

# Cliptography: Clipping the Power of Kleptographic Attacks

Alexander Russell\*    Qiang Tang<sup>†</sup>    Moti Yung<sup>‡</sup>    Hong-Sheng Zhou<sup>§</sup>

May 21, 2016

## Abstract

Kleptography, introduced 20 years ago by Young and Yung [Crypto '96], considers the (in)security of malicious implementations (or instantiations) of standard cryptographic primitives that embed a “backdoor” into the system. Remarkably, crippling subliminal attacks are possible even if the subverted cryptosystem produces output indistinguishable from a truly secure “reference implementation.” Bellare, Paterson, and Rogaway [Crypto '14] recently initiated a formal study of such attacks on symmetric key encryption algorithms, demonstrating a kleptographic attack can be mounted in broad generality against randomized components of cryptographic systems.

We enlarge the scope of current work on the problem by permitting adversarial subversion of (randomized) key generation; in particular, we initiate the study of cryptography in the *complete subversion model*, where *all* relevant cryptographic primitives are subject to kleptographic attacks. We construct secure one-way permutations and trapdoor one-way permutations in this “complete subversion” model, describing a general, rigorous immunization strategy to clip the power of kleptographic subversions. Our strategy can be viewed as a formal treatment of the folklore “nothing up my sleeve” wisdom in cryptographic practice. We also describe a related “split program” model that can directly inform practical deployment. We additionally apply our general immunization strategy to directly yield a backdoor-free PRG. This notably amplifies previous results of Dodis, Ganesh, Golovnev, Juels, and Ristenpart [Eurocrypt '15], which require an honestly generated random key.

We then examine two standard applications of (trapdoor) one-way permutations in this complete subversion model and construct “higher level” primitives via black-box reductions. We showcase a digital signature scheme that preserves existential unforgeability when *all* algorithms (including key generation, which was not considered to be under attack before) are subject to kleptographic attacks. Additionally, we demonstrate that the classic Blum–Micali pseudorandom generator (PRG), using an “immunized” one-way permutation, yields a backdoor-free PRG.

Alongside development of these secure primitives, we set down a hierarchy of kleptographic attack models which we use to organize past results and our new contributions; this taxonomy may be valuable for future work.

---

\*University of Connecticut, acr@cse.uconn.edu

<sup>†</sup>Cornell University, qt44@cornell.edu

<sup>‡</sup>Snapchat Inc. & Columbia University, moti@cs.columbia.edu

<sup>§</sup>Virginia Commonwealth University, hszhou@vcu.edu

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>A Definitional Framework for Cliptography</b>	<b>5</b>
2.1	From Cryptography to Cliptography . . . . .	5
2.2	A Formal Definition . . . . .	7
<b>3</b>	<b>Subversion-Resistant One-Way Permutations</b>	<b>11</b>
3.1	Defining subversion resistant OWP/TDOWP . . . . .	12
3.2	Constructing subversion-resistant <sup>A</sup> OWP . . . . .	15
3.3	Constructing subversion-resistant <sup>SP</sup> OWP/TDOWP . . . . .	18
<b>4</b>	<b>Subversion-Resistant Signatures</b>	<b>19</b>
<b>5</b>	<b>Subversion-Resistant Pseudorandom Generators</b>	<b>21</b>
5.1	Preliminaries: The definition of a subversion-resistant <sup>A</sup> PRG . . . . .	22
5.2	Constructing $q$ -PRG <sup>A</sup> from a OWP <sup>A</sup> . . . . .	23
5.3	A general public immunization strategy for PRG <sup>A</sup> . . . . .	24
<b>A</b>	<b>Supporting Material: Omitted Definitions</b>	<b>27</b>
<b>B</b>	<b>Supporting Material: Omitted Constructions</b>	<b>28</b>
B.1	Full Fledged Subversion Resistant PRG <sup>A</sup> . . . . .	28
<b>C</b>	<b>Supporting Material: Omitted Proofs</b>	<b>29</b>
C.1	Proof of Lemma 3.2 . . . . .	29
C.2	Proof of Theorem 3.6 . . . . .	30
C.3	Proof of Theorem 4.3 . . . . .	31
C.4	Proof of Lemma 5.2 . . . . .	32

