p>FIN(i^*, m^*, σ^*):</p> <ol style="list-style-type: none"> 6 Return $((i^*, m^*) \notin \mathcal{Q})$ $\wedge (v_{i^*}(m^*, \sigma^*) = 1)$ 	<p><u>Game $\mathbf{G}_{\text{TS}, \tilde{\mathbf{V}}, n}^{\text{ffe}}$</u></p> <p>INIT:</p> <ol style="list-style-type: none"> 1 For $i = 1, 2, \dots, n$ do $(s_i, v_i) \leftarrow_{\text{S}} \text{TS}$ 2 $((\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_n), e) \leftarrow_{\text{S}} \tilde{\mathbf{V}}(v_1, v_2, \dots, v_n)$ 3 $\mathcal{X} \leftarrow \emptyset$ 4 Return $((v_1, v_2, \dots, v_n), (\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_n))$ <p>ESIGN(i, m'):</p> <ol style="list-style-type: none"> 5 $(m, \sigma) \leftarrow_{\text{S}} e(i, m', 1)$; $\mathcal{X} \leftarrow \mathcal{X} \cup \{(i, m)\}$ 6 Return σ <p>SIGN(i, m):</p> <ol style="list-style-type: none"> 7 $\sigma \leftarrow_{\text{S}} s_i(m)$; $\mathcal{X} \leftarrow \mathcal{X} \cup \{(i, m)\}$ 8 Return σ <p>FIN(i^*, m^*, σ^*):</p> <ol style="list-style-type: none"> 9 Return $((i^*, m^*) \notin \mathcal{X}) \wedge (\tilde{v}_{i^*}(m^*, \sigma^*) = 1)$
---	---

Figure 14: **Left:** The multi-user uf-cma game over n users. **Right:** Forgery-finding exclusivity (ffe) of an ASA on signature verification.

addition of the exploit-finding ESIGN oracle, following the framework of the prior two sections. If A is an adversary, we let $\text{Adv}_{\text{TS}, \tilde{\mathbf{V}}, n}^{\text{ffe}}(A) = \Pr \left[\mathbf{G}_{\text{TS}, \tilde{\mathbf{V}}, n}^{\text{ffe}}(A) \right]$ be its ffe advantage.

CONSTRUCTION. We construct a multi-user ASA on signatures using a slight modification to our main transform. Given a signature scheme S , a family of verification functions \mathbf{V} for target signature scheme TS , and an embedding function Emb_{sig} that is compatible with S, \mathbf{V} and is correct for predicate P_{sig} , we let $\tilde{\mathbf{V}} = \mathbf{ASA}_n[\mathbf{V}, \text{S}, \text{Emb}_{\text{sig}}]$. This is specified in Figure 15, both for the signature case and for the more general mu case. $\tilde{\mathbf{V}}$ is described on the left side of Figure 15 and is constructed using \mathbf{V}, S , and Emb_{sig} .

In the following proposition and theorem statement, we prove that this is an effective ASA, satisfying utility and ffe exclusivity.

<p>$\tilde{V}(v_1, v_2, \dots, v_n)$:</p> <ol style="list-style-type: none"> 1 $(vk, sk) \leftarrow \text{S.Kg}$ 2 Define as below: <ul style="list-style-type: none"> $\tilde{v}_i : \{0, 1\}^* \rightarrow \{0, 1\}$ $e : \{0, 1\}^{\lceil \log n \rceil} \times \{0, 1\}^* \times \{0, 1\} \rightarrow \{0, 1\}^*$ 3 Return $((\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_n), e)$ <p>$\tilde{v}_i(x)$:</p> <ol style="list-style-type: none"> 4 $((1 \parallel \sigma_S), m) \leftarrow \text{Emblnv}_{\text{sig}}(x)$ 5 If $\text{S.Vfy}(vk, (i, 1, m), \sigma_S)$ then return 1 6 Else return $v_i(x)$ <p>$e(i, u = m', y)$:</p> <ol style="list-style-type: none"> 7 Require: $y = 1$ 8 $\sigma_S \leftarrow \text{S.Sign}(sk, (i, 1, m'))$ 9 $x \leftarrow \text{Emb}_{\text{sig}}((1 \parallel \sigma_S), m')$ 10 Return x 	<p>$\tilde{F}(f_1, f_2, \dots, f_n)$:</p> <ol style="list-style-type: none"> 1 $(vk, sk) \leftarrow \text{S.Kg}$ 2 Define as below: <ul style="list-style-type: none"> $\tilde{f}_i : \{0, 1\}^* \rightarrow \{0, 1\}^{\text{F.ol}}$ $e : \{0, 1\}^{\lceil \log n \rceil} \times \{0, 1\}^* \times \{0, 1\}^{\text{F.ol}} \rightarrow \{0, 1\}^*$ 3 Return $((\tilde{f}_1, \tilde{f}_2, \dots, \tilde{f}_n), e)$ <p>$\tilde{f}_i(x)$:</p> <ol style="list-style-type: none"> 4 $w \leftarrow \text{Emb}^{-1}(x)$ 5 If $(w = \perp)$ then return $f_i(x)$ 6 $((y \parallel \sigma), u) \leftarrow w$ 7 If $\text{S.Vfy}(vk, (i, y, u), \sigma)$ then return y 8 Else return $f_i(x)$ <p>$e(i, u, y)$:</p> <ol style="list-style-type: none"> 9 Require: $y \in \{0, 1\}^{\text{F.ol}}$ 10 $\sigma \leftarrow \text{S.Sign}(sk, (i, y, u))$ 11 $x \leftarrow \text{Emb}((y \parallel \sigma), u)$ 12 Return x
---	--

Figure 15: **Left:** Construction of an ASA on multi-user signature verification. **Right:** Our general ASA on public functions from Section 4, adapted to multiple users. Recall that vk is considered hardcoded in \tilde{v}_i (or \tilde{f}_i) and sk is hardcoded in e .

Proposition 7.1 *Let S be a signature scheme, TS a target signature scheme with verification family V , and Emb_{sig} an embedding function for predicate P_{sig} which is compatible with S, V . Assume $S.\text{sl} \leq TS.\text{sl}$. Let $\tilde{V} = \text{ASA}_n[V, S, \text{Emb}_{\text{sig}}]$ as in Figure 15. If S and Emb_{sig} are correct, then \tilde{V} achieves utility for P_{sig} .*

Proof of Proposition 7.1: For $i = 1, 2, \dots, n$ let $(s_i, v_i) \leftarrow \text{S}$ and let $((\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_n), e) \leftarrow \text{S}$. Let $m' \in \{0, 1\}^*$ be any message and i be any user $1 \leq i \leq n$. Utility asks that the exploit algorithm can produce a signature on m' which validates for user i .

As defined in Figure 15, e produces $x \leftarrow e(i, u = m', y = 1)$ where $x = \text{Emb}_{\text{sig}}((1 \parallel \sigma_S), m')$ and $\sigma_S \leftarrow \text{S.Sign}(sk, (i, 1, m'))$. That is, e outputs some verification input x which embeds both m' and σ_S . If Emb_{sig} is a correct embedding function then $P_{\text{sig}}(x, u = m') = 1$. (An example of a correct embedding is in Figure 13. Recall that P_{sig} ensures that m' is in fact the message in x .)

