

A preliminary version of this paper appears in the proceedings of the 22nd ACM Computer and Communications Security (CCS) Conference, 2015.

# Mass-surveillance without the State: Strongly Undetectable Algorithm-Substitution Attacks

MIHIR BELLARE<sup>1</sup>

JOSEPH JAEGER<sup>2</sup>

DANIEL KANE<sup>3</sup>

August 12, 2015

## Abstract

We present new algorithm-substitution attacks (ASAs) on symmetric encryption that improve over prior ones in two ways. First, while prior attacks only broke a sub-class of randomized schemes having a property called coin injectivity, our attacks break *all* randomized schemes. Second, while prior attacks are stateful, ours are stateless, achieving a notion of strong undetectability that we formalize. Together this shows that ASAs are an even more dangerous and powerful mass surveillance method than previously thought. Our work serves to increase awareness about what is possible with ASAs and to spur the search for deterrents and counter-measures.

---

<sup>1</sup> Department of Computer Science & Engineering, University of California San Diego, 9500 Gilman Drive, La Jolla, California 92093, USA. Email: [mihir@eng.ucsd.edu](mailto:mihir@eng.ucsd.edu). URL: <http://cseweb.ucsd.edu/~mihir/>. Supported in part by NSF grants CNS-1116800 and CNS-1228890 and a gift from Microsoft. This work was done in part while Bellare was visiting the Simons Institute for the Theory of Computing, supported by the Simons Foundation and by the DIMACS/Simons Collaboration in Cryptography through NSF grant CNS-1523467.

<sup>2</sup> Department of Computer Science & Engineering, University of California San Diego, 9500 Gilman Drive, La Jolla, California 92093, USA. Email: [jsjaeger@eng.ucsd.edu](mailto:jsjaeger@eng.ucsd.edu). URL: <http://cseweb.ucsd.edu/~jsjaeger/>. Supported in part by NSF grants CNS-1116800 and CNS-1228890.

<sup>3</sup> Department of Computer Science & Engineering, University of California San Diego, 9500 Gilman Drive, La Jolla, California 92093, USA. Email: [dakane@eng.ucsd.edu](mailto:dakane@eng.ucsd.edu). URL: <http://cseweb.ucsd.edu/~dakane/>.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Notation and Definitions</b>	<b>6</b>
<b>3</b>	<b>ASA Framework</b>	<b>7</b>
<b>4</b>	<b>Attack and Analysis</b>	<b>10</b>
<b>5</b>	<b>Defenses</b>	<b>16</b>

# 1 Introduction

The Snowden revelations have exposed a pervasive program of mass surveillance, one of whose potential mechanisms is an algorithm substitution attack (ASA) [19, 20, 3]. In the context of symmetric encryption which is our focus, your encryption code (which implements the prescribed encryption algorithm of the scheme) is replaced by big-brother created malware that aims to undetectably exfiltrate your key through ciphertexts [3].

In this paper we present a new ASA on symmetric encryption that improves over the prior one of BPR [3] in two ways. First, while the prior ASA only broke a sub-class of randomized schemes having a property called coin injectivity, ours breaks *all* randomized schemes. Second, while the prior ASA was stateful, ours is stateless, achieving a stronger notion of undetectability that we motivate and formalize. Together this shows that ASAs are even more dangerous and powerful than previously thought.

ASAs. Alice expects that her encryption code is implementing the encryption algorithm `SE.Enc` prescribed by her chosen symmetric encryption scheme `SE`. In an ASA [3], the code of `SE.Enc` has been replaced by big-brother designed malware `A.Enc`. When Alice calls her encryption routine, giving it key  $K_{SE}$  and message  $m$  to encrypt, it is `A.Enc` that runs, producing a ciphertext  $c$  based not only on the given inputs but also on another, hardwired key we denote  $K_A$ . Colluding with `A.Enc` is an accomplice `A.Ext` who knows  $K_A$  and will pick up  $c$  by eavesdropping on the channel. `A.Enc` aims to produce  $c$  in such a way that `A.Ext` can violate privacy of Alice’s communications, the most effective (from the point of view of big brother) attack being one that recovers  $K_{SE}$  from  $c$ .

One’s first reaction may be that a successful ASA is trivial. For example, let the subverted code `A.Enc`, given  $K_A, K_{SE}, m$ , set ciphertext  $c$  equal to the encryption key  $K_{SE}$ . BPR [3] argue that this is attack is unattractive to big brother because it is *detectable*. Anyone (in particular the decrypter), obtaining  $c$ , will see that something is not right. BPR [3] put forth the goal of an ASA as being to violate privacy while remaining undetectable.

UNDETECTABILITY. To rigorously explore ASAs, BPR [3] provide a formalization of undetectability. Their viewpoint is that the detector is the decrypter. It thus knows  $K_{SE}$  and can check that a ciphertext decrypts to an intended message. But it can do other checks as well. The definition models detector  $\mathcal{D}$  as an adversary given input  $K_{SE}$  and access to an encryption oracle `ENC` that, given a message  $m$ , returns ciphertext  $c$  computed either via `SE.Enc` on inputs  $K_{SE}, m$  or via `A.Enc` on inputs  $K_A, K_{SE}, m$ . Undetectability requires that no efficient  $\mathcal{D}$  should be able to tell which.

THE BPR ATTACK. The question that emerges is whether it is possible to mount a successful but undetectable ASA. BPR [3] address this in depth when `SE.Enc` is randomized. (Here and henceforth, a randomized scheme means a stateless scheme with non-trivial entropy in the ciphertexts. See Section 4 for a formal definition. A deterministic scheme, which has no entropy in the ciphertexts, does not qualify!) Randomized schemes are the most basic and common, making this question important. Their strongest result is represented by their *biased ciphertext attack*. We will first explain the attack and then discuss our work.

Let us view `SE.Enc` as a deterministic function of  $K_{SE}, m$  and coins  $r$ , producing a ciphertext as  $c \leftarrow \text{SE.Enc}(K_{SE}, m; r)$ . In the true encryption process,  $r$  is chosen at random. In the biased ciphertext attack of BPR [3], `A.Enc` also produces  $c$  as `SE.Enc`( $K_{SE}, m; r$ ) but picks  $r$ , not uniformly at random, but at random from a subset  $S$  of the set of all possible coins, where  $S$  is defined as the set of  $r$  so that a PRF  $F$  under  $K_A$ , when applied to input  $c = \text{SE.Enc}(K_{SE}, m; r)$ , returns the first bit  $K_{SE}[1]$  of  $K_{SE}$ . This allows `A.Ext`, given  $K_A, c$ , to retrieve  $K_{SE}[1] = F(K_A, c)$ . The process is repeated to exfiltrate  $K_{SE}$  bit by bit.

The difficulty is assessing undetectability. The ciphertext  $c$  produced by `A.Enc` is not distributed the same way as one produced by the true encryption process. Can a detection adversary  $\mathcal{D}$  tell? This question turns out to be surprisingly difficult. BPR [3] provide a partial answer. Calling a scheme `SE` coin injective if the mapping  $r \mapsto \text{SE.Enc}(K_{\text{SE}}, m; r)$  is injective for all  $K_{\text{SE}}, m$ , they use a combinatorial analysis [3, Lemma 1] to prove undetectability in this case. Based on this they move away from randomized, stateless schemes to obtain security via deterministic, stateful ones.

SHORTCOMINGS OF PRIOR WORK. The biased ciphertext attack has two shortcomings. Both are noted in BPR [3] and left as open problems which we will resolve.

The first shortcoming is that the attack only applies to coin injective schemes. The basic question left un-answered here is, does there exist an ASA-secure randomized scheme, or are *all* randomized schemes subject to attack?

The second shortcoming is that in the biased ciphertext ASA, the subverted encryption algorithm `A.Enc` is *stateful*, maintaining state  $\sigma$  across invocations. (The state in this case is an integer representing either which bit of the key `A.Enc` is trying to exfiltrate or taking a special value to indicate that exfiltration is complete and encryption should be done as usual.) This renders the attack detectable in practice in simple ways. For example, a state reset, as can happen with a reboot or cloning to create a virtual machine, leads, in their attack, to detection. However, this is not captured by their definition of detectability, under which they prove their attack undetectable (for coin injective schemes). BPR [3] recognize this and comment that a stateless attack would be better, but that they do not know how to make their attack stateless.