# 1 Introduction

Consider the conventional use of a cryptographic primitive, such as an encryption scheme: To encrypt a desired plaintext, one simply runs an implementation (or an instantiation with particular parameters) of the encryption algorithm obtained from a hardware or software provider. Although the underlying algorithms may be well-studied and proven secure, malicious implementations and instantiations may cleverly “backdoor” the system or directly embed sensitive information—such as the secret key—into the ciphertext in a fashion that permits recovery by the provider/manufacturer but is undetectable to other parties. It is notable that *such leakage is possible even if the implementation produces “functionally and statistically clean” output that is indistinguishable from that of a faithful implementation*. While the underlying concept of kleptography was proposed by Young and Yung two decades ago [YY96, YY97], striking recent examples—including those of the Snowden revelations [PLS13]—have reawakened the security community to the seriousness of these issues [Rog15]. As a result, the topic has received renewed attention; see, e.g., [BPR14, BH15, DGG<sup>+</sup>15, MS15, AMV15]. In particular, Bellare, Paterson, and Rogaway [BPR14] studied algorithm substitution attacks—with a focus on symmetric key encryption—and demonstrated a devastating framework for such attacks that apply in broad generality to randomized algorithms. These results were later amplified [BJK15] to show that such attacks can be carried out even if the adversarial implementation is *stateless*. Soon after, Dodis, Ganesh, Golovnev, Juels, and Ristenpart [DGG<sup>+</sup>15] formalized the subversion of Dual\_EC pseudorandom generators (PRG) and studied backdoored PRGs in generality; they additionally studied methods for “immunizing” PRG in such hostile settings.

**Our contributions.** We continue this line of pursuit. Specifically, we are motivated to develop cryptographic schemes in a *complete subversion model*, in which *all* algorithms of a scheme are potentially subverted by the adversary. This model provides a conceptually simple abstraction of the adversary’s power, and significantly amplifies previously studied settings, which rely on trusted key generation or clean randomness that is assumed private from the adversary.

In particular, motivated by the question of defending against the kleptographic attacks on **key generation** as demonstrated in the original paper of [YY96, YY97], we study two fundamental cryptographic primitives in the complete subversion model—(trapdoor) one-way permutations (OWP)—and apply these primitives to construct other cryptographic tools such as digital signatures and PRGs. Along the way, we identify novel generic defending strategies and a hierarchy of attack models. We hope to stimulate a systematic study of “**cliptography**,” the challenge of developing a broad class of familiar cryptographic tools that remain secure in such kleptographic settings. As mentioned above, prior to our work kleptographic attacks on various primitives have been addressed in weaker models; see *Related work* in Section 1. In detail, we show the following:

- Inspired by the model of Bellare, Paterson and Rogaway [BPR14], we set down a hierarchy of three security models that capture practical kleptographic settings. The models are characterized by three parties: an adversary, who may provide potentially subverted implementations of *all* cryptographic algorithms; a “watchdog,” who either certifies or rejects the implementations by subjecting them to (black-box) interrogation;<sup>1</sup> and a challenger, who plays a conventional security game (but now using the potentially subverted algorithms) with the adversary. Armed with the “specification” of the cryptographic algorithms and oracle

---

<sup>1</sup>Without the watchdog, it is elusive to achieve interesting cryptographic functionalities in those stringent settings.

access to the implementations provided by the adversary, the watchdog attempts to detect any subversion in the implementations. Various models arise by adjusting the supervisory power of the watchdog; see Section 2.

- We study (trapdoor) one-way permutations in the presence of kleptographic attacks, introducing notions of subversion-resistance that can survive various natural kleptographic attacks. We first give a simple example of a OWP that can be proven secure in the conventional sense, but can be completely broken under the kleptographic attack. This demonstrates the need for judicious design of cryptographic primitives to defend against kleptographic attacks.

We then construct subversion-resistant (trapdoor) one way permutations via a general transformation that “sanitizes” arbitrary OWPs by *randomizing the function index*. This transformation clips potential correlation between the function and the possible backdoor that the adversary may possess. Additionally, we introduce a *split-program* model to make the general method above applicable using only standard hash functions (see Section 3.3).