The second requirement of utility is that $\tilde{v}_i(x) = 1$; that is, that a correct signature on m' for user i is produced. This is ensured by lines 4,5 of \tilde{v}_i in Figure 15. On this input x , which is a correctly produced embedding, line 4 parses $((1 \parallel \sigma_S), m) \leftarrow \text{Emblnv}_{\text{sig}}(x)$, where $m = m'$ by correctness of the embedding. Now line 5 computes $\text{S.Vfy}(vk, (i, 1, m'), \sigma_S)$ which passes because, as in the paragraph above, σ_S is a signature on $(i, 1, m')$ using scheme S . Verification thus returns 1 on line 5, proving that utility is achieved. \blacksquare

Theorem 7.2 *Let S be a signature scheme, TS a target signature scheme with verification family V , and Emb_{sig} a correct embedding function for predicate P_{sig} which is compatible with S, V . Assume $S.\text{sl} \leq TS.\text{sl}$. Let $\tilde{V} = \text{ASA}_n[V, S, \text{Emb}_{\text{sig}}]$ as in Figure 15. Given an adversary A against the ffe*

exclusivity of \tilde{V} we can build adversaries A_S, A_{TS} such that

$$\mathbf{Adv}_{TS, \tilde{V}, n}^{\text{ffe}}(A) \leq \mathbf{Adv}_S^{\text{suf-cma}}(A_S) + \mathbf{Adv}_{TS, n}^{\text{mu-ufcma}}(A_{TS}). \quad (11)$$

If A makes q_s, q_e SIGN, ESIGN queries, respectively, then A_S makes q_e SIGN queries and A_{TS} makes q_s SIGN queries. The running times of A_S, A_{TS} are close to that of A .

We present the proof of Theorem 7.2 in Appendix B.

SINGLE-USER FFE $\not\Rightarrow$ MULTI-USER FFE. Above, we defined game $\mathbf{G}_{TS, \tilde{V}, n}^{\text{ffe}}$ in a multi-user setting. In contrast to the unforgeability of a normal signature scheme, ffe unforgeability in the single-user setting does not generally imply unforgeability in the multi-user setting. The intuitive reason for this is that the signature created using the exploit algorithm depends on the target user i for which the signature is forged. As a concrete counterexample, we can compare construction **ASA** to **ASA_n** as defined in Figure 15. In the latter, the exploit algorithm computes $\sigma_S \leftarrow \text{S.Sign}(sk, (i, 1, m'))$ while i is not included in the signed value in the former. In the multi-user setting, if i is not included, one could use the fact that one exploit-produced signature will verify for different (if not all) users. Multi-user considerations thus arise for this ASA setting.

8 Application: Forged certificates

CERTIFICATES AND PKI. In this section we discuss an application and embedding method which is relevant to hash functions and signatures as used in public-key infrastructure (PKI). The usual realization of PKI uses X.509 certificates and certificate authorities. For simplicity, suppose there is one honest CA who is operating with signature scheme S_{ca} and hash function $h \in H_{ca}$. Let (vk_{ca}, sk_{ca}) be the verification and signing keys of the CA; in our ASA syntax the signature verification algorithm is $v(\cdot) = S_{ca}.Vfy(vk_{ca}, \cdot)$. Thus all users who use the PKI with this CA have implementations of h and v (along with vk_{ca}) on their devices.

As specified in RFC 5280 [23], a certificate C consists of a sequence of key-value pairs. The important fields for our discussion are:

- `C.tbsCert`, consisting of the certificate’s identifying, validity, and other certificate data. At a minimum, this specifies the CA who signed the certificate and includes information to recover vk_{ca} .
- `C.sigAlg`, the name of the signature algorithm, such as “PKCS #1 SHA-256 With RSA Encryption.” Along with the first item above, this allows one to recover algorithms h and v .
- `C.sigValue`, a signature on message `C.tbsCert`, using the algorithm specified in `C.sigAlg` and the CA’s signing key sk_{ca} .

Issuance of a certificate takes as input a `tbsCert'` and auxiliary information `csr` (representing a certificate signing request) to produce either \perp , or a signed certificate C . Deterministic validation of a certificate takes as input a certificate C and auxiliary information `aux` (representing a certificate chain, and local store of root certificates) to produce a bit $d \in \{0, 1\}$.

In our discussion, we are more interested in validation than issuance. At a high level, validation `Validate` proceeds as in Figure 16. We now discuss two applications of ASAs on hash functions and signatures, though we also remark that the certificate validation process has other “general public functions” to which our general framework could apply.

ASA ON THE HASH FUNCTION. Let $\tilde{H} = \mathbf{ASA}[H_{ca}, S, \text{Emb}_{cert}]$ given a signature scheme S with *short* signatures and an embedding function that will be specified shortly. Let $(\tilde{h}, e) \leftarrow \text{S} \tilde{H}(h)$. As in Figure 16, an attacker substitutes h with \tilde{h} on a user’s device.

<p><u>Validate(C, aux):</u></p> <ol style="list-style-type: none"> 1 First, perform expiry, revocation, and other checks. 2 Extract vk_{ca} from (C, aux) 3 Extract (h, v) from C.sigAlg and vk_{ca} 4 $y \leftarrow h(\text{C.tbsCert})$ 5 Return $v(y, \text{C.sigValue})$ <p><u>ASA on h changes line 4:</u></p> <ol style="list-style-type: none"> 6 $y \leftarrow \tilde{h}(\text{C.tbsCert})$ <p><u>ASA on v changes line 5:</u></p> <ol style="list-style-type: none"> 7 Return $\tilde{v}(y, \text{C.sigValue})$
--

Figure 16: Validation of an X.509 certificate, simplified. Recall that $v(\cdot) = S_{ca}.Vfy(vk_{ca}, \cdot)$.

<p><u>$P_{cert}(x = \text{tbsCert}^*, u = \text{tbsCert}')$:</u></p> <ol style="list-style-type: none"> 1 If $((\text{tbsCert}^*.f_h = \perp) \vee (\text{tbsCert}^*.f_\sigma = \perp))$ then return false 2 If $((\text{tbsCert}^*.f = \text{tbsCert}'.f)$ for all other fields f) then return true 3 Else return false <p><u>$\text{Emb}_{cert}((y \parallel \sigma), u = \text{tbsCert}')$:</u></p> <ol style="list-style-type: none"> 4 For all fields f, do: $\text{tbsCert}^*.f \leftarrow \text{tbsCert}'.f$ 5 $\text{tbsCert}^*.f_h \leftarrow y$; $\text{tbsCert}^*.f_\sigma \leftarrow \sigma$ 6 Return tbsCert^* <p><u>$\text{Emb}_{cert}^{-1}(x = \text{tbsCert}^*)$:</u></p> <ol style="list-style-type: none"> 7 $y \leftarrow \text{tbsCert}^*.f_h$; $\sigma \leftarrow \text{tbsCert}^*.f_\sigma$ 8 $\text{tbsCert}'.f \leftarrow \text{tbsCert}^*.f$ for all fields f except f_h, f_σ 9 Return $((y \parallel \sigma), \text{tbsCert}')$ <hr/> <p><u>Certificate forgery using \tilde{H}:</u></p> <ol style="list-style-type: none"> 10 $(\tilde{h}, e) \leftarrow \tilde{s} \tilde{H}(h)$ 11 Choose a target honest certificate C with hash y. 12 Choose any $\text{tbsCert}'$ which specifies the same issuer CA and vk_{ca} as C. 13 $\text{tbsCert}^* \leftarrow \tilde{s} e(u = \text{tbsCert}', y)$ 14 $\text{C}^*.tbsCert \leftarrow \text{tbsCert}^*$ 15 $\text{C}^*.sigAlg$ specifies (h, v) 16 $\text{C}^*.sigValue \leftarrow \text{C.sigValue}$ 17 Return C^*
--

Figure 17: **Above:** Predicate P_{cert} which captures whether data tbsCert^* embeds all of the data from chosen $\text{tbsCert}'$, along with a correct embedding function for this predicate. This captures the attacker’s goal of using arbitrary data in their forgery. We give an explicit example of this embedding in Appendix C. **Below:** How an ASA on a hash function h allows the attacker to forge certificates. We use f to denote any X.509 field name. Field labels f_h, f_σ are fixed.