CONTRIBUTIONS IN BRIEF. We contribute (1) new definitions (2) new attacks and (3) new analyses. The definitions are for strong undetectability (which captures the above issues and can only be met by stateless schemes) and key recovery security, the latter reflecting the need to formalize not just security goals but attack goals. The attack is a new ASA shown to achieve both. The analyses establishing this resolve technical issues from BPR [3] via a different approach. As a by product, we resolve the above questions, presenting a stateless ASA that breaks all randomized encryption schemes and showing that the move to deterministic schemes for security, made by BPR [3], is a necessary one.

STRONG UNDETECTABILITY. Our first contribution is to introduce and formalize a stronger notion of undetectability called strong undetectability. Our formalization, in Section 3, uses the framework of the above discussed definition of BPR [3] but makes the oracle `ENC` more powerful with regard both to inputs and outputs. It now takes input not just a message, but a key, allowing the detector to pick, rather than merely know, the encryption key  $K_{\text{SE}}$ . It returns not just the ciphertext, but the current state  $\sigma$  of the encryption code.

The last means that as long as the good encryption algorithm `SE.Enc` is stateless (the case we are interested in here), statelessness of `A.Enc` is a necessary (but not sufficient) condition for strong undetectability. That is, a strongly undetectable ASA must be stateless. A consequence of this is that the ASA will not be detectable upon system reset or cloning, meaning strong undetectability excludes the easy and natural methods of detection allowed by undetectability under BPR [3].

Ensuring this was the main purpose of the definition. Beyond that, while BPR [3] took the view that detection is performed by the decrypter, strong undetectability suggests that detection may to some extent be performed by the encryptor. It imagines that the detector may have blackbox access to the encryption code and can experiment with it, so that it can feed it keys of its choice and see whether or not it maintains state. However one must be careful to not take this interpretation to an extreme, for there are certainly detection methods that someone with blackbox access to the code could employ that strong undetectability does not capture, for example timing the responses and comparing this to the time the real algorithm would be expected to take. Indeed, it is not

possible to evade *all* forms of detection. We are trying to evade more forms of detection than prior work and figure out how “undetectable” a successful ASA can be.

The biased ciphertext ASA of BPR [3] will not be strongly undetectable due to its use of non-trivial state. The question this raises is whether strongly undetectable (and in particular stateless) ASAs are possible, and on which randomized encryption schemes. We will show that they are possible on all randomized schemes.

KEY RECOVERY. Beyond introducing and formalizing strong undetectability, we take a more rigorous approach to attack goals. As is conventional in cryptography, BPR [3] formalize security goals, not attacker goals. Thus, they provide an indistinguishability style surveillance advantage, viewing the goal as being to ensure that this advantage stays small for undetectable ASAs. This is a good goal for the scheme designer to achieve, but its violation by an attacker means little. Big brother is after more, specifically recovery of the target key. We carefully formalize a framework for key recovery attacks in which we can compare different ones to differentiate between the effectiveness of different attacks and make rigorous claims about what our attacks achieve.

Our formalization introduces as a parameter a message distribution  $\mathcal{M}$ , distinct from the adversary, that represents the messages encrypted by the encrypter. The larger the class of message distributions for which the attack succeeds, the more powerful the attack. Our ASA will succeed for *all* message distributions, which means that big brother recovers the target key regardless of what messages the sender sends. In particular, big brother would not need to pick or control the messages to succeed.

ASA. We now sketch our new ASA which we will show to achieve key recovery (for all message distributions) and strong undetectability for all randomized encryption schemes. Let the good key  $K_{SE}$  have length  $k$ , and let PRF  $F$  take outputs in the set  $\{0, 1\} \times [k]$  where  $[k] = \{1, \dots, k\}$ , meaning an output is a pair  $(v, i)$  where  $v$  is a bit and  $i$  is an index into the key. In our new ASA, the subverted encryption algorithm  $A.Enc$ , given  $K_A, K_{SE}, m$ , aims to pick  $r$  such that the  $(v, i)$  returned by  $F(K_A, SE.Enc(K_{SE}, m; r))$  satisfies  $v = K_{SE}[i]$ , returning  $c \leftarrow SE.Enc(K_{SE}, m; r)$ .  $A.Ext$  can retrieve  $(K_{SE}[i], i) = F(K_A, c)$ . In this way it gets the  $i$ -th bit of the key.

However, there are some difficulties with the above description. First, it may not be possible to directly and efficiently pick  $r$  in the way described. In our ASA, which is presented in detail in Section 4,  $A.Enc$  uses rejection sampling, repeatedly picking  $r$  at random until either the desired condition is met or it times out, when the latter happens being controlled by a parameter  $s$  of the attack. Second, we must use a series of  $q$  encryptions to allow  $A.Ext$  to recover the entire key, so that we must determine how to pick  $q$ , which is another parameter of the attack. The challenge is to prove both key recovery and strong undetectability without making any assumptions on the encryption scheme other than that it is randomized, and with good, concrete bounds enabling concrete and practical choices of  $q, s$  that make the attack both efficient and effective. The analyses we discuss next resolve these challenges.

ANALYSES. Theorem 4.1 proves that big brother, with our ASA, will indeed recover the target key, with high probability even for relatively small values of the attack parameters  $q, s$ . Given  $q, s$  and the key length  $k$  of the targeted randomized encryption scheme  $SE$ , the theorem gives a concrete lower bound on the key recovery advantage of our adversary as a function of  $q, s, k$ . The proof is a sequence of games. We begin, in a standard way, by exploiting the assumed PRF security of  $F$  to move to a game where it is replaced by a random function. We then exploit the assumption that the scheme is randomized (ciphertexts have non-trivial min-entropy) to move to a game where the  $(v, i)$  values are all picked independently at random. This allow us to move to a game where the sampling continues infinitely, reaching the conceptually correct distribution of the above attack idea. A coupon collector argument provides the bound for the last game. Each game transition

accumulates an error, the sum of these being our final bound. We show with a numerical example that a 128 bit key is recovered quite quickly.

Our approach to proving (strong) undetectability is entirely different from that used by BPR [3] to prove undetectability of their biased ciphertext attack, and it is this novel analysis that allows us to avoid extra assumptions like coin injectivity on the encryption scheme SE. While they use a combinatorial analysis [3, Lemma 1], we use a game sequence. Theorem 4.2 proves strong undetectability of our ASA assuming only that SE is randomized (ciphertexts have non-trivial min entropy) and F is a PRF.

RELATED WORK. Simmons work on subliminal channels [15, 16, 17, 18] was an early indication of the danger of ASAs. ASAs are part of the broader framework of kleptography studied by Young and Yung [19, 20, 21, 23, 22]. Their ideas seem prescient now. Back-doored blockciphers were studied in [12, 10, 11]. Goh, Boneh, Pinkas and Golle [7] show how to add key recovery to the SSL/TLS and SSH protocols. Cryptographic reverse firewalls [9] represent an architecture to counter ASAs via trusted code in network perimeter filters. Dodis, Ganesh, Golovnev, Juels and Ristenpart [6] provide a formal treatment of backdooring of PRGs, another form of subversion. Russell, Tang, Yung and Zhou [13] consider ASAs on one-way and trapdoor one-way functions. The survey by Schneier, Fredrikson, Kohno and Ristenpart [14] takes a broader look at subversion, providing useful categorizations. In independent work that we will discuss in more depth after giving our definitions, Degabriele, Farshim and Poettering [5] critique and refine the definitions of BPR [3]. Ateniese, Magri and Venturi [1] study ASAs on signature schemes in particular giving stateful attacks based on the methods of BPR [3]. We note that our methods can be employed to make their attacks stateless as well.