- In Section 4, we then observe that subversion-resistant trapdoor OWPs give us a way to construct key generation algorithms (for digital signature schemes) against kleptographic attacks. We then showcase a concrete example of a digital signature scheme in the complete subversion model. More concretely, we achieve this result by (1) using the subversion-resistant trapdoor one way permutation directly as a key generation algorithm, and then (2) instantiating the unique signature generation mechanism via full domain hash (FDH). We stress that the reduction of the standard FDH signature scheme does not go through in the kleptographic setting. To resolve this issue, we slightly modify the FDH approach by hashing the message *together with the public key*. We remark that the original kleptographic attacks [YY96, YY97] were indeed applied to the key generation algorithm, while recent efforts [BPR14, AMV15] shift focus to other algorithmic aspects of encryption or digital signature schemes, assuming that key generation is honest. Our result is the first digital signature scheme allowing the adversary to sabotage all algorithms, including key generation.
- We then turn our attention to PRGs. Previous work of Dodis et al. [DGG<sup>+</sup>15] investigated a notion of “backdoored PRG” in which the adversary sets up a PRG instance (i.e., the public parameter), and is able to distinguish the output from uniform with a backdoor. They then proposed powerful immunizing strategies which apply a keyed hash function to the output—*assuming the key is unknown to the adversary*—in the public parameter generation phase. Motivated by their success, we focus on constructing backdoor-free PRGs in the complete subversion model (where such clean randomness is not permitted). Our first construction is based on the classic Blum-Micali construction, using our subversion-resistant OWP and the Goldreich-Levin hardcore predicate. Dodis et al. [DGG<sup>+</sup>15] additionally show that it is impossible to achieve a public immunizing strategy for all PRGs by applying a public function to the PRG output. We sidestep this impossibility result via an alternative public immunizing strategy: Rather than randomizing the output of the PRG, we randomize the public parameter of PRG, which yields a general construction for PRG in the complete subversion model. See Section 5.

Finally, we remark that black-box constructions and reductions do not, in general, survive in the kleptographic model. However, two of the results above—the Blum-Micali construction and the signature scheme—give explicit examples of reductions that can be salvaged.

**Remarks: Our techniques and the “nothing up my sleeve” principle; single use of randomized algorithms and subliminal channels.** We remark that our general defending technique significantly differs from known methods: We use a—potentially subverted—hash function to “randomize” the index and public parameter of a (perhaps randomized) algorithm so that any correlation with some potential backdoor can be eliminated. This can be seen as an instance of the folklore wisdom of a “nothing up my sleeve number” [Wik] which has been widely used in practical cryptographic designs. The basic principle calls for constants appearing in the development of cryptographic algorithms to be drawn from a “rigid” source, like the digits of  $\pi$ ; the idea is that this prevents them from possessing hidden properties that might give advantage to an attacker (or the designer). In our setting, the fact that a given value  $v$  is *supplied along with a preimage  $x$  so that  $h(x) = v$*  (for a hash function  $h$ ) is a evidence that  $v$  has “nothing up its sleeve.” In fact, the situation is complicated: While this does effectively mean that  $v$  is generated by selecting  $x$  and computing  $h(x)$  and, thus, severely restricts the possibility for tampering with  $v$ , it does not eliminate subliminal channels introduced by rejection sampling or entirely “clean”  $v$ . In particular, detailed analysis is still required to control the behavior of  $v$ .

Previous results either use a trusted random source to re-randomize the output of a randomized algorithm, or consider only deterministic algorithms. Permitting randomized algorithms in a kleptographic framework immediately invites the (devastating) general “steganochannel” attack of Bellare et al. [BPR14, BJK15]. While the prospect of full “immunization” for general randomized algorithms (in particular, generic destruction of a steganochannel) is a—presumably challenging—direction of future work, our primitives do permit randomized algorithms, although the security games we analyze invoke them only once (to, e.g., derive a key).

For simplicity, we focus on (potentially subverted) algorithms that do not maintain “internal state” between invocations. We remark that typical steganographic attacks can indeed be carried out in a stateless model [BJK15]. Moreover, this restriction can be lifted for the constructions in the paper; see Remark 2.5.

**Related work.** The concept of *kleptography*—subverting cryptographic algorithms by modifying their implementations to leak secrets covertly—was proposed by Young and Yung [YY96, YY97] in 1996. They gave concrete examples showing that backdoors can be embedded into the public keys of commonly used cryptographic schemes; while the resulting public keys appear normal to every user, the adversary is nevertheless capable of learning the secret keys. It may not be surprising that defending against such deliberate attacks is challenging and only limited feasibility results exist. We next briefly describe these existing results.

In [JG02], Juels and Guajardo suggested the following idea: the user and a trusted certificate authority (CA) jointly generate the public key; as a part of this process, the user proves to the CA that the public key is generated honestly. This contrasts markedly with our setting, where the user does not have any secret, and every component is provided by the adversary.