At a high level, certificate forgery proceeds by choosing a target hash y for which a CA signature is already known, on an honest certificate C. Then the attacker finds a preimage certificate C^* which maps to y under \tilde{h} . (Thus the signature on C may be reused.) The ability to specify what this preimage certificate “looks like” is captured by a predicate and embedding. We give one such example in Figure 17. Here, the predicate enforces that the attacker can choose any desired certificate data $\text{tbsCert}'$, and can find a forgery C^* which differs only by adding $(H_{ca}.ol + S.sl)$

bits. If \tilde{H} is effective for this P_{cert} , and if $\text{Validate}(C, \text{aux}) = 1$, then $\text{Validate}(C^*, \text{aux}) = 1$. The steps are written in more detail in Figure 17. Moreover, we give an explicit example of a certificate embedding in Appendix C but we remark that our transform works for *any* correct predicate and embedding. In particular, one could design other ways to embed $(H_{\text{ca.ol}} + S.\text{sl})$ bits in a certificate.

Although we consider the ASA setting, we note that a variety of work has already studied the repercussions of PKI allowing weak hash functions, in particular MD5 [48, 49, 61, 65, 69]. At Eurocrypt 2007 [60], Stevens, Lenstra, and de Weger presented two X.509 certificates with the same MD5 hash value and thus signature. The two certificates were produced at the same time, and the cost was estimated to be “2 months real time” [60]. In the ASA model, an attacker wants to forge (almost) arbitrary certificates, at arbitrary times, and to do so easily.

ASA ON SIGNATURE VERIFICATION. Rather than attacking the hashing step of validation, one could attack the signature verification. Let $(\tilde{v}, e) \leftarrow_s \tilde{V}(v)$ for an ASA \tilde{V} . As in Figure 16, an attacker substitutes v with \tilde{v} on a user’s device. Given any certificate data $C^*.\text{tbsCert}$, the task is now to set $y \leftarrow h(C^*.\text{tbsCert})$ and find $C^*.\text{sigValue}$ such that $\tilde{v}(y, C^*.\text{sigValue}) = 1$. Note that the embedding task is already accomplished by the utility definition for a signature ASA; the message y may be arbitrarily selected, and the exploit algorithm returns $(y, \sigma) \leftarrow_s e(y, 1)$ such that $\tilde{v}(y, \sigma) = 1$. Now we simply set $C^*.\text{sigValue} \leftarrow \sigma$.

Thus, if we instantiate $\tilde{V} = \mathbf{ASA}[V_{\text{ca}}, S, \text{Emb}_{\text{sig}}]$ with a signature scheme such that $S.\text{sl} \leq S_{\text{ca.}}.\text{sl}$, the results of Section 7 show that the ASA can construct arbitrary certificates with validating signatures under \tilde{v} . (No special certificate embedding is needed for this ASA.)

MANY USERS AND CAS. We have only described a single-user setting in this section. In Section 7, we pointed out the difference between single-user and multi-user due to exclusivity. This arises in practice in PKI, as many users are trusting several CAs. The observation in the signature setting was that the exploit algorithm e must consider the user i (or their identifier) when producing a preimage, which should depend on i . In our construction, this meant that i was included in the signature embedded in the found preimage. The same observation would apply in a many-user, many-CA PKI setting if the ASA aims to prevent exploit-produced preimages from being reused for different users.

References

- [1] H. Abelson, R. Anderson, S. M. Bellovin, J. Benaloh, M. Blaze, W. Diffie, J. Gilmore, M. Green, S. Landau, P. G. Neumann, R. L. Rivest, J. I. Schiller, B. Schneier, M. Specter, and D. J. Weitzner. Keys under doormats: Mandating insecurity by requiring government access to all data and communications, 2015. <https://dspace.mit.edu/bitstream/handle/1721.1/97690/MIT-CSAIL-TR-2015-026.pdf>. (Cited on 5.)
- [2] A. Albertini, J.-P. Aumasson, M. Eichlseder, F. Mendel, and M. Schl affer. Malicious hashing: Eve’s variant of SHA-1. In A. Joux and A. M. Youssef, editors, *SAC 2014*, volume 8781 of *LNCS*, pages 1–19. Springer, Heidelberg, Aug. 2014. (Cited on 4, 6, 14.)
- [3] M. Armour and B. Poettering. Subverting decryption in AEAD. In M. Albrecht, editor, *17th IMA International Conference on Cryptography and Coding*, volume 11929 of *LNCS*, pages 22–41. Springer, Heidelberg, Dec. 2019. (Cited on 2, 5.)
- [4] G. Ateniese, D. Francati, B. Magri, and D. Venturi. Public immunization against complete subversion without random oracles. In R. H. Deng, V. Gauthier-Uma na, M. Ochoa, and

- M. Yung, editors, *ACNS 19*, volume 11464 of *LNCS*, pages 465–485. Springer, Heidelberg, June 2019. (Cited on 5.)
- [5] G. Ateniese, B. Magri, and D. Venturi. Subversion-resilient signature schemes. In I. Ray, N. Li, and C. Kruegel, editors, *ACM CCS 2015*, pages 364–375. ACM Press, Oct. 2015. (Cited on 2, 5, 20.)
- [6] B. Auerbach, M. Bellare, and E. Kiltz. Public-key encryption resistant to parameter subversion and its realization from efficiently-embeddable groups. In M. Abdalla and R. Dahab, editors, *PKC 2018, Part I*, volume 10769 of *LNCS*, pages 348–377. Springer, Heidelberg, Mar. 2018. (Cited on 6.)
- [7] F. Banfi, K. Gegier, M. Hirt, U. Maurer, and G. Rito. Anamorphic encryption, revisited. In M. Joye and G. Leander, editors, *Advances in Cryptology – EUROCRYPT 2024*, pages 3–32, Cham, 2024. Springer Nature Switzerland. (Cited on 6.)
- [8] M. Bellare, G. Fuchsbauer, and A. Scafuro. NIZKs with an untrusted CRS: Security in the face of parameter subversion. In J. H. Cheon and T. Takagi, editors, *ASIACRYPT 2016, Part II*, volume 10032 of *LNCS*, pages 777–804. Springer, Heidelberg, Dec. 2016. (Cited on 6, 17.)
- [9] M. Bellare, J. Jaeger, and D. Kane. Mass-surveillance without the state: Strongly undetectable algorithm-substitution attacks. In I. Ray, N. Li, and C. Kruegel, editors, *ACM CCS 2015*, pages 1431–1440. ACM Press, Oct. 2015. (Cited on 2, 5.)
- [10] M. Bellare, K. G. Paterson, and P. Rogaway. Security of symmetric encryption against mass surveillance. In J. A. Garay and R. Gennaro, editors, *CRYPTO 2014, Part I*, volume 8616 of *LNCS*, pages 1–19. Springer, Heidelberg, Aug. 2014. (Cited on 2, 3, 5.)
- [11] M. Bellare and P. Rogaway. The security of triple encryption and a framework for code-based game-playing proofs. In S. Vaudenay, editor, *EUROCRYPT 2006*, volume 4004 of *LNCS*, pages 409–426. Springer, Heidelberg, May / June 2006. (Cited on 6, 9, 13, 15, 32.)
- [12] P. Bemmman, S. Berndt, and R. Chen. Subversion-resilient signatures without random oracles. In C. Pöpper and L. Batina, editors, *Applied Cryptography and Network Security*, pages 351–375, Cham, 2024. Springer Nature Switzerland. (Cited on 14.)
- [13] P. Bemmman, R. Chen, and T. Jager. Subversion-resilient public key encryption with practical watchdogs. In J. Garay, editor, *PKC 2021, Part I*, volume 12710 of *LNCS*, pages 627–658. Springer, Heidelberg, May 2021. (Cited on 5.)
- [14] S. Berndt, J. Wichelmann, C. Pott, T.-H. Traving, and T. Eisenbarth. ASAP: Algorithm substitution attacks on cryptographic protocols. In Y. Suga, K. Sakurai, X. Ding, and K. Sako, editors, *ASIACCS 22*, pages 712–726. ACM Press, May / June 2022. (Cited on 2, 5.)
- [15] D. J. Bernstein, T. Lange, and R. Niederhagen. Dual EC: A standardized back door. In P. Ryan, D. Naccache, and J.-J. Quisquater, editors, *The New Codebreakers*. Springer, Heidelberg, 2016. (Cited on 5.)
- [16] D. Boneh, B. Lynn, and H. Shacham. Short signatures from the Weil pairing. In C. Boyd, editor, *ASIACRYPT 2001*, volume 2248 of *LNCS*, pages 514–532. Springer, Heidelberg, Dec. 2001. (Cited on 34.)