## 2 Notation and Definitions

NOTATION. If  $n$  is an integer then we let  $[n] = \{1, \dots, n\}$ . If  $x$  is a string then  $|x|$  denotes its length while if  $S$  is a set then  $|S|$  denotes its size. By  $\varepsilon$  we denote the empty string. By  $x[i]$  we denote the  $i$ -th bit of a string  $x$ , for  $i \in [|x|]$ . By  $s \leftarrow_s S$  we denote picking  $s$  at random from set  $S$ . If  $A$  is a randomized algorithm then  $y \leftarrow A(x_1, \dots; r)$  denotes running  $A$  on inputs  $x_1, \dots$  and coins  $r$  to deterministically obtain output  $y$ , and  $y \leftarrow_s A(x_1, \dots)$  denotes picking  $r$  at random and letting  $y \leftarrow A(x_1, \dots; r)$ . Definitions and proofs use code-based games [4]. See Fig. 1 for an example of a game. If  $G$  is a game then  $\Pr[G]$  denotes the probability that it returns `true`. We adopt the convention that the running time of an adversary means the worst case execution time of the adversary in the game that executes it, so that time for game setup steps and time to compute answers to oracle queries is included.

PRFS. We recall the definition as per [2, 8]. Let  $F : \{0, 1\}^\ell \times \{0, 1\}^* \rightarrow R$  be a function taking a key  $L \in \{0, 1\}^\ell$  and input  $c \in \{0, 1\}^*$  to return an output  $F(L, c) \in R$ . Consider game  $\text{PRF}_F^{\mathcal{F}}$  associated to  $F$  and adversary  $\mathcal{F}$ . It provides the adversary with an oracle FN as shown. Let

$$\text{Adv}_F^{\text{prf}}(\mathcal{F}) = 2 \Pr[\text{PRF}_F^{\mathcal{F}}] - 1$$

be the prf advantage of adversary  $\mathcal{F}$  against function  $F$ . In proofs we will use the standard fact that it can also be expressed as

$$\text{Adv}_F^{\text{prf}}(\mathcal{F}) = \Pr[\text{PRF}_F^{\mathcal{F}} \mid b_{\text{prf}} = 1] - (1 - \Pr[\text{PRF}_F^{\mathcal{F}} \mid b_{\text{prf}} = 0]) . \quad (1)$$

<p><b>Game</b> <math>\text{PRF}_F^{\mathcal{F}}</math></p> <p><math>L \leftarrow_{\\$} \{0, 1\}^\ell; b_{\text{prf}} \leftarrow_{\\$} \{0, 1\}; C \leftarrow \emptyset</math></p> <p><math>b'_{\text{prf}} \leftarrow_{\\$} \mathcal{F}^{\text{FN}}</math></p> <p>Return <math>(b_{\text{prf}} = b'_{\text{prf}})</math></p> <p><u>FN</u>(<math>c</math>)</p> <p>If <math>(b_{\text{prf}} = 1)</math> then <math>y_c \leftarrow F(L, c)</math></p> <p>Else</p> <p style="padding-left: 20px;">If <math>c \notin C</math> then <math>y_c \leftarrow_{\\$} R</math></p> <p style="padding-left: 20px;"><math>C \leftarrow C \cup \{c\}</math></p> <p>Return <math>c</math></p>
--

Figure 1: Game used to define prf advantage of adversary  $\mathcal{F}$  against function  $F$ .

### 3 ASA Framework

We recall basic syntax of symmetric encryption. We then provide our novel definitions for ASAs, namely strong undetectability as well as adversary advantage in key recovery.

SYMMETRIC ENCRYPTION. A (symmetric) encryption scheme  $\text{SE}$  specifies the following. Via  $K_{\text{SE}} \leftarrow_{\$} \{0, 1\}^{\text{SE.kl}}$ , one selects a key of length  $\text{SE.kl}$ . Encryption algorithm  $\text{SE.Enc}$  takes  $K_{\text{SE}}$ , message  $m$  and coins  $r \in \{0, 1\}^{\text{SE.rl}}$  to deterministically obtain ciphertext  $c \leftarrow \text{SE.Enc}(K_{\text{SE}}, m; r)$ . By  $c \leftarrow_{\$} \text{SE.Enc}(K_{\text{SE}}, m)$  we denote the operation  $r \leftarrow_{\$} \{0, 1\}^{\text{SE.rl}}; c \leftarrow \text{SE.Enc}(K_{\text{SE}}, m; r)$ . Decryption algorithm  $\text{SE.Dec}$  is deterministic and we require that

$$\text{SE.Dec}(K_{\text{SE}}, \text{SE.Enc}(K_{\text{SE}}, m; r)) = m$$

for all  $K_{\text{SE}}, m, r$ .

ASAs. An *algorithm substitution attack* (ASA)  $\mathbf{A}$  specifies the following. Via  $K_{\mathbf{A}} \leftarrow_{\$} \{0, 1\}^{\mathbf{A.kl}}$ , one selects a key of length  $\mathbf{A.kl}$ . The subverted encryption algorithm  $\mathbf{A.Enc}$  takes  $K_{\mathbf{A}}, K_{\text{SE}}, m$  and current state  $\sigma$  to produce ciphertext  $c$  and updated state,  $(c, \sigma) \leftarrow_{\$} \mathbf{A.Enc}(K_{\mathbf{A}}, K_{\text{SE}}, m, \sigma)$ . (Denoting the state the same in input and output simply indicates an update of this variable.) The idea is that  $\mathbf{A}$  is specified by big brother. Key  $K_{\mathbf{A}}$  is shared between the subverted encryption algorithm  $\mathbf{A.Enc}$  and its external accomplice  $\mathbf{A.Ext}$  who aims, from subverted ciphertexts, to violate security of  $\text{SE}$ . Note that conceptually,  $\mathbf{A}$  is an adversary (the ASA) rather than a scheme, a departure in perspective from BPR [3]. We say that  $\mathbf{A.Enc}$  is stateless if the updated state it returns is always the empty string  $\varepsilon$ , meaning its output for any inputs  $K_{\mathbf{A}}, K_{\text{SE}}, m, \sigma$  has the form  $(c, \varepsilon)$ . In this case we might drop  $\sigma$  in both input and output, writing  $c \leftarrow_{\$} \mathbf{A.Enc}(K_{\mathbf{A}}, K_{\text{SE}}, m)$ .

We said that  $\mathbf{A}$ 's goal is to violate security of  $\text{SE}$ . With no further conditions, this is too easy, as explained in Section 1. The goal for big brother is an attack that is undetectable yet violates security of  $\text{SE}$ . We now turn to formalizing each of these components.

STRONG UNDETECTABILITY. Consider game  $\text{SDET}$  of Fig. 2 associated to encryption scheme  $\text{SE}$ , ASA  $\mathbf{A}$  and a detection adversary  $\mathcal{D}$ . It can obtain via oracle  $\text{ENC}$  encryptions of messages of its choice under keys of its choice, produced either via  $\text{SE.Enc}$  (when  $b = 1$ ) or via  $\mathbf{A.Enc}$  (when  $b = 0$ ), and aims from examination of these ciphertexts to determine  $b$ . Let

$$\text{Adv}_{\text{SE}, \mathbf{A}}^{\text{sdet}}(\mathcal{D}) = 2 \Pr[\text{SDET}_{\text{SE}, \mathbf{A}}^{\mathcal{D}}] - 1 .$$

This is  $\mathcal{D}$ 's advantage in detecting the subversion.