Bellare et al. considered a powerful family of kleptographic attacks that they call *algorithm substitution attacks*, and explore these in both symmetric-key [BPR14, BJK15] and public-key [BH15] settings. They first proposed a generic attack, highlighting the relevance of steganographic techniques in this framework: specifically, a sabotaged randomized algorithm can leak a secret bit-by-bit by invoking steganographic rejection-sampling; then an adversary possessing the backdoor can identify the leaked bits from the biased output, which appears unmolested to other observers. The attack and analysis relies on the effectiveness of subliminal channels [Sim83, Sim86, HLv02], and is particularly striking because it can be applied in such generality. They then introduced

a framework for defending against such attacks by focusing on algorithms that having a unique output for each input: relevant examples of such algorithms include unique ciphertext encryption algorithms. These results were later refined by [DFP15]. Their defending mechanism does not, however, address the (necessarily randomized) process of key generation—it implicitly assumes key generation to be honest. This state of affairs is the direct motivation of the current article: we adopt a significantly amplified *complete subversion model* where *all* cryptographic algorithms—including key generation—are subject to kleptographic (i.e., substitution) attacks. This forces us to manage certain randomized algorithms (such as key generation) in a kleptographic setting. The details of the model, with associated commentary about its relevance to practice, appear below.

Dodis et al. [DGG<sup>+</sup>15] pioneered the rigorous study of pseudorandom generators in such settings, developing an alternative family of kleptographic attacks on pseudorandom generators in order to formalize the notorious Dual.EC PRG subversion [NIS12, CNE<sup>+</sup>14]. In their model, the adversary subverts the security of the PRG by opportunistically setting the public parameter while privately keeping some backdoor information (instead of providing an implementation). They then demonstrate an impossibility result: backdoored PRGs cannot be immunized by applying a public function—even a trusted random oracle—to the output. They also proposed and analyzed immunizing strategies obtained by applying a keyed hash function to the output (of the PRG). Note that the (hash) key plays a special role in their model: it is selected uniformly and is unknown to the adversary during the public parameter generation phase. These results likewise inspire our adoption of the amplified *complete subversion model*, which excludes such reliance on public randomness beyond the reach of the adversary. In particular, our general immunizing strategy (randomizing the public parameter of a backdoored PRG instead of randomizing the PRG output) permits us to bypass the impossibility result. Additionally, our results on subversion-resistant OWFs can be applied to construct a specific “backdoor-free” PRG following the classic Blum-Micali framework.

Other works suggest different angles of defense against mass surveillance. For example, in [MS15, DMSD15] the authors proposed a general framework of safeguarding protocols by randomizing the incoming/outgoing messages via a trusted (reverse) firewall. Their results demonstrate that with a trusted random source, many tasks become achievable. As they rely on a “subversion-free” firewall, these results require a more generous setting than provided by our *complete subversion model*.

Ateniese et al. [AMV15] continued the study of algorithm substitution attacks on signatures and propose two defending mechanisms, one utilizes a unique signature scheme assuming the key generation and verify algorithms to be honest; the other adopts the reverse firewall model that assumes trusted randomness. We construct a signature scheme that can be proven secure in the complete subversion model which does not make assumptions on honesty or require trusted randomness. We remark that the strength of the “watchdog” that is required for the signature scheme is, however, stronger than that required for the other primitives; it must be permitted a transcript of the security game. See Section 4.

## 2 A Definitional Framework for Cliptography

### 2.1 From Cryptography to Cliptography

In this section, we lay down a definitional framework for cliptography. The adversary in this new setting is “proud-but-malicious”: as reflected in the models of [BPR14], the adversary aims at supplying subverted implementation to break security while keeping the subversion “under the radar” of any detection. Thus the basic framework reflects the ability of the adversary to provide (potentially subverted) implementations of the cryptographic algorithms of interest, the ability of an efficient “watchdog” to interrogate such implementations in order to check their veracity, and a classical “challenger-adversary” security game. Specifically, the model considers an adversary that commences activities by supplying a (potentially subverted) implementation of the cryptographic primitive; one then considers two parallel procedures: a classical challenger-adversary security game in which the challenger must use only (oracle access to) the adversary’s implementations, and a process in which the “watchdog” compares—also via oracle access—the adversary’s implementations against a specification of the primitives. (For entertainment, we occasionally refer to the adversary as “big brother.”)

**Cryptographic games.** We express the security of (standard) cryptographic schemes via *cryptographic games* between a challenger  $\mathcal{C}$  and an adversary  $\mathcal{A}$ .