- [17] D. Boneh, E. Shen, and B. Waters. Strongly unforgeable signatures based on computational Diffie-Hellman. In M. Yung, Y. Dodis, A. Kiayias, and T. Malkin, editors, *PKC 2006*, volume 3958 of *LNCS*, pages 229–240. Springer, Heidelberg, Apr. 2006. (Cited on 7.)
- [18] S. Chakraborty, C. Ganesh, and P. Sarkar. Reverse firewalls for oblivious transfer extension and applications to zero-knowledge. In C. Hazay and M. Stam, editors, *EUROCRYPT 2023, Part I*, volume 14004 of *LNCS*, pages 239–270. Springer, Heidelberg, Apr. 2023. (Cited on 17.)
- [19] R. Chen, X. Huang, and M. Yung. Subvert KEM to break DEM: Practical algorithm-substitution attacks on public-key encryption. In S. Moriai and H. Wang, editors, *ASIACRYPT 2020, Part II*, volume 12492 of *LNCS*, pages 98–128. Springer, Heidelberg, Dec. 2020. (Cited on 2, 5.)
- [20] S. S. M. Chow, A. Russell, Q. Tang, M. Yung, Y. Zhao, and H.-S. Zhou. Let a non-barking watchdog bite: Cryptographic signatures with an offline watchdog. In D. Lin and K. Sako, editors, *PKC 2019, Part I*, volume 11442 of *LNCS*, pages 221–251. Springer, Heidelberg, Apr. 2019. (Cited on 2.)
- [21] B. Cogliati, J. Ethan, and A. Jha. Subverting telegram’s end-to-end encryption. *IACR Transactions on Symmetric Cryptology*, 2023. <https://tosc.iacr.org/index.php/ToSC/article/view/10302/9747>. (Cited on 2, 5.)
- [22] S. Contini, A. K. Lenstra, and R. Steinfeld. VSH, an efficient and provable collision-resistant hash function. In S. Vaudenay, editor, *EUROCRYPT 2006*, volume 4004 of *LNCS*, pages 165–182. Springer, Heidelberg, May / June 2006. (Cited on 4, 14.)
- [23] D. Cooper, S. Santesson, S. Farrell, S. Boeyen, R. Housley, and W. Polk. RFC 5280: Internet X.509 public key infrastructure certificate and certificate revocation list (CRL) profile, May 2008. (Cited on 23.)
- [24] R. Cramer and V. Shoup. Signature schemes based on the strong RSA assumption. In J. Motiwalla and G. Tsudik, editors, *ACM CCS 99*, pages 46–51. ACM Press, Nov. 1999. (Cited on 7.)
- [25] C. Crépeau and A. Slakmon. Simple backdoors for RSA key generation. In M. Joye, editor, *CT-RSA 2003*, volume 2612 of *LNCS*, pages 403–416. Springer, Heidelberg, Apr. 2003. (Cited on 20.)
- [26] J. P. Degabriele, P. Farshim, and B. Poettering. A more cautious approach to security against mass surveillance. In G. Leander, editor, *FSE 2015*, volume 9054 of *LNCS*, pages 579–598. Springer, Heidelberg, Mar. 2015. (Cited on 2, 5.)
- [27] J. P. Degabriele, K. G. Paterson, J. C. N. Schuldt, and J. Woodage. Backdoors in pseudorandom number generators: Possibility and impossibility results. In M. Robshaw and J. Katz, editors, *CRYPTO 2016, Part I*, volume 9814 of *LNCS*, pages 403–432. Springer, Heidelberg, Aug. 2016. (Cited on 6.)
- [28] D. Derler, K. Samelin, and D. Slamanig. Bringing order to chaos: The case of collision-resistant chameleon-hashes. In A. Kiayias, M. Kohlweiss, P. Wallden, and V. Zikas, editors, *PKC 2020, Part I*, volume 12110 of *LNCS*, pages 462–492. Springer, Heidelberg, May 2020. (Cited on 14.)

- [29] Y. Dodis, C. Ganesh, A. Golovnev, A. Juels, and T. Ristenpart. A formal treatment of backdoored pseudorandom generators. In E. Oswald and M. Fischlin, editors, *EUROCRYPT 2015, Part I*, volume 9056 of *LNCS*, pages 101–126. Springer, Heidelberg, Apr. 2015. (Cited on 6.)
- [30] Y. Dodis, I. Mironov, and N. Stephens-Davidowitz. Message transmission with reverse firewalls—secure communication on corrupted machines. In M. Robshaw and J. Katz, editors, *CRYPTO 2016, Part I*, volume 9814 of *LNCS*, pages 341–372. Springer, Heidelberg, Aug. 2016. (Cited on 5.)
- [31] N. Döttling, S. Garg, Y. Ishai, G. Malavolta, T. Mour, and R. Ostrovsky. Trapdoor hash functions and their applications. In A. Boldyreva and D. Micciancio, editors, *CRYPTO 2019, Part III*, volume 11694 of *LNCS*, pages 3–32. Springer, Heidelberg, Aug. 2019. (Cited on 14.)
- [32] E. Felten. The linux backdoor attempt of 2003, 2013. <https://freedom-to-tinker.com/2013/10/09/the-linux-backdoor-attempt-of-2003/>. (Cited on 5.)
- [33] M. Fischlin and F. Günther. Verifiable verification in cryptographic protocols. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security, CCS '23*, page 3239–3253, New York, NY, USA, 2023. Association for Computing Machinery. (Cited on 17.)
- [34] M. Fischlin, C. Janson, and S. Mazaheri. Backdoored hash functions: Immunizing HMAC and HKDF. In S. Chong and S. Delaune, editors, *CSF 2018 Computer Security Foundations Symposium*, pages 105–118. IEEE Computer Society Press, 2018. (Cited on 2, 4, 6, 10, 14.)
- [35] A. Freund. Openwall oss-security mailing list, 29 Mar. 2024. <https://www.openwall.com/lists/oss-security/2024/03/29/4>. (Cited on 6.)
- [36] G. Fuchsbauer. Subversion-zero-knowledge SNARKs. In M. Abdalla and R. Dahab, editors, *PKC 2018, Part I*, volume 10769 of *LNCS*, pages 315–347. Springer, Heidelberg, Mar. 2018. (Cited on 17.)
- [37] R. Gennaro, S. Halevi, and T. Rabin. Secure hash-and-sign signatures without the random oracle. In J. Stern, editor, *EUROCRYPT'99*, volume 1592 of *LNCS*, pages 123–139. Springer, Heidelberg, May 1999. (Cited on 7.)
- [38] E.-J. Goh, D. Boneh, B. Pinkas, and P. Golle. The design and implementation of protocol-based hidden key recovery. In C. Boyd and W. Mao, editors, *ISC 2003*, volume 2851 of *LNCS*, pages 165–179. Springer, Heidelberg, Oct. 2003. (Cited on 2, 5.)
- [39] O. Goldreich. *Foundations of Cryptography: Basic Applications*, volume 2. Cambridge University Press, Cambridge, UK, 2004. (Cited on 7.)
- [40] S. Goldwasser, M. P. Kim, V. Vaikuntanathan, and O. Zamir. Planting undetectable backdoors in machine learning models : [extended abstract]. In *63rd FOCS*, pages 931–942. IEEE Computer Society Press, Oct. / Nov. 2022. (Cited on 2, 3, 4, 6.)
- [41] P. Hodges and D. Stebila. Algorithm substitution attacks: State reset detection and asymmetric modifications. *IACR Trans. Symm. Cryptol.*, 2021(2):389–422, 2021. (Cited on 2.)
- [42] D. Hofheinz and E. Kiltz. Programmable hash functions and their applications. *Journal of Cryptology*, 25(3):484–527, July 2012. (Cited on 14.)