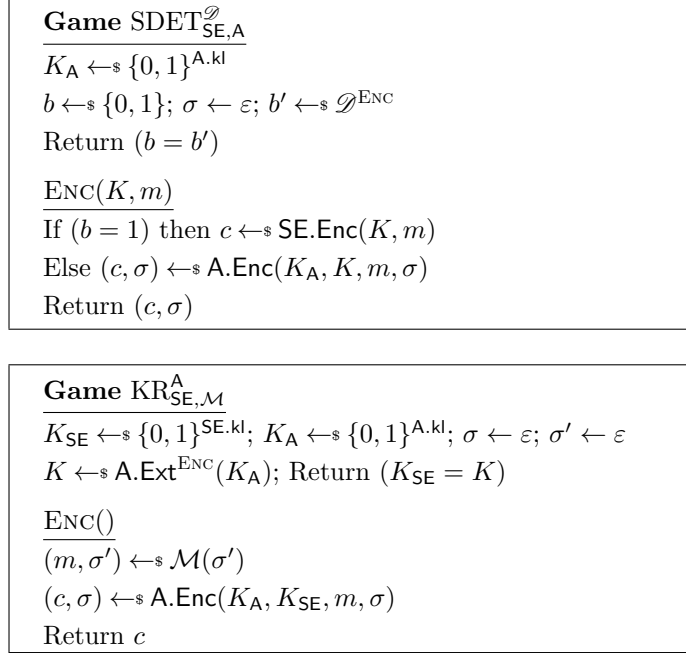


Figure 2: Games used to define detection and key recovery security of subversion A of encryption scheme SE.

The state  $\sigma$  returned by ENC will always be the empty string  $\varepsilon$  in the case  $b = 1$  because SE.Enc is stateless. Thus, if  $\mathcal{D}$  ever sees  $\sigma \neq \varepsilon$  in a reply to a ENC query, it knows that  $b = 0$ , meaning the subversion has been detected. Our formalization thus effectively implies that a subversion must be stateless to be undetectable. A consequence of this is that state reset, as can arise due to system reboot or cloning to create a virtual machine, will not allow the ASA to be detected, unlike for undetectability as per BPR [3]. This exclusion of some simple and natural forms of detection not covered by BPR [3] is the main consequence and intent of the new definition.

However there are other directions as well. The view of BPR [3] was that detection is performed by the decryptor. Detection through state reset continues that perspective, for this can in many cases lead to detection from the ciphertexts alone. But strong undetectability also moves towards a perspective where the encryptor, not just the decryptor, may be making some attempts at detection. We view the detector as having some sort of blackbox access to the encryption code, so that it can pick inputs and see all outputs written to memory. In particular, the detector not just knows, but can pick, the encryption key  $K$ , and it can see any state  $\sigma$  that the encryption code tries to maintain across invocations. Thus  $\mathcal{D}$  can supply, in calls to ENC, not just the message, but *any* key  $K \in \{0, 1\}^{\text{SE.kl}}$  of its choice for SE.Enc, getting in response not just the ciphertext  $c$ , but also  $\sigma$ . This obviously significantly increases the power of the detector.

However we must be careful to note that we do not capture *all* possible detection strategies that a “real” detector in such a position could mount. For example the detector could measure the CPU time of an execution of the code and compare this with the expected CPU time of the real code. Or it could look at the number of calls to the underlying pseudo-random number generator that is being used to obtain coins. These and other detection methods are not covered by our definition.

The fact is that it impossible for an ASA to evade *all* forms of detection. Our work aims to understand how far we can push the boundary. Security in this domain is a tradeoff between



detection effort and attack success. Our results indicate that detecting ASAs takes more effort than was previously thought.

We note that strong undetectability implies the BPR notion of undetectability. A proof may be given by reduction, the idea being as follows. Recall that in the detection game of BPR [3], the undetectability adversary  $\mathcal{D}'$  has access to an oracle KEY that, given  $i$ , returns a key  $K_i \leftarrow_s \{0, 1\}^{\text{SE.kl}}$  for user  $i$ , and an oracle ENC that, given message  $M$  and  $i$ , returns a ciphertext produced by running either SE.Enc or A.Enc to encrypt  $M$  under  $K_i$ , the choice depending on the challenge bit that  $\mathcal{D}'$  is trying to guess. Given such a  $\mathcal{D}'$ , we build a strong undetectability adversary  $\mathcal{D}$  with strong undetectability advantage at least the undetectability advantage of  $\mathcal{D}'$ . Our adversary  $\mathcal{D}$  runs  $\mathcal{D}'$ . When the latter makes a query KEY( $i$ ), adversary  $\mathcal{D}$  itself picks  $K_i \leftarrow_s \{0, 1\}^{\text{SE.kl}}$  and returns it to  $\mathcal{D}'$ . When the latter makes a query ENC( $M, i$ ), adversary  $\mathcal{D}$  queries its own ENC oracle with  $K_i, M$  to get back  $(c, \sigma)$  and returns  $c$  to  $\mathcal{D}'$ . Finally  $\mathcal{D}$  returns the same decision as  $\mathcal{D}'$ .

We note that the detectability game of BPR was explicitly multi-user, meaning it involved multiple keys. The strong undetectability game instead allows the adversary to query ENC with any key of its choice. One advantage of this formulation, as we have just illustrated, is that multi-user security is a consequence.

KEY RECOVERY. The other side of the coin is what it means for an ASA to succeed, meaning violate the security of the users of SE. One measure, formalized in BPR [3], is distinguishability, meaning A succeeds if it can distinguish encryptions of messages of its choice. Another measure is key recovery, meaning A succeeds if it finds  $K_{\text{SE}}$ . Distinguishing is a strong measure for security but a weak one for attacks (achieving it provides high security, but violating it entails little loss) while key recovery is a weak measure for security but a strong one for attacks (violating it is very damaging but achieving it means little for security). Since our focus is attacks, we target and formalize key recovery.

Game KR of Fig. 2 is parameterized by a message sampler algorithm  $\mathcal{M}$  that, given its current state  $\sigma'$  returns the next message  $m$  to be encrypted and updated state. It represents the choice of messages made by the sender. A wins if A.Ext recovers the key  $K_{\text{SE}}$  from the ciphertexts produced by the subverted encryption algorithm A.Enc on messages produced by  $\mathcal{M}$ . (The state  $\sigma$  maintained by the latter is entirely distinct from the state  $\sigma'$  of  $\mathcal{M}$ .) The key recovery advantage of A is

$$\text{Adv}_{\text{SE}, \mathcal{M}}^{\text{kr}}(\text{A}) = \Pr[\text{KR}_{\text{SE}, \mathcal{M}}^{\text{A}}].$$

Parametrization by  $\mathcal{M}$  allows a fine-grained taxonomy of key recovery attacks. The less they assume about  $\mathcal{M}$ , the stronger they are. The strongest attack is one that works for any  $\mathcal{M}$ . This corresponds to an attack that works regardless of what messages the encrypter chooses to encrypt. A somewhat weaker attack might work for certain message sequences, meaning some restricted class of samplers  $\mathcal{M}$ . Our ASA is of the stronger type.

DISCUSSION AND EXTENSIONS. Game SDET of Fig. 2 has been written for the case where encryption scheme SE is stateless, since this is the case of interest for our attacks. It can be extended to the case where SE is stateful. In this case, when  $b = 1$ , algorithm SE.Enc takes the current state  $\sigma$  as an additional input and returns not just ciphertext  $c$  but also an updated state  $\sigma$ . Additionally one should provide the adversary with a reset oracle RESET that resets the state  $\sigma$  to  $\varepsilon$ . The notion is thus requiring that the detector cannot tell whether it is talking to the real or subverted encryption algorithm even if it sees the state and can reset it. Note that the reset oracle is redundant (and hence has been omitted) in our present context because with stateless schemes the definition already implies that a subversion must be stateless as well, and with both stateless a reset oracle is vacuous.

BPR [3] had required any ASA  $A$  to satisfy a decryptability condition which asked that ciphertexts produced under  $A.\text{Enc}(\cdot, K, \cdot, \cdot)$  decrypted correctly under  $\text{SE}.\text{Dec}(K, \cdot)$ . We have dropped this condition, so that decryptability holds only to the extent that it is implied by strong undetectability, which we think is more realistic from a detection perspective. However, the ASAs we provide do meet the decryptability condition of BPR [3].