**Definition 2.1 (Cryptographic Game [HH09])** A cryptographic game  $G = (\mathcal{C}, \delta)$  is defined by a random system  $\mathcal{C}$ , called the challenger, and a constant  $\delta \in [0, 1)$ . On security parameter  $\lambda$ , the challenger  $\mathcal{C}(1^\lambda)$  interacts with some adversary  $\mathcal{A}(1^\lambda)$  and outputs a bit  $b$ . We denote this interaction by  $b = (\mathcal{A}(1^\lambda) \Leftrightarrow \mathcal{C}(1^\lambda))$ . The advantage of an attacker  $\mathcal{A}$  in the game  $G$  is defined as

$$\mathbf{Adv}_{\mathcal{A}, G}(1^\lambda) = \Pr[(\mathcal{A}(1^\lambda) \Leftrightarrow \mathcal{C}(1^\lambda)) = 1] - \delta.$$

We say a cryptographic game  $G$  is secure if for all PPT attackers  $\mathcal{A}$ , the advantage  $\mathbf{Adv}_{\mathcal{A}, G}(1^\lambda)$  is negligible.

The above conventional security notions are formulated under the assumption that the relevant algorithms of the cryptographic scheme are faithfully implemented and, moreover, that participants of the task have access to truly private randomness (thus have, e.g., truly random keys).

**The complete subversion model.** A basic question that must be addressed by a kleptographic model concerns the selection of algorithms the adversary is permitted to subvert. We work exclusively in a setting where the adversary is permitted to provide implementations of *all* the relevant cryptographic elements of a scheme, a setting we refer as the *complete subversion* model. Thus, all guarantees about the quality of the algorithms are delivered by the watchdog’s testing activities. This contrasts with previous work, which explicitly protected some of the algorithms from subversion, or assumed clean randomness. Such a setting we refer to as *partial subversion* model.

**Choosing the right watchdog.** By varying the information provided to the watchdog, one obtains different models that reflect various settings of practical interest. The weakest (and perhaps most attractive) model is the *offline* watchdog, which simply interrogates the supplied implementations, comparing them with the specification of the primitives, and declares them to be “fit” or “unfit.” Of course, we must insist that such watchdogs find the actual specification “fit”: formally,

the definition is formulated in terms of distinguishing an adversarial implementation from the specification.

One can strengthen the watchdog by permitting it access to the full transcript of the challenger-adversary security game, resulting in the *online* watchdog. Finally, we describe an even more powerful *omniscient* watchdog, which is even privy to private state of the challenger. (While we do not use such a powerful watchdog in our results, it is convenient for discussing previous work.)

We remark these various watchdogs reflect various levels of “checking” that a society might entertain for cryptographic algorithms (and conversely, various levels of tolerance that an adversary may have to exposure): the offline watchdog reflects a “one-time” laboratory that attempts to check the implementations; an online watchdog actually crawls public transcripts of cryptographic protocols to detect errors; the omniscient watchdog requires even more, involving (at least) individuals effectively checking their results against the specification.

## 2.2 A Formal Definition

Having specified the power of the big brother (the adversary) and that of the watchdog, we are ready to introduce *cliptographic games* to formulate security. To simplify the presentation, we here initially consider *complete* subversion with an *offline* watchdog. In the next section, we will discuss the other variants.

A cryptographic scheme  $\Pi$  consists of a set of (possibly randomized) algorithms  $(F^1, \dots, F^k)$  where each  $F^i$  is with input space  $\mathcal{X}^i$ , output space  $\mathcal{Y}^i$ , and randomness space  $\mathcal{R}^i$ . For example, a digital signature scheme consists of three algorithms, a (randomized) key generation algorithm, a signing algorithm, and deterministic verification algorithm. The definition of a scheme  $\Pi = (F^1, \dots, F^k)$  results in a specification of the associated algorithms; for concreteness, we label these as  $\Pi_{\text{SPEC}} = (F_{\text{SPEC}}^1, \dots, F_{\text{SPEC}}^k)$ ; when a scheme is (perhaps adversarially) implemented, we denote the implementation as  $\Pi_{\text{IMPL}} = (F_{\text{IMPL}}^1, \dots, F_{\text{IMPL}}^k)$ . If the implementation honestly follows the specification of the scheme, we overload the notation and represent them interchangeably with the specification as  $\Pi_{\text{SPEC}}$ .

In our definition, the adversary  $\mathcal{A}$  will interact with both the challenger  $\mathcal{C}$  and the watchdog  $\mathcal{W}$ . (In the offline case, these interactions are independent; in the online case,  $\mathcal{W}$  is provided a transcript of the interaction with  $\mathcal{C}$ .) Following the definition of cryptographic game, we use  $b_{\mathcal{C}} = (\mathcal{A}(1^\lambda) \leftrightarrow \mathcal{C}^{F_{\text{IMPL}}^1, \dots, F_{\text{IMPL}}^k}(1^\lambda))$  to denote the interaction between  $\mathcal{A}$  and  $\mathcal{C}$ ;  $b_{\mathcal{C}}$  denotes the bit returned by the challenger  $\mathcal{C}$ . (Note that the challenger must use the implementation of  $\Pi$  provided by the adversary, thus the interaction between  $\mathcal{A}, \mathcal{C}$  now is denoted by  $\leftrightarrow$  instead of  $\Leftrightarrow$  in the classical cryptographic game in Definition 2.1.)