- [43] T. Horel, S. Park, S. Richelson, and V. Vaikuntanathan. How to subvert backdoored encryption: Security against adversaries that decrypt all ciphertexts. In A. Blum, editor, *ITCS 2019*, volume 124, pages 42:1–42:20. LIPIcs, Jan. 2019. (Cited on 6.)
- [44] A. Joux, J. Loss, and B. Wagner. Kleptographic attacks against implicit rejection. Cryptology ePrint Archive, Report 2024/260, 2024. <https://eprint.iacr.org/2024/260>. (Cited on 2, 5.)
- [45] B. Kaliski. PKCS #5: Password-based cryptography specification. RSA Laboratories, Sept. 2000. Version 2.0. (Cited on 17.)
- [46] H. Krawczyk and T. Rabin. Chameleon signatures. In *NDSS 2000*. The Internet Society, Feb. 2000. (Cited on 14.)
- [47] M. Kutylowski, G. Persiano, D. H. Phan, M. Yung, and M. Zawada. Anamorphic signatures: Secrecy from a dictator who only permits authentication! In H. Handschuh and A. Lysyanskaya, editors, *CRYPTO 2023, Part II*, volume 14082 of *LNCS*, pages 759–790. Springer, Heidelberg, Aug. 2023. (Cited on 6.)
- [48] A. Lenstra, X. Wang, and B. de Weger. Colliding x.509 certificates. Cryptology ePrint Archive, Report 2005/067, 2005. <https://eprint.iacr.org/2005/067>. (Cited on 25.)
- [49] A. K. Lenstra and B. de Weger. On the possibility of constructing meaningful hash collisions for public keys. In C. Boyd and J. M. G. Nieto, editors, *ACISP 05*, volume 3574 of *LNCS*, pages 267–279. Springer, Heidelberg, July 2005. (Cited on 25.)
- [50] I. Mironov and N. Stephens-Davidowitz. Cryptographic reverse firewalls. In E. Oswald and M. Fischlin, editors, *EUROCRYPT 2015, Part II*, volume 9057 of *LNCS*, pages 657–686. Springer, Heidelberg, Apr. 2015. (Cited on 5.)
- [51] OpenSSL. <https://www.openssl.org/>. (Cited on 34.)
- [52] G. Persiano, D. H. Phan, and M. Yung. Anamorphic encryption: Private communication against a dictator. In O. Dunkelman and S. Dziembowski, editors, *EUROCRYPT 2022, Part II*, volume 13276 of *LNCS*, pages 34–63. Springer, Heidelberg, May / June 2022. (Cited on 6.)
- [53] P. Ravi, S. Bhasin, A. Chattopadhyay, Aikata, and S. S. Roy. Backdooring post-quantum cryptography: Kleptographic attacks on lattice-based KEMs. Cryptology ePrint Archive, Report 2022/1681, 2022. <https://eprint.iacr.org/2022/1681>. (Cited on 2, 5.)
- [54] A. Russell, Q. Tang, M. Yung, and H.-S. Zhou. Cliptography: Clipping the power of kleptographic attacks. In J. H. Cheon and T. Takagi, editors, *ASIACRYPT 2016, Part II*, volume 10032 of *LNCS*, pages 34–64. Springer, Heidelberg, Dec. 2016. (Cited on 2, 5.)
- [55] A. Russell, Q. Tang, M. Yung, and H.-S. Zhou. Generic semantic security against a kleptographic adversary. In B. M. Thuraisingham, D. Evans, T. Malkin, and D. Xu, editors, *ACM CCS 2017*, pages 907–922. ACM Press, Oct. / Nov. 2017. (Cited on 2, 5.)
- [56] A. Russell, Q. Tang, M. Yung, and H.-S. Zhou. Correcting subverted random oracles. In H. Shacham and A. Boldyreva, editors, *CRYPTO 2018, Part II*, volume 10992 of *LNCS*, pages 241–271. Springer, Heidelberg, Aug. 2018. (Cited on 14.)
- [57] D. Shumow and N. Ferguson. On the possibility of a back door in the NIST sp800-90 dual EC PRNG. *CRYPTO’07 Rump Session*, 2007. <https://rump2007.cr.yp.to/15-shumow.pdf>. (Cited on 5.)

- [58] G. J. Simmons. The prisoners’ problem and the subliminal channel. In D. Chaum, editor, *CRYPTO’83*, pages 51–67. Plenum Press, New York, USA, 1983. (Cited on 6.)
- [59] G. J. Simmons. The subliminal channel and digital signature. In T. Beth, N. Cot, and I. Ingemarsson, editors, *EUROCRYPT’84*, volume 209 of *LNCS*, pages 364–378. Springer, Heidelberg, Apr. 1985. (Cited on 6.)
- [60] M. Stevens, A. K. Lenstra, and B. de Weger. Chosen-prefix collisions for MD5 and colliding X.509 certificates for different identities. In M. Naor, editor, *EUROCRYPT 2007*, volume 4515 of *LNCS*, pages 1–22. Springer, Heidelberg, May 2007. (Cited on 4, 25.)
- [61] M. Stevens, A. Sotirov, J. Appelbaum, A. K. Lenstra, D. Molnar, D. A. Osvik, and B. de Weger. Short chosen-prefix collisions for MD5 and the creation of a rogue CA certificate. In S. Halevi, editor, *CRYPTO 2009*, volume 5677 of *LNCS*, pages 55–69. Springer, Heidelberg, Aug. 2009. (Cited on 4, 25.)
- [62] T. Tiemann, S. Berndt, T. Eisenbarth, and M. Liskiewicz. “Act natural!”: Having a private chat on a public blockchain. Cryptology ePrint Archive, Report 2021/1073, 2021. <https://eprint.iacr.org/2021/1073>. (Cited on 2, 5, 20.)
- [63] P. R. Tiwari and M. Green. Subverting cryptographic hardware used in blockchain consensus. In J. Clark and E. Shi, editors, *FC 2024*, 2024. (Cited on 2, 5, 17.)
- [64] F. Valsorda. Bluesky post, 30 Mar. 2024. <https://bsky.app/profile/did:plc:x2nsupeeo52oznrmplwapppl/post/3kowjkkx2njy2b>. (Cited on 6.)
- [65] X. Wang and H. Yu. How to break MD5 and other hash functions. In R. Cramer, editor, *EUROCRYPT 2005*, volume 3494 of *LNCS*, pages 19–35. Springer, Heidelberg, May 2005. (Cited on 25.)
- [66] A. Young and M. Yung. The dark side of “black-box” cryptography, or: Should we trust capstone? In N. Kobitz, editor, *CRYPTO’96*, volume 1109 of *LNCS*, pages 89–103. Springer, Heidelberg, Aug. 1996. (Cited on 2, 5.)
- [67] A. Young and M. Yung. Kleptography: Using cryptography against cryptography. In W. Fumy, editor, *EUROCRYPT’97*, volume 1233 of *LNCS*, pages 62–74. Springer, Heidelberg, May 1997. (Cited on 2, 5.)
- [68] A. Young and M. Yung. The prevalence of kleptographic attacks on discrete-log based cryptosystems. In B. S. Kaliski Jr., editor, *CRYPTO’97*, volume 1294 of *LNCS*, pages 264–276. Springer, Heidelberg, Aug. 1997. (Cited on 5.)
- [69] G. Zaverucha and D. Shumow. Are certificate thumbprints unique? Cryptology ePrint Archive, Report 2019/130, 2019. <https://eprint.iacr.org/2019/130>. (Cited on 4, 25.)