DFP [5] suggest that the BPR decryptability condition is too strong. They suggest a relaxation, and then provide an ASA that is undetectable under the BPR definition, meets their relaxed decryptability condition, and yet succeeds in that it violates the BPR subversion security definition. However the DFP attack does not succeed in violating our key recovery notion, and from a practical perspective, is weak. In the DFP attack,  $A.\text{Enc}(K_A, K, m, \sigma)$  returns  $K$  when  $m = K_A$ . That is, the attack requires that the attacker can induce the encryptor to encrypt the message  $m = K_A$ . But in practice this is quite hard and it is unlikely that a “real” sender will ever encrypt this message. Our attacks are much more powerful since the attacker succeeds regardless of what messages the encryptor encrypts. This is captured by the introduction of  $\mathcal{M}$  in our key recovery definition and the fact that the attack succeeds regardless of its choice. In fact in our key recovery definition, the messages being encrypted in the attack cannot even depend on the subverter’s key  $K_A$ , since the latter is not given to  $\mathcal{M}$ . The latter is why the DFP attack [5] does not succeed under our key recovery definition.

The issues here reflect, as we discussed above, that achieving a certain notion may provide good practical security, but violating the same notion may not constitute an effective practical attack. (That is, insecurity is not really the opposite of security. There is a lot in between.) Thus, achieving security under the BPR subversion notion is a good target for a scheme. But violating security under this notion need not be (and, in the case of DFP [5], is not) an effective attack. Conversely, our key recovery notion is a good one for attacks but not a good one for security: violating it constitutes a powerful attack, but achieving it provides only weak security. As an analogy, indistinguishability (semantic security) is accepted as a good security target for encryption schemes that one wants constructions to reach, but violating it is a weak attack.

## 4 Attack and Analysis

We present our ASA, having the following properties: (1) It is stateless and proven *strongly* undetectable (2) It breaks *any* given encryption scheme that has a non-trivial amount of randomization, and (3) break here is in the strong sense of key recovery for arbitrary message distributions. The attack is simple to specify but the analysis is more challenging. We provide a careful analysis to establish both key recovery and strong undetectability with concrete bounds. From these results we can extract concrete values of the parameters for a concrete attack.

PRELIMINARIES. As we discussed, we need SE to have non-trivial randomization. (Otherwise, a deterministic scheme is just a special case of a randomized scheme.) Formally, this means that ciphertexts have some min-entropy. To measure this we define the min-entropy  $\mathbf{H}_\infty(\text{SE})$  of the scheme SE via

$$2^{-\mathbf{H}_\infty(\text{SE})} = \max_{K_{\text{SE}}, m, c} \Pr[\text{SE}.\text{Enc}(K_{\text{SE}}, m; r) = c]$$

Here, with  $K_{\text{SE}}, m, c$  fixed, the probability is over a random choice of  $r$  from  $\{0, 1\}^{\text{SE}.r}$ . Our results will assume that  $2^{-\mathbf{H}_\infty(\text{SE})}$  is negligible. We note that we could use collision entropy in place of min entropy.

ATTACK DESCRIPTION. Let SE denote the target encryption scheme and let  $F: \{0, 1\}^{\text{F.kl}} \times \{0, 1\}^* \rightarrow$

```

A.Enc( $K_A, K_{SE}, m, \sigma$ )
 $j \leftarrow 0$ 
Repeat
   $j \leftarrow j + 1$ 
   $r \leftarrow_{\$} \{0, 1\}^{SE.rl}$ 
   $c \leftarrow SE.Enc(K_{SE}, m; r)$ 
   $(v, t) \leftarrow F(K_A, c)$ 
  success  $\leftarrow (K_{SE}[t] = v)$ 
  out-of-time  $\leftarrow (j = s)$ 
Until (success OR
out-of-time)
Return  $(c, \epsilon)$ 

A.Ext( $K_A$ )
 $K \leftarrow 0^{SE.kl}$ 
For  $i = 1, \dots, q$  do
   $c \leftarrow_{\$} ENC()$ 
   $(v, t) \leftarrow F(K_A, c)$ 
   $K[t] \leftarrow v$ 
Return  $K$ 

```

Figure 3: ASA A for encryption scheme SE has algorithms A.Enc and A.Ext as shown above. Here  $F: \{0, 1\}^{A.kl} \times \{0, 1\}^* \rightarrow \{0, 1\} \times [SE.kl]$  is a PRF used in the attack and  $q, s \geq 1$  are parameters of the attack.

$\{0, 1\} \times [SE.kl]$  be a PRF and key length  $A.kl = F.kl$ , meaning  $K_A$  will be a key for the PRF  $F$  defining the function  $f(\cdot) = F(K_A, \cdot)$ . The algorithms A.Enc and A.Ext are shown in Fig. 3. Here  $s, q \geq 1$  are parameters of the attack. (We will show that quite small values of these suffice.)

Subverted encryption algorithm A.Enc is given  $K_A, K_{SE}$ , message  $m$  and a state  $\sigma$  that it ignores. (It will be stateless and thus it will always be that  $\sigma = \epsilon$ .) Its goal is to pick  $r$  so that if  $c = SE.Enc(K_{SE}, m; r)$  and  $(v, t) = F(K_A, c)$  then  $K_{SE}[t] = v$ . It would then return  $c$ . When A.Ext picks up  $c$ , it can compute  $(v, t) = F(K_{SE}, c)$  and set  $K[t] = v$ . An appropriate choice of  $q$  will ensure that A.Ext eventually gets all bits of the key, meaning  $K = K_{SE}$  except with tiny probability.

Conceptually, then, we imagine A.Enc as trying to pick  $r$  at random subject to the constraint that if  $c = SE.Enc(K_{SE}, m; r)$  and  $(v, t) = F(K_A, c)$  then  $K_{SE}[t] = v$ . However it is not clear how to directly pick  $r$  in this way. Indeed, such an  $r$  may not even exist. To resolve this, A.Enc in Fig. 3 samples by picking  $r$  at random from the full space  $\{0, 1\}^{SE.rl}$  until either  $r$  satisfies the desired condition (at which point the flag `success` becomes true) or the process exceeds the number of allowed sampling attempts  $s$  (at which point the flag `out-of-time` becomes true). The tradeoff is that for efficiency we want  $s$  to be quite small but the smaller it is the farther is the distribution of  $r$  from the conceptually desired one. The theorems and analyses that follow will deal with these issues and show how to appropriately pick both  $s$  and  $q$  for an effective yet efficient attack.

**KEY RECOVERY.** The following theorem lower bounds the key recovery advantage of our ASA A of Fig. 3, showing that it is close to one for reasonable values of the parameters, meaning the attack can efficiently and successfully recover the target key.