As for the watchdog  $\mathcal{W}$ , the adversary provides  $\mathcal{W}$  his potentially subverted implementations  $\Pi_{\text{IMPL}}$  of the primitive (as oracles);  $\mathcal{W}$  may then interrogate them in an attempt to detect divergence from the specification, which he possesses. On the basis of these tests, the watchdog produces a bit (Intuitively, the bit indicates whether the implementations passed whatever tests the watchdog carried out to detect inconsistencies with the specification.)

**Definition 2.2 (Cliptographic Game)** A cliptographic game  $\widehat{\mathcal{G}} = (\mathcal{C}, \Pi_{\text{SPEC}}, \delta)$  is defined by a challenger  $\mathcal{C}$ , a specification  $\Pi_{\text{SPEC}}$ , and a constant  $\delta \in [0, 1)$ . Given an adversary  $\mathcal{A}$ , a watchdog  $\mathcal{W}$ , and a security parameter  $\lambda$ , we define the detection probability of the watchdog  $\mathcal{W}$  with respect to  $\mathcal{A}$  to be

$$\text{Det}_{\mathcal{W}, \mathcal{A}}(1^\lambda) = \left| \Pr[\mathcal{W}^{F_{\text{IMPL}}^1, \dots, F_{\text{IMPL}}^k}(1^\lambda) = 1] - \Pr[\mathcal{W}^{F_{\text{SPEC}}^1, \dots, F_{\text{SPEC}}^k}(1^\lambda) = 1] \right|,$$



where  $\Pi_{\text{IMPL}} = (F_{\text{IMPL}}^1, \dots, F_{\text{IMPL}}^k)$  denotes the implementation produced by  $\mathcal{A}$ . The advantage of the adversary is defined to be

$$\text{Adv}_{\mathcal{A}}(1^\lambda) = \left| \Pr \left[ (\mathcal{A}(1^\lambda) \rightsquigarrow \mathcal{C}^{F_{\text{IMPL}}^1, \dots, F_{\text{IMPL}}^k}(1^\lambda)) = 1 \right] - \delta \right|.$$

We say that a game is **subversion-resistant** if there exists a PPT watchdog  $\mathcal{W}$  such that for all PPT adversaries  $\mathcal{A}$ , either  $\text{Det}_{\mathcal{W}, \mathcal{A}}$  is non-negligible, or  $\text{Adv}_{\mathcal{A}}$  is negligible, in the security parameter  $\lambda$ .

**Discussion: random oracles.** We remark that the definition of a cryptographic game requires setting down both the challenger *and* the specification of the primitive  $\Pi$ , as this latter data determines the notion of “detection probability.” In many settings, we establish results in the conventional *random oracle* model which requires some special treatment in the model above. In general, we consider a random oracle to be an (extremely powerful) heuristic substitute for a deterministic function with strong cryptographic properties. In a kleptographic setting with complete subversion, we must explicitly permit the adversary to tamper with the “implementation” of the random oracle supplied to the challenger. In such settings, we provide the watchdog—as usual—oracle access to both the “specification” of the random oracle (simply a random function) and the adversary’s “implementation” of the random oracle, which may deviate from the random oracle itself. For concreteness, we permit the adversary to “tamper” with a random oracle  $h$  by providing an efficient algorithm  $T^h(x)$  (with oracle access to the random oracle  $h$ ) which computes the “implementation”  $\tilde{h}$ —thus the implementation  $\tilde{h}(x)$  is given by  $T^h(x)$  for all  $x$ . Likewise, during the security game, the challenger is provided oracle access only to the potentially subverted implementation  $\tilde{h}$  of the random oracle. As usual, the probabilities defining the security (and detection) games are taken over the choice of the random oracle. Fortunately, when the random oracle is applied to a known input distribution, an offline watchdog can ensure that the implementation is almost consistent with its specification (see Lemma 2.3).

**Other watchdog variants.** In the above definition, we chose the strongest model: the watchdog is *universal* and *offline*. In particular, primitives secure in this model are secure in any of the other models considered. The detection algorithm of the watchdog must be designed for a *given specification*, regardless of how the adversary subverts the implementation; furthermore, it may only carry out a one-time check on the implementation (and may not supervise the security game). To permit a broader class of feasibility results, it is possible to extend the basic model in both directions.