A Proof of Theorem 6.2

Proof of Theorem 6.2: Consider game G_0 of Figure 18. We rewrite the winning condition by splitting it into three checks. Similar to the proof of Theorem 5.2, we already include line 11 which sets a flag `bad`. We have

$$\text{Adv}_{\text{PS}, \tilde{V}}^{\text{pfe}}(A) = \Pr [G_0(A)].$$

<p>Games $G_0, \boxed{G_1}$</p> <p>INIT:</p> <ol style="list-style-type: none"> 1 $(p, v) \leftarrow \text{PS}$; $(vk, sk) \leftarrow \text{S.Kg}$ 2 Define \tilde{v} as in Fig. 11 (using vk selected above) 3 $\mathcal{X} \leftarrow \emptyset$ 4 Return (p, v, \tilde{v}) <p>SIM(ϕ'):</p> <ol style="list-style-type: none"> 5 $\sigma \leftarrow \text{S.Sign}(sk, (1, \phi'))$ 6 $x \leftarrow \text{Emb}_{\text{PS}}((1 \parallel \sigma), \phi')$ 7 $\mathcal{X} \leftarrow \mathcal{X} \cup \{\phi'\}$ 8 Return x <p>FIN(ϕ^*, π^*):</p> <ol style="list-style-type: none"> 9 If $\phi^* \in \mathcal{X}$ then return false 10 If $\phi^* \in L_R$ then return false 11 If $(v(\phi^*, \pi^*) = 0 \wedge \tilde{v}(\phi^*, \pi^*) = 1)$ then bad \leftarrow true ; return false 12 Return $v(\phi^*, \pi^*)$
--

Figure 18: Games G_0, G_1 for the proof of Theorem 6.2. G_1 contains the boxed code and G_0 does not. Oracle SIM is as in Figure 12.

<p><u>Adversary $A_S(vk)$:</u></p> <ol style="list-style-type: none"> 1 $(p, v) \leftarrow \text{PS}$ 2 Define \tilde{v} as in Fig. 11 (using the input vk) 3 $\mathcal{X} \leftarrow \emptyset$; $\mathcal{Q}_S \leftarrow \emptyset$ 4 $(\phi^*, \pi^*) \leftarrow A[\text{SIM}_S](p, v, \tilde{v})$ 5 $x^* \leftarrow (\phi^*, \pi^*)$ 6 $((1 \parallel \sigma), \phi) \leftarrow \text{EmbInVPS}(x^*)$ 7 Return $(m = (1, \phi), \sigma)$ <p><u>Oracle $\text{SIM}_S(\phi')$:</u></p> <ol style="list-style-type: none"> 8 $\sigma \leftarrow \text{SIGN}((1, \phi'))$ 9 $x \leftarrow \text{Emb}_{\text{PS}}((1 \parallel \sigma), \phi')$ 10 $\mathcal{X} \leftarrow \mathcal{X} \cup \{\phi'\}$ 11 $\mathcal{Q}_S \leftarrow \mathcal{Q}_S \cup \{((1, \phi'), \sigma)\}$ 12 Return x 	<p><u>Adversary $A_{\text{PS}}(p, v)$:</u></p> <ol style="list-style-type: none"> 1 $(\tilde{v}, e) \leftarrow \tilde{\text{V}}(v)$; $\mathcal{X} \leftarrow \emptyset$ 2 $(\phi^*, \pi^*) \leftarrow A[\text{SIM}_{\text{PS}}](p, v, \tilde{v})$ 3 Return (ϕ^*, π^*) <p><u>Oracle $\text{SIM}_{\text{PS}}(\phi')$:</u></p> <ol style="list-style-type: none"> 4 $(\phi, \pi) \leftarrow e(\phi', 1)$; $\mathcal{X} \leftarrow \mathcal{X} \cup \{\phi'\}$ 5 Return π
---	--

Figure 19: Adversaries A_S (left) and A_{PS} (right) for the proof of Theorem 6.2.

We next turn to game G_1 which outputs false whenever **bad** is set; that is, when $\tilde{v}(\phi^*, \pi^*)$ is true, but $v(\phi^*, \pi^*)$ is not. Therefore, games G_0, G_1 are identical-until-bad and we have

$$\begin{aligned} \Pr[G_0(A)] &= \Pr[G_1(A)] + (\Pr[G_0(A)] - \Pr[G_1(A)]) \\ &\leq \Pr[G_1(A)] + \Pr[G_1(A) \text{ sets bad}] . \end{aligned}$$

We now construct adversaries A_S, A_{PS} such that the following two equations hold:

$$\Pr[G_1(A) \text{ sets bad}] \leq \text{Adv}_{\text{S}}^{\text{suf-cma}}(A_S) \tag{12}$$

$$\Pr[G_1(A)] \leq \text{Adv}_{\text{PS}}^{\text{snd}}(A_{\text{PS}}) , \tag{13}$$

which will complete the proof of Eq. (10) in the theorem statement.

We begin with A_S which is in game $\mathbf{G}_S^{\text{suf-cma}}$ and operates according to the description in Figure 19. A_S runs A , responding to oracle queries via SIM_S . These follow the responses in game G_1 , with A_S running its own SIGN oracle on line 8. These queries to its own oracle are tracked as message-signature pairs in set \mathcal{Q}_S . Now suppose that **bad** is set in G_1 .

Then $\phi^* \notin \mathcal{X}$, $v(\phi^*, \pi^*) = 0$ and $\tilde{v}(\phi^*, \pi^*) = 1$. In particular, $\phi^* \notin \mathcal{X}$ implies that $(1, \phi^*)$ is not in any message-signature pair in set \mathcal{Q}_S (else both would have been added on lines 10-11 of Figure 19). Moreover, since $\tilde{v}(\phi^*, \pi^*) = 1$ the parsed signature σ is a valid signature such that $\text{S.Vfy}(vk, (1, \phi^*), \sigma)$. A_S thus outputs this winning pair $((1, m^*), \sigma)$ in its suf-cma game. This justifies Eq. (12). Further, if A issues q queries to its simulation oracle SIM_S , then A_S also issues q queries to SIGN .

We next turn to adversary A_{PS} which is in game $\mathbf{G}_{PS}^{\text{snd}}$ and is depicted on the right side of Figure 19. A 's view is that of game G_1 ; initialization and SIM_{PS} return the same responses as in G_1 . Now, if $G_1(A)$ returns **true**, and since the boxed code is executed in G_1 , then it must be that the same conditions are satisfied as in $\mathbf{G}_{PS}^{\text{snd}}$ and A returns a winning output which proves Eq. (13).