<p><b>Game <math>\mathbb{H}_0, \mathbb{H}_1</math></b>  <math>K_{SE} \leftarrow_s \{0, 1\}^{SE.kl}; \sigma' \leftarrow \varepsilon</math>  <math>C \leftarrow \emptyset; T \leftarrow \emptyset</math>  For <math>i = 1, \dots, q</math> do  <math>j \leftarrow 0</math>; out-of-time <math>\leftarrow</math> false  <math>(m, \sigma') \leftarrow_s \mathcal{M}(\sigma')</math>  Repeat  <math>j \leftarrow j + 1</math>; <math>r \leftarrow_s \{0, 1\}^{SE.rl}</math>  <math>c \leftarrow SE.Enc(K_{SE}, m; r)</math>  <math>(v, t) \leftarrow_s \{0, 1\} \times [k]</math>  If <math>(c \in C)</math> then  <math>bad \leftarrow true</math>  <math>(v, t) \leftarrow (v_{c'}, t_{c'})</math>  <math>(v_{c'}, t_{c'}) \leftarrow (v, t)</math>  <math>C \leftarrow C \cup \{c\}</math>  success <math>\leftarrow (K_{SE}[t_c] = v_c)</math>  out-of-time <math>\leftarrow (j = s)</math>  Until (success OR out-of-time)  If success then <math>T \leftarrow T \cup \{t\}</math>  Else <math>T \leftarrow T \setminus \{t\}</math>  Return <math>(T \neq [k])</math></p>	<p><b>Adversary <math>\mathcal{F}</math></b>  <math>K_{SE} \leftarrow_s \{0, 1\}^{SE.kl}; \sigma' \leftarrow \varepsilon</math>  <math>T \leftarrow \emptyset</math>  For <math>i = 1, \dots, q</math> do  <math>j \leftarrow 0</math>; out-of-time <math>\leftarrow</math> false  <math>(m, \sigma') \leftarrow_s \mathcal{M}(\sigma')</math>  Repeat  <math>j \leftarrow j + 1</math>; <math>r \leftarrow_s \{0, 1\}^{SE.rl}</math>  <math>c \leftarrow SE.Enc(K_{SE}, m; r)</math>  <math>(v_c, t_c) \leftarrow_s FN(c)</math>  success <math>\leftarrow (K_{SE}[t_c] = v_c)</math>  out-of-time <math>\leftarrow (j = s)</math>  Until (success OR out-of-time)  If success then <math>T \leftarrow T \cup \{t\}</math>  Else <math>T \leftarrow T \setminus \{t\}</math>  If <math>(T \neq [k])</math> then return 1  Else return 0</p> <p><b>Game <math>\mathbb{H}_4</math></b>  <math>K_{SE} \leftarrow_s \{0, 1\}^{SE.kl}; T \leftarrow \emptyset</math>  For <math>i = 1, \dots, q</math> do  Repeat  <math>(v, t) \leftarrow_s \{0, 1\} \times [k]</math>  Until <math>(K_{SE}[t] = v)</math>  <math>T \leftarrow T \cup \{t\}</math>  Return <math>(T \neq [k])</math></p>	<p><b>Game <math>\mathbb{H}_2, \mathbb{H}_3</math></b>  <math>K_{SE} \leftarrow_s \{0, 1\}^{SE.kl}; T \leftarrow \emptyset</math>  For <math>i = 1, \dots, q</math> do  <math>j \leftarrow 0</math>; out-of-time <math>\leftarrow</math> false  Repeat  <math>j \leftarrow j + 1</math>  <math>(v, t) \leftarrow_s \{0, 1\} \times [k]</math>  success <math>\leftarrow (K_{SE}[t] = v)</math>  If <math>(j = s</math> AND success = false) then  <math>bad \leftarrow true</math>  <math>out-of-time \leftarrow true</math>  Until (success OR out-of-time)  If success then <math>T \leftarrow T \cup \{t\}</math>  Else <math>T \leftarrow T \setminus \{t\}</math>  Return <math>(T \neq [k])</math></p> <p><b>Game <math>\mathbb{H}_5</math></b>  <math>K_{SE} \leftarrow_s \{0, 1\}^{SE.kl}; T \leftarrow \emptyset</math>  For <math>i = 1, \dots, q</math> do  <math>t \leftarrow_s [k]; T \leftarrow T \cup \{t\}</math>  For <math>\ell = 1, \dots, k</math> do <math>bad_\ell \leftarrow (\ell \notin T)</math>  <math>bad \leftarrow (bad_1 \vee \dots \vee bad_k)</math>  Return bad</p>
---	---	---

Figure 4: Games and adversary for proof of Theorem 4.1.

We recall that our convention is that adversary running time refers to the time of the game in which the adversary executes, so that the time of oracle calls is included. Thus for example the reason the time for  $\mathcal{M}$  does not show up in the time of  $\mathcal{F}$  below is that KR runs  $\mathcal{M}$  already.

**Theorem 4.1** *Let SE be a symmetric encryption scheme and let  $k = SE.kl$ . Let  $F : \{0, 1\}^{F.kl} \times \{0, 1\}^* \rightarrow \{0, 1\} \times [k]$  be a PRF. Let  $q, s \geq 1$  and let  $\mathbf{A}$  be defined as in Fig. 3. Let  $\mathcal{M}$  be an arbitrary message distribution. Then we can build PRF adversary  $\mathcal{F}$  such that*

$$\mathbf{Adv}_{SE, \mathcal{M}}^{kr}(\mathbf{A}) \geq 1 - \mathbf{Adv}_F^{prf}(\mathcal{F}) - \epsilon(q, s, k) \quad (2)$$

where

$$\epsilon(q, s, k) \leq ke^{-q/k} + q2^{-s} + q^2s^2 \cdot 2^{-\mathbf{H}_\infty(SE)-1}. \quad (3)$$

The running time of  $\mathcal{F}$  is about the sum of the running times of  $\mathbf{A}.Enc$  and  $\mathbf{A}.Ext$ , and it makes at most  $qs$  oracle queries.

For example, the key for SE is typically an AES key so that  $k = SE.kl = 128$ . We could let  $q = 128 \cdot 7 = 896$  and  $s = 13$ , which implies  $ke^{-q/k} + q2^{-s} \leq 1/4$ . We can assume  $\mathbf{Adv}_F^{prf}(\mathcal{F})$  is negligible. Then as long as  $\mathbf{H}_\infty(SE) \geq 28$ , our attack has advantage around 1/2. The  $q, s$  values are quite small, making the attack quite practical.

That  $\mathcal{M}$  may be *any* distribution makes the ASA particularly strong. The distribution  $\mathcal{M}$  represents the choice of messages made by the encryptor. We are saying the attack works regardless of these choices. It is thus a known message attack, not a chosen message attack.

**Proof of of Theorem 4.1:** Game  $H_0$  of Fig. 4 includes the boxed code while  $H_1$  does not. Game  $H_0$  implements game KR of Fig. 2 with  $A$  being as defined in Fig. 3, with two changes: (1)  $F(K_A, \cdot)$  is replaced by a lazily sampled random function, and (2) the game returns `true` when key recovery fails rather than when it succeeds. Let prf adversary  $\mathcal{F}$  be as defined in the second panel of Fig. 4. It uses its FN oracle where  $A$  would use  $F(K_A, \cdot)$ . It returns 1 when key recovery fails and 0 when it succeeds. Letting  $b_{\text{prf}}$  denote the challenge bit in game PRF, we have

$$\begin{aligned} \Pr[\text{PRF}_{\mathbb{F}}^{\mathcal{F}} \mid b_{\text{prf}} = 1] &\geq 1 - \mathbf{Adv}_{\text{SE}}^{\text{kr}}(A) \\ 1 - \Pr[\text{PRF}_{\mathbb{F}}^{\mathcal{F}} \mid b_{\text{prf}} = 0] &= \Pr[H_0] . \end{aligned}$$

The reason the first equation above is an inequality rather than an equality is that  $A.\text{Ext}$  initially sets all bits of  $K$  to 0 and, so, for a particular bit, it may by chance end up having the right value even when that bit is not set in the `For` loop. Now, using Eq. (1), we have

$$\begin{aligned} \mathbf{Adv}_{\mathbb{F}}^{\text{prf}}(\mathcal{F}) &= \Pr[\text{PRF}_{\mathbb{F}}^{\mathcal{F}} \mid b_{\text{prf}} = 1] - (1 - \Pr[\text{PRF}_{\mathbb{F}}^{\mathcal{F}} \mid b_{\text{prf}} = 0]) \\ &\geq 1 - \mathbf{Adv}_{\text{SE}}^{\text{kr}}(A) - \Pr[H_0] . \end{aligned}$$

Re-arranging terms, we have

$$\mathbf{Adv}_{\text{SE}}^{\text{kr}}(A) \geq 1 - \mathbf{Adv}_{\mathbb{F}}^{\text{prf}}(\mathcal{F}) - \Pr[H_0] .$$

We let

$$\epsilon(q, s, k) = \Pr[H_0] .$$

This establishes Eq. (2). We now proceed to upper bound  $\epsilon(q, s, k)$ . Since games  $H_0, H_1$  are identical until `bad`, the fundamental lemma of game playing [4] says that

$$\epsilon(q, s, k) = \Pr[H_0] \leq \Pr[H_1] + \Pr[H_1 \text{ sets bad}] .$$

Game  $H_1$  sets `bad` when there is a collision in the ciphertexts. Since at most  $qs$  ciphertexts are created we have

$$\begin{aligned} \Pr[H_1 \text{ sets bad}] &\leq \binom{qs}{2} \cdot 2^{-\mathbf{H}_{\infty}(\text{SE})} \\ &\leq q^2 s^2 \cdot 2^{-\mathbf{H}_{\infty}(\text{SE})-1} . \end{aligned}$$