**Swapping the quantifiers.** It is also reasonable to consider a watchdog that may be tailored to the adversary, i.e., the quantifiers are changed to be  $\forall \mathcal{A}, \exists \mathcal{W}$ . Indeed, such quantification (or even more generous settings, see below) was considered implicitly in previous works, e.g., [BPR14, DGG<sup>+</sup>15, BJK15]. We remark that such a model is still highly non-trivial in that the adversary can be randomized by, e.g., selecting a random backdoor. (Thus knowing the code of the adversary does not necessarily help the watchdog to identify a faulty implementation which might be based on a random backdoor that is only known to  $\mathcal{A}$ .) Note that such a model is particularly interesting for evaluating *attacks*, where one would like to guarantee that the attack is undetectable even by a watchdog privy to the details of the adversary: specifically, when establishing security, weak watchdogs are preferable; when establishing the value of an attack, strong watchdogs are preferable.

We develop one-way permutations and pseudorandom generators in the offline model. However, it appears that richer primitives may require qualitatively stronger watchdogs regarding the inputs

they may get. Considering that an offline watchdog cannot ensure *exact* equality for deterministic algorithms, we remark that a clever adversary may be able to launch attacks by altering such deterministic functions at only a few locations. Imagine a security game where the adversary supplies a string  $m$  to which the challenger is expected to apply one of the subverted algorithms; this takes place, e.g., in the typical signature security game. The adversary may now select a random string  $w$  and implement the deterministic algorithm in such a way that it diverges from the specification at (only) this preselected point. Observe that such inconsistencies are (essentially) undetectable by an offline watchdog; however, the adversary can ensure that the subverted algorithm is indeed queried at  $w$  during the security game. Such attacks have been studied in various settings, e.g., input-triggering attacks in [BPR14, DFP15, AMV15] motivated them to consider the extra “decryptability condition” and “verifiability condition” assumptions.

An *online watchdog* can guard against this possibility; he is permitted to monitor the public interactions between users. More precisely, the online watchdog is permitted to certify both the implementations *and* the transcript between the challenger and adversary. The security game is then altered by considering  $\mathcal{W}^{\Pi_{\text{IMPL}}}(1^\lambda, \tau)$ , identical to the offline case except that the watchdog is provided the transcript  $\tau$ <sup>2</sup> of the security game ( $\mathcal{C} \rightsquigarrow \mathcal{A}$ ). (We use the shorthand notation  $\Pi_{\text{IMPL}}$  here to denote the collection of oracles  $F_{\text{IMPL}}^1, \dots, F_{\text{IMPL}}^k$ .) The detection game must then be adjusted, guaranteeing that the transcripts produced when the challenger uses  $\Pi_{\text{IMPL}}$  are indistinguishable from those produced when the challenger uses  $\Pi_{\text{SPEC}}$ . Our results on digital signature schemes will require such a watchdog. We remark that previous work on subversion-resistant digital signatures [AMV15] assumes a verifiability condition: every message-signature pair produced by the subverted sign algorithm (at least the responses to the signing queries) can pass the verification of the specification of the verify algorithm. This extra assumption can be guaranteed by an online watchdog (and, indeed, it demands either an absolute universal guarantee or an on-line guarantee for those pairs that appear in the security game).

An *omniscient watchdog* is even stronger. In addition to access to the transcript, the omniscient watchdog is aware of the entire internal state of the challenger (and can monitor the interactions between users and the subverted implementations). Similarly, by replacing  $\mathcal{W}$  in Definition 2.2 above with an omniscient watchdog, we obtain cryptographic games with omniscient watchdog. As mentioned, omniscient watchdog has been considered in literature [BPR14, DFP15]. In those works, they assume the extra decryptability condition such that ciphertext generated by the subverted encryption algorithm decrypts correctly with the honest decryption algorithm. Again, without allowing the watchdog to input the whole transcript and the decryption key, this assumption cannot be supported. While we do not require omniscient watchdog in our constructions in this paper, we discuss omniscient watchdog here as part of definitional framework.

**Discussion: the guarantees provided by an offline watchdog.** We make some general observations about the guarantees that an offline watchdog provides.

Consider a *deterministic* algorithm implemented by the adversary; an offline watchdog cannot ensure that such an algorithm is perfectly implemented. But in some cases, it can ensure that the implementation almost always agrees with the specification. One important case is that the algorithm is with public input distributions. As the adversary would keep the subverted implementation to be indistinguishable from the specification, and an offline watchdog can simply sample input from the input distribution and interrogate the implementation and compare the output with that from the specification. Formally,

---

<sup>2</sup>We remark that the transcript  $\tau$  includes the final output bit of the challenger in the security game.