We conclude the proof by observing that A_S and A_{PS} maintain running times close to that of A . \blacksquare

B Proof of Theorem 7.2

Proof of Theorem 7.2: Consider game G_0 of Figure 20. We claim that G_0 is equivalent to game $\mathbf{G}_{TS, \tilde{V}, n}^{\text{ffe}}$ when instantiated with $\tilde{V} = \mathbf{ASA}_n[\mathbf{V}, \mathbf{S}, \text{Emb}_{\text{sig}}]$. Note that the INIT and SIGN oracles are running the same steps, only with additional accounting. The ESIGN oracle in G_0 has replaced the computation $e(i, m', 1)$ with its execution as per the definition of \tilde{V} , and the FIN oracle of G_0 returns true if $(i^*, m^*) \notin \mathcal{X}$ and if (including the boxed code) either of $\text{S.Vfy}(vk, (i^*, 1, m^*), \sigma_S) = 1$ or $v_{i^*}(m^*, \sigma^*) = 1$. Recall that FIN of $\mathbf{G}_{TS, \tilde{V}, n}^{\text{ffe}}$ returns true if $(i^*, m^*) \notin \mathcal{X}$ and $\tilde{v}_{i^*}(m^*, \sigma^*) = 1$. The definition of \tilde{v}_i shows that $\tilde{v}_{i^*}(m^*, \sigma^*) = 1$ precisely when either $\text{S.Vfy}(vk, (i^*, 1, m^*), \sigma_S) = 1$ or $v_{i^*}(m^*, \sigma^*) = 1$. This completes the justification that

$$\mathbf{Adv}_{TS, \tilde{V}, n}^{\text{ffe}}(A) = \Pr[G_0(A)].$$

Since games G_0, G_1 are identical-until-**bad**, the Fundamental Lemma of Game Playing [11] implies

$$\begin{aligned} \Pr[G_0(A)] &= \Pr[G_1(A)] + (\Pr[G_0(A)] - \Pr[G_1(A)]) \\ &\leq \Pr[G_1(A)] + \Pr[G_1(A) \text{ sets bad}]. \end{aligned}$$

We now construct adversaries A_S, A_{TS} such that the following two equations hold:

$$\Pr[G_1(A)] \leq \mathbf{Adv}_S^{\text{suf-cma}}(A_S) \tag{14}$$

$$\Pr[G_1(A) \text{ sets bad}] \leq \mathbf{Adv}_{TS, n}^{\text{mu-ufcma}}(A_{TS}), \tag{15}$$

which will complete the proof of Eq. (11) in the theorem statement.

We begin with A_S which is in game $\mathbf{G}_S^{\text{suf-cma}}$ and operates according to the description in Figure 21. A_S runs A , responding to oracle queries via ESIGN_S and SIGN_S . These follow the responses in game G_1 , with A_S running its own SIGN oracle on line 8 in responses to ESIGN_S queries. These queries to its own oracle are tracked as message-signature pairs in set \mathcal{Q}_S . Now suppose that $G_1(A)$ returns

```

Games  $\boxed{G_0}, G_1$ 

INIT:
1 For  $i = 1, 2, \dots, n$  do  $(s_i, v_i) \leftarrow \text{TS}$ 
2  $(vk, sk) \leftarrow \text{S.Kg}$ 
3 Define  $(\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_n)$  as in Fig. 15 (using  $vk$  selected above)
4  $\mathcal{X} \leftarrow \emptyset$ ;  $\mathcal{Q}_S \leftarrow \emptyset$ 
5 Return  $((v_1, v_2, \dots, v_n), (\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_n))$ 

ESIGN( $i, m'$ ):
6  $\sigma_S \leftarrow \text{S.Sign}(sk, (i, 1, m'))$ ;  $x \leftarrow \text{Emb}_{\text{sig}}((1 \parallel \sigma_S), m')$ 
7  $\mathcal{X} \leftarrow \mathcal{X} \cup \{(i, m')\}$ ;  $\mathcal{Q}_S \leftarrow \mathcal{Q}_S \cup \{((i, 1, m'), \sigma_S)\}$ 
8 Return  $x$ 

SIGN( $i, m$ ):
9  $\sigma_{\text{TS}} \leftarrow s_i(m)$ 
10  $\mathcal{X} \leftarrow \mathcal{X} \cup \{(i, m)\}$ 
11 Return  $\sigma_{\text{TS}}$ 

FIN( $i^*, m^*, \sigma^*$ ):
12 If  $((i^*, m^*) \in \mathcal{X})$  then return false
13  $x^* \leftarrow (m^*, \sigma^*)$ 
14  $((1 \parallel \sigma_S), m) \leftarrow \text{Embln}_{\text{sig}}(x^*)$ 
15  $r_1 \leftarrow \text{S.Vfy}(vk, (i^*, 1, m^*), \sigma_S)$ 
16  $r_2 \leftarrow v_{i^*}(m^*, \sigma^*)$ 
17 If  $(r_1 = 0 \wedge r_2 = 1)$  then bad  $\leftarrow$  true;  $\boxed{\text{return true}}$ 
18 Return  $(r_1 = 1)$ 

```

Figure 20: Games G_0, G_1 for the proof of Theorem 7.2. G_0 contains the boxed code and G_1 does not.

```

Adversary  $A_S(vk)$ :
1 For  $i = 1, 2, \dots, n$  do  $s_i, v_i \leftarrow \text{TS}$ 
2 Define  $(\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_n)$  as in Fig. 15 (using the provided input  $vk$ )
3  $\mathcal{X} \leftarrow \emptyset$ ;  $\mathcal{Q}_S \leftarrow \emptyset$ 
4  $(i^*, m^*, \sigma^*) \leftarrow A[\text{ESIGN}_S, \text{SIGN}_S]((v_1, v_2, \dots, v_n), (\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_n))$ 
5  $x^* \leftarrow (m^*, \sigma^*)$ 
6  $((1 \parallel \sigma_S), m) \leftarrow \text{Embln}_{\text{sig}}(x^*)$ 
7 Return  $((i^*, 1, m^*), \sigma_S)$ 

Oracle  $\text{ESIGN}_S(i, m')$ :
8  $\sigma_S \leftarrow \text{SIGN}((i, 1, m'))$ ;  $x \leftarrow \text{Emb}_{\text{sig}}((1 \parallel \sigma_S), m')$ 
9  $\mathcal{X} \leftarrow \mathcal{X} \cup \{(i, m')\}$ ;  $\mathcal{Q}_S \leftarrow \mathcal{Q}_S \cup \{((i, 1, m'), \sigma_S)\}$ 
10 Return  $x$ 

Oracle  $\text{SIGN}_S(i, m)$ :
11  $\sigma_{\text{TS}} \leftarrow s_i(m)$ 
12  $\mathcal{X} \leftarrow \mathcal{X} \cup \{(i, m)\}$ 
13 Return  $\sigma_{\text{TS}}$ 

```

Figure 21: Adversary A_S for the proof of Theorem 7.2.

true. Then $(i^*, m^*) \notin \mathcal{X}$ and $r_1 = 1$. In particular, $(i^*, m^*) \notin \mathcal{X}$ implies that $(i^*, 1, m^*)$ is not in any message-signature pair in set \mathcal{Q}_S (else both would have been added on line 9 of Figure 21). Moreover,