We now proceed to upper bound  $\Pr[H_1]$ . The ciphertexts being chosen in this game are not relevant to its outcome, since  $(v_c, t_c)$  is chosen at random each time. Furthermore, the setting of the `out-of-time` flag only has an effect on the behavior of the game if the `success` flag is not already set at the time. This leads to game  $H_2$  of Fig. 4, which includes the boxed code. We have

$$\Pr[H_1] = \Pr[H_2] \leq \Pr[H_3] + \Pr[H_3 \text{ sets bad}] ,$$

the inequality by the fundamental lemma of game playing [4] because games  $H_2, H_3$  are identical until `bad`. In game  $H_3$  the boxed code is not included, with the result that the repeat loop continues until success. Since each iteration is successful with probability  $1/2$  we have

$$\Pr[H_3 \text{ sets bad}] \leq q2^{-s} .$$

<p><b>Game I<sub>0</sub></b>  <math>b \leftarrow_s \{0, 1\}; C \leftarrow \emptyset</math>  <math>b' \leftarrow_s \mathcal{D}^{\text{ENC}}; \text{Return } (b = b')</math>  <u>ENC(K, m)</u>  If (b = 1) then  <math>c \leftarrow_s \text{SE.Enc}(K, m)</math>  Else  <math>j \leftarrow 0; \text{out-of-time} \leftarrow \text{false}</math>  Repeat  <math>j \leftarrow j + 1; r \leftarrow_s \{0, 1\}^{\text{SE.rl}}</math>  <math>c \leftarrow \text{SE.Enc}(K, m; r)</math>  <math>(v, t) \leftarrow_s \{0, 1\} \times [k]</math>  If (<math>c \in C</math>) then  <math>(v, t) \leftarrow (v_{c'}, t_{c'})</math>  <math>(v_c, t_c) \leftarrow (v, t)</math>  <math>C \leftarrow C \cup \{c\}</math>  <math>\text{success} \leftarrow (K[t_c] = v_c)</math>  <math>\text{out-of-time} \leftarrow (j = s)</math>  Until (success OR out-of-time)  Return (c, <math>\epsilon</math>)</p>	<p><b>Adversary <math>\mathcal{F}</math></b>  <math>b \leftarrow_s \{0, 1\}</math>  <math>b' \leftarrow_s \mathcal{D}^{\text{ENC}}</math>  If (b = b') then return 1  Else return 0  <u>ENC(K, m)</u>  If (b = 1) then  <math>c \leftarrow_s \text{SE.Enc}(K, m)</math>  Else  <math>j \leftarrow 0; \text{out-of-time} \leftarrow \text{false}</math>  Repeat  <math>j \leftarrow j + 1; r \leftarrow_s \{0, 1\}^{\text{SE.rl}}</math>  <math>c \leftarrow \text{SE.Enc}(K, m; r)</math>  <math>(v_c, t_c) \leftarrow_s \text{FN}(c)</math>  <math>\text{success} \leftarrow (K[t_c] = v_c)</math>  <math>\text{out-of-time} \leftarrow (j = s)</math>  Until (success OR out-of-time)  Return (c, <math>\epsilon</math>)</p>	<p><b>Game I<sub>1</sub>, I<sub>2</sub></b>  <math>C \leftarrow \emptyset</math>  <math>b' \leftarrow_s \mathcal{D}^{\text{ENC}}; \text{Return } (b' = 1)</math>  <u>ENC(K, m)</u>  <math>j \leftarrow 0;</math>  <math>\text{out-of-time} \leftarrow \text{false}</math>  Repeat  <math>j \leftarrow j + 1; r \leftarrow_s \{0, 1\}^{\text{SE.rl}}</math>  <math>c \leftarrow \text{SE.Enc}(K, m; r)</math>  <math>(v, t) \leftarrow_s \{0, 1\} \times [k]</math>  If (<math>c \in C</math>) then  <math>\text{bad} \leftarrow_s \text{true}</math>  <math>(v, t) \leftarrow (v_{c'}, t_{c'})</math>  <math>(v_c, t_c) \leftarrow (v, t)</math>  <math>C \leftarrow C \cup \{c\}</math>  <math>\text{success} \leftarrow (K[t_c] = v_c)</math>  <math>\text{out-of-time} \leftarrow (j = s)</math>  Until (success OR out-of-time)  Return (c, <math>\epsilon</math>)</p>
--	--	--

Figure 5: Games an adversary for proof of Theorem 4.2.

We proceed to upper bound  $\Pr[\text{H}_3]$ . Since the sampling in  $\text{H}_3$  continues until success we have

$$\Pr[\text{H}_3] = \Pr[\text{H}_4] = \Pr[\text{H}_5].$$

We proceed to upper bound  $\Pr[\text{H}_5]$ . Let  $p_\ell = \Pr[\text{H}_5 \text{ sets } \text{bad}_\ell]$  for  $\ell \in [k]$ . Now we can use a standard coupon collector problem analysis. For any particular  $\ell \in [k]$  we have

$$p_\ell = \left(1 - \frac{1}{k}\right)^q \leq e^{-q/k}.$$

Thus

$$\Pr[\text{H}_5] \leq p_1 + \dots + p_k \leq ke^{-q/k}.$$

Putting the inequalities together yields Eq. (3). ■

**STRONG UNDETECTABILITY.** The following theorem says that the ASA of our adversary  $\mathbf{A}$  of Fig. 3 is strongly undetectable.

**Theorem 4.2** *Let SE be a symmetric encryption scheme and let  $k = \text{SE.kl}$ . Let  $\text{F} : \{0, 1\}^{\text{F.kl}} \times \{0, 1\}^* \rightarrow \{0, 1\} \times [k]$  be a PRF. Let  $q, s \geq 1$  and let  $\mathbf{A}$  be defined as in Fig. 3. Let  $\mathcal{D}$  be an adversary against the strong undetectability of  $\mathbf{A}$  that makes at most  $n$  queries to its ENC oracle. Then we can build PRF adversary  $\mathcal{F}$  such that*

$$\text{Adv}_{\text{SE}, \mathbf{A}}^{\text{sdet}}(\mathcal{D}) \leq 2\text{Adv}_{\text{F}}^{\text{prf}}(\mathcal{F}) + n^2 s^2 \cdot 2^{-\mathbf{H}_\infty(\text{SE})}. \quad (4)$$

The running time of  $\mathcal{F}$  is about that of  $\mathcal{D}$  and it makes at most  $ns$  oracle queries.

**Proof of of Theorem 4.2:** Game I<sub>0</sub> of Fig. 5 implements SDET with  $\text{F}(K_{\mathbf{A}}, \cdot)$  replaced by a

lazily sampled random function. Let prf adversary  $\mathcal{F}$  be as defined in the second panel of Fig. 5. It runs  $\mathcal{D}$ , itself simulating the latter's ENC oracle, in the process appealing to its own FN oracle. Letting  $b_{\text{prf}}$  denote the challenge bit in game PRF, we have

$$\begin{aligned} \Pr[\text{PRF}_{\mathbb{F}}^{\mathcal{F}} \mid b_{\text{prf}} = 1] &= \frac{1}{2} + \frac{1}{2} \mathbf{Adv}_{\text{SE}, \mathbb{A}}^{\text{sdet}}(\mathcal{D}) \\ 1 - \Pr[\text{PRF}_{\mathbb{F}}^{\mathcal{F}} \mid b_{\text{prf}} = 0] &= \Pr[\text{I}_0]. \end{aligned}$$

Thus, using Eq. (1), we have

$$\begin{aligned} \mathbf{Adv}_{\mathbb{F}}^{\text{prf}}(\mathcal{F}) &= \Pr[\text{PRF}_{\mathbb{F}}^{\mathcal{F}} \mid b_{\text{prf}} = 1] - (1 - \Pr[\text{PRF}_{\mathbb{F}}^{\mathcal{F}} \mid b_{\text{prf}} = 0]) \\ &= \frac{1}{2} + \frac{1}{2} \mathbf{Adv}_{\text{SE}, \mathbb{A}}^{\text{sdet}}(\mathcal{D}) - \Pr[\text{I}_0]. \end{aligned}$$

Re-arranging terms, we have

$$\mathbf{Adv}_{\text{SE}, \mathbb{A}}^{\text{sdet}}(\mathcal{D}) \leq 2 \mathbf{Adv}_{\mathbb{F}}^{\text{prf}}(\mathcal{F}) + (2 \Pr[\text{I}_0] - 1).$$

We proceed to upper bound  $2 \Pr[\text{I}_0] - 1$ . Consider games  $\text{I}_1, \text{I}_2$  of Fig. 5. Game  $\text{I}_2$  includes the boxed code while  $\text{I}_1$  does not. In both,  $(v_c, t_c)$  is picked at random. However, in  $\text{I}_2$ , if it turns out that  $c$  was already seen, then the game corrects, resetting  $(v_c, t_c)$  to its prior value. As a result,  $\text{I}_2$  is equivalent to  $\text{I}_0$  with the challenge bit  $b$  in the latter set to  $b = 1$ . On the other hand, in  $\text{I}_1$ , the choice of  $(v_c, t_c)$  is independent of the choice of  $c$  and does not impact the distribution of the latter, making it equivalent to  $\text{I}_0$  with  $b = 0$ . Thus we have

$$\begin{aligned} \Pr[\text{I}_2] &= \Pr[\text{I}_0 \mid b = 1] \\ \Pr[\text{I}_1] &= \Pr[\text{I}_0 \mid b = 0]. \end{aligned}$$

Additionally, games  $\text{I}_0, \text{I}_1$  are identical until **bad**. Via the fundamental lemma of game playing [4] we have

$$\begin{aligned} 2 \Pr[\text{I}_0] - 1 &= \Pr[\text{I}_0 \mid b = 1] - \Pr[\text{I}_0 \mid b = 0] \\ &= \Pr[\text{I}_2] - \Pr[\text{I}_1] \\ &\leq \Pr[\text{I}_1 \text{ sets bad}]. \end{aligned}$$

Game  $\text{I}_1$  sets **bad** whenever the same value of  $c$  is returned twice by SE.Enc. There are  $ns$  queries to ENC. Thus

$$\begin{aligned} \Pr[\text{I}_1 \text{ sets bad}] &\leq \binom{ns}{2} \cdot 2^{-\mathbf{H}_{\infty}(\text{SE})} \\ &\leq n^2 s^2 \cdot 2^{-\mathbf{H}_{\infty}(\text{SE})-1}. \end{aligned}$$

Putting everything together gives the desired Eq. (4).  $\blacksquare$

DISCUSSION. We do not expect that most users will mount involved detection efforts. But the weakness of the ASA of BPR [3] is that it is relatively easily detectable, even without much effort, due to its use of state. The most obvious way for this to happen is that the state is reset, for example due to a system reboot. When this happens, the decrypter will be able to detect the subversion. This, however, is not captured by their notion of detectability. Strong undetectability

fills the gap. Our new ASA, achieving this and in particular being stateless, is harder to detect and thus more dangerous.

In a world of subversion, there are no panaceas. The extent to which big brother will risk detection in order to successfully recover a key is not known, but it would be natural that, all else being equal, big brother will pick the ASA that minimizes the chance of detection. This means an ASA like ours. Our work is intended to increase awareness and spur the search for deterrents.

## 5 Defenses

BPR [3] present schemes that are subversion resilient. They define unique ciphertext schemes, show they are subversion resilient, and then provide several ways to build them. These schemes are deterministic and stateful, meaning the encryptor and decryptor have to maintain a synchronized state.

Having strengthened the BPR undetectability condition to strong undetectability, a natural question is whether subversion resilience can still be achieved. In fact this strengthening makes no difference, and the schemes of BPR continue to achieve subversion resilience. This is because a subversion resilient scheme according to BPR [3] is one where any subversion satisfying the decryptability condition cannot succeed in their subversion game. Decryptability means that ciphertexts created by the subverted encryption algorithm decrypt properly under the real decryption algorithm, and it can be seen as represented a particular and fixed form of undetectability.

DFP [5] critique the BPR decryptability condition and instead suggest an alternative formulation of detectability represented by their  $\overline{\text{DETECT}}$  game in which the subversion and detection adversaries run together, the latter getting a transcript of the interaction of the former with its oracles. They show that unique ciphertext schemes continue to achieve their notion. In this context, consideration of strong undetectability would involve modifying their game to a new one,  $\overline{\text{SDETECT}}$ . For example one might return the state to the detection adversary as part of the transcript, and also add a reset oracle  $\text{RESET}$  to allow state reset. Then the question is whether unique ciphertext schemes continue to meet this new definition. They do.

## References

- [1] G. Ateniese, B. Magri, and D. Venturi. Subversion-resilient signature schemes. Cryptology ePrint Archive, Report 2015/517, 2015. <http://eprint.iacr.org/2015/517>.
- [2] M. Bellare, J. Kilian, and P. Rogaway. The security of cipher block chaining. In Y. Desmedt, editor, *CRYPTO'94*, volume 839 of *LNCS*, pages 341–358. Springer, Aug. 1994.
- [3] 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, Aug. 2014.
- [4] 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, May / June 2006.
- [5] 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*. Springer, Mar. 2015.



- [6] 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, Apr. 2015.
- [7] 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, Oct. 2003.
- [8] O. Goldreich, S. Goldwasser, and S. Micali. How to construct random functions. *Journal of the ACM*, 33(4):792–807, Oct. 1986.
- [9] 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, Apr. 2015.
- [10] J. Patarin and L. Goubin. Asymmetric cryptography with S-boxes. In Y. Han, T. Okamoto, and S. Qing, editors, *ICICS 97*, volume 1334 of *LNCS*, pages 369–380. Springer, Nov. 1997.
- [11] K. G. Paterson. Imprimitive permutation groups and trapdoors in iterated block ciphers. In L. R. Knudsen, editor, *FSE'99*, volume 1636 of *LNCS*, pages 201–214. Springer, Mar. 1999.
- [12] V. Rijmen and B. Preneel. A family of trapdoor ciphers. In E. Biham, editor, *FSE'97*, volume 1267 of *LNCS*, pages 139–148. Springer, Jan. 1997.
- [13] A. Russell, Q. Tang, M. Yung, and H.-S. Zhou. Cliptography: Clipping the power of kleptographic attacks. Cryptology ePrint Archive, Report 2015/695, 2015. <http://eprint.iacr.org/>.
- [14] B. Schneier, M. Fredrikson, T. Kohno, and T. Ristenpart. Surreptitiously weakening cryptographic systems. Cryptology ePrint Archive, Report 2015/097, 2015. <http://eprint.iacr.org/2015/097>.
- [15] 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.
- [16] 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, Apr. 1985.
- [17] G. J. Simmons. A secure subliminal channel (?). In H. C. Williams, editor, *CRYPTO'85*, volume 218 of *LNCS*, pages 33–41. Springer, Aug. 1986.
- [18] G. J. Simmons. Subliminal communication is easy using the DSA. In T. Helleseth, editor, *EUROCRYPT'93*, volume 765 of *LNCS*, pages 218–232. Springer, May 1994.
- [19] 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, Aug. 1996.
- [20] A. Young and M. Yung. Kleptography: Using cryptography against cryptography. In W. Fumy, editor, *EUROCRYPT'97*, volume 1233 of *LNCS*, pages 62–74. Springer, May 1997.
- [21] A. Young and M. Yung. Monkey: Black-Box symmetric ciphers designed for MONopolizing KEYS. In S. Vaudenay, editor, *FSE'98*, volume 1372 of *LNCS*, pages 122–133. Springer, Mar. 1998.
- [22] A. Young and M. Yung. A subliminal channel in secret block ciphers. In H. Handschuh and A. Hasan, editors, *SAC 2004*, volume 3357 of *LNCS*, pages 198–211. Springer, Aug. 2004.
- [23] A. L. Young and M. Yung. Backdoor attacks on black-box ciphers exploiting low-entropy plaintexts. In R. Safavi-Naini and J. Seberry, editors, *ACISP 03*, volume 2727 of *LNCS*, pages 297–311. Springer, July 2003.