<p><u>Adversary $A_{\text{T}\mathcal{S}}(v_1, v_2, \dots, v_n)$:</u></p> <ol style="list-style-type: none"> 1 $(vk, sk) \leftarrow \text{S.Kg}$ 2 Define $(\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_n)$ as in Fig. 15 (using vk selected above) 3 $\mathcal{X} \leftarrow \emptyset$ 4 $(i^*, m^*, \sigma^*) \leftarrow A[\text{ESIGN}_{\text{T}\mathcal{S}}, \text{SIGN}_{\text{T}\mathcal{S}}]((v_1, v_2, \dots, v_n), (\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_n))$ 5 Return (i^*, m^*, σ^*) <p><u>Oracle $\text{ESIGN}_{\text{T}\mathcal{S}}(i, m')$:</u></p> <ol style="list-style-type: none"> 6 $\sigma_{\mathcal{S}} \leftarrow \text{S.Sign}(sk, (i, 1, m'))$; $x \leftarrow \text{Emb}_{\text{sig}}((1 \parallel \sigma_{\mathcal{S}}), m')$ 7 $\mathcal{X} \leftarrow \mathcal{X} \cup \{(i, m')\}$ 8 Return x <p><u>Oracle $\text{SIGN}_{\text{T}\mathcal{S}}(i, m)$:</u></p> <ol style="list-style-type: none"> 9 $\sigma_{\text{T}\mathcal{S}} \leftarrow \text{SIGN}(i, m)$ 10 $\mathcal{X} \leftarrow \mathcal{X} \cup \{(i, m)\}$ 11 Return $\sigma_{\text{T}\mathcal{S}}$
--

Figure 22: Adversary $A_{\text{T}\mathcal{S}}$ for the proof of Theorem 7.2.

since $r_1 = 1$ the parsed signature $\sigma_{\mathcal{S}}$ is a valid signature such that $\text{S.Vfy}(vk, (i^*, 1, m^*), \sigma_{\mathcal{S}})$. $A_{\mathcal{S}}$ thus outputs this winning pair $((i^*, 1, m^*), \sigma_{\mathcal{S}})$ in its suf-cma game. This justifies Eq. (14).

Next we turn to adversary $A_{\text{T}\mathcal{S}}$ which is in game $\mathbf{G}_{\text{T}\mathcal{S},n}^{\text{mu-ufcma}}$ and is depicted in Figure 22. $A_{\text{T}\mathcal{S}}$ runs A , responding to oracle queries according to $\text{ESIGN}_{\text{T}\mathcal{S}}$ and $\text{SIGN}_{\text{T}\mathcal{S}}$. Oracle $\text{ESIGN}_{\text{T}\mathcal{S}}$ is the same as in game G_1 while oracle $\text{SIGN}_{\text{T}\mathcal{S}}$ calls $A_{\text{T}\mathcal{S}}$'s own SIGN oracle to match G_1 . We claim that if the flag `bad` is set during this execution (as per line 17 of G_1) then $A_{\text{T}\mathcal{S}}$ returns a valid forgery in game $\mathbf{G}_{\text{T}\mathcal{S},n}^{\text{mu-ufcma}}$. If `bad` is set then we have $(i^*, m^*) \notin \mathcal{X}$, $r_1 = 0$ and $r_2 = 1$. The latter means that $v_{i^*}(m^*, \sigma^*)$ and thus (i^*, m^*, σ^*) is a winning forgery for $A_{\text{T}\mathcal{S}}$, proving Eq. (15).

We complete the proof of the theorem statement by noting that $A_{\mathcal{S}}$ makes q_e SIGN queries while $A_{\text{T}\mathcal{S}}$ makes q_s SIGN queries, as needed. ■

C Forged certificate embedding

In Section 8 we gave certificate forgery as an example of a realistic embedding scenario. An ASA could be mounted either on the hash function or on the signature verification steps of certificate validation. Recall that effectiveness of an ASA on a hash function is defined relative to a predicate P_{cert} , for which our construction assumes a correct embedding function. This predicate is defined in pseudocode in Figure 17 but we here give a concrete example. Moreover, effectiveness of an ASA on signature verification already allows for *any* target message to be signed, as we cover in Section 7. Nonetheless, feasibility of an attack where forgeries “look innocuous” may ask for more sophisticated embeddings, so we provide this illustration. We focus on the hash function case.

Suppose we have selected target hash $y = 680f8b1123be39f4451430d6267a8159033034403ce0df1abdf11c105031d719$. This corresponds to a public certificate C with a valid signature $C.\text{sigValue}$ where $C.\text{sigAlg}$ specifies “PKCS #1 SHA-256 With RSA Encryption.” The aim is now to construct $C^* \neq C$ with the same hash y ; thus the signature on C can be reused.

An ASA attacker, following the specification of Section 8, does the following. Suppose they intend to use $\hat{H} = \text{ASA}[H_{\text{ca}}, S, \text{Emb}_{\text{cert}}]$ where H_{ca} is SHA256, S is ECDSA over `secp256k1`, and the embedding is as in Figure 17 of Section 8. We chose ECDSA because it is easily used in OpenSSL [51]; however a shorter signature scheme like BLS [16] may be easier to embed.

```

Certificate C*:
1 Data:
2   Version: 3 (0x2)
3   Serial Number: ...
4   Signature Algorithm: PKCS #1 SHA-256 With RSA Encryption
5   Issuer: C.Issuer
6   Validity
7     Not Before: Jan 1 08:00:00 2024 GMT
8     Not After : Jun 1 08:00:00 2024 GMT
9   Subject: O = Big Brother, CN = *.bigbrother.com
10  Subject Public Key Info: ...
11  X509v3 extensions:
12    X509v3 Basic Constraints: critical
13      CA:TRUE
14    fh: 680f8b1123be39f4451430d6267a8159033034403ce0df1abdf11c105031d719
15    fσ: 304502202b978f95a853dfa2d2574ff9...56ff5dbdeed8948eb7570089e12d5
16  Signature Algorithm: PKCS #1 SHA-256 With RSA Encryption
17  C.sigValue

```

Figure 23: A certificate forgery for an ASA on SHA256; we want $\tilde{h}(C^*) = \tilde{h}(C) = \text{SHA256}(C)$ so that a signature on C can be reused for forgery C^* . The highlighted lines are the overhead in constructing the hash collision; that is, these are determined by C or cannot be arbitrarily chosen. The remainder of the certificate, and in particular the “...” sections may be arbitrarily set by the ASA attacker. On a user’s device, SHA256 is replaced by \tilde{h} .

The forgery proceeds as follows. First, the substitution is generated via $(\tilde{h}, e) \leftarrow_{\mathcal{S}} \tilde{H}(h)$, which in particular means generating (vk^*, sk^*) for ECDSA. (Recall that vk^* is hardcoded in \tilde{h} while sk^* is hardcoded in e .) We select:

```

vk* = 04d0722759460447f1719ac66a1734054651f7c557a96166583d686
      ad405ca9b6f5fe47a7e425a8722edfa13be606fcbe4053ecacb27f2
      b0bc3dd1e83152c9a8a3 .

```

Next $\text{tbsCert}'$ is chosen, which is the certificate data to be contained in the forgery. Section 8 discussed arbitrary data and in the Introduction, we discussed impersonation of a legitimate website by swapping out the public key in their certificate. For this example, we suppose that the attacker is aiming to forge arbitrary $\text{tbsCert}'$. Now they use e to find a preimage tbsCert^* of target hash y , where the constraint is that tbsCert^* is “close to” $\text{tbsCert}'$. Concretely, tbsCert^* adds two additional fields. In the first, $\text{tbsCert}^*.f_h = y$. In the second, $\text{tbsCert}^*.f_\sigma = \sigma$, where $\sigma \leftarrow_{\mathcal{S}} \text{S.Sign}(sk^*, (y, \text{tbsCert}'))$. For our chosen data and y , we find

```

σ = 304502202b978f95a853dfa2d2574ff9980a4351e7d6c9c4fcc0529
    d636c750fdf4c16a8022100efbb50c105df2a4766cfa94910d3a190
    19656ff5dbdeed8948eb7570089e12d5 .

```

Now the forgery is ready to be put together: it includes data tbsCert^* , signature C.sigValue , and algorithm specification “PKCS #1 SHA-256 With RSA Encryption” (where $h = \text{SHA256}$ is substituted by \tilde{h} on the user’s device.) This is shown in Figure 23.