**Lemma 2.3** Consider an adversarial implementation  $\Pi_{\text{IMPL}} := (F_{\text{IMPL}}^1, \dots, F_{\text{IMPL}}^k)$  of a specification  $\Pi_{\text{SPEC}} = (F_{\text{SPEC}}^1, \dots, F_{\text{SPEC}}^k)$  in a cryptographic game, where  $F^1, \dots, F^k$  are deterministic algorithms. Additionally, for each security parameter  $\lambda$ , public input distributions  $X_\lambda^1, \dots, X_\lambda^k$  are defined respectively. If  $\forall j \in [k], \Pr[F_{\text{IMPL}}^j(x) \neq F_{\text{SPEC}}^j(x) : x \leftarrow X_\lambda^j] \leq \text{negl}(\lambda)$  does not hold, Then, there is a PPT offline watchdog can detect with a non-negligible probability.

The above includes the cases that the deterministic algorithm is with a known input distribution (e.g, uniform distribution), or with an input distribution that is generated by other (adversarial) implementations. Jumping ahead, the evaluation function of a one way permutation takes a uniform input distribution; and a pseudorandom generator stretch function takes  $\mathcal{K} \times \mathcal{U}$  as (public) input distribution, where  $\mathcal{K}$  is the output distribution of a parameter generation algorithm implemented by the big brother, and  $\mathcal{U}$  is the uniform seed distribution.

In our analysis, we will use this simple observation extensively. In particular, when a hash specification is modeled as a random oracle, on most of the input points, the implementation has to be consistent with it, and thus the corresponding random oracle queries have to be asked.

Next, consider a *randomized* algorithm (with fixed inputs) that is supposed to output a high-entropy distribution. The offline watchdog can simply run the algorithm twice to see whether there is collision. If some particular value can be predicated by the adversary, it would appear with a non-negligible probability when the watchdog can reset the implementation to a same initial state.

**Lemma 2.4** Consider an adversary  $\mathcal{A}$  prepares implementation  $F_{\text{IMPL}}$  of a specification  $F_{\text{SPEC}}$  in a cryptographic game, where  $F$  is a randomized algorithm that produces an output distribution with  $\omega(\log \lambda)$  entropy. If  $\Pr[x = x' : x \leftarrow \mathcal{A}(\lambda), x' \leftarrow F_{\text{IMPL}}] \leq \text{negl}(\lambda)$  does not hold, then there is a PPT offline watchdog can detect with a non-negligible probability.

**Remark 2.5 Stateless/stateful implementations.** In principle, algorithms in the specification of a cryptographic scheme or implementations provided by an adversary could be stateful; for simplicity, we focus on stateless implementations in the above lemmas. However, to jump ahead a bit, those results still hold (with simple modification) in the stateful settings that matter in this paper. To see this, (1) for randomized algorithm to produce high-entropy output distribution, in the case of stateful implementation which maintains a local state, the unpredictability requirement can still be ensured by an offline watchdog who can rewind the implementation. The watchdog simply samples (rewinds to the same state and then samples) from the randomized algorithm to see whether there is a collision. (2) For deterministic algorithm with a state, as an example, we consider a stateful PRG, where the seed is updated in each iteration. In this case, the public input distribution is evolving during the iterations. One may wonder whether an offline watchdog will be ineffective. However, this is too pessimistic. Observe that the offline watchdog can ensure the consistency of the implementation of PRG when the input is chosen from a uniform distribution. This means the “bad” input set that the implementation deviates from its specification could be at most negligibly small. Note that starting from a uniform seed, any polynomially number of PRG iterations will yield poly-many pseudorandom bits. Thus the probability for any of them falls into the “bad” input set would still be negligible.

**Schemes with augmented system parameter.** Often, deployment of a cryptographic scheme may involve a *system parameter generation* algorithm  $pp \leftarrow \text{Gen}(1^\lambda)$ . When we consider such an augmented scheme  $\Pi = (\text{Gen}, F^1, F^2, F^3)$  in our setting, we can treat the system parameter  $pp$  in two natural ways: (1.) as in Definition 2.2, the adversary simply provides the implementation  $\text{Gen}_{\text{IMPL}}$

to  $\mathcal{W}$  (and  $\mathcal{C}$ ) as usual and the challenger computes  $pp$  by running  $\text{Gen}_{\text{IMPL}}$  during the security game; (2) the *adversary provides*  $pp$  directly to the watchdog  $\mathcal{W}$  (and  $\mathcal{C}$ ); we write  $\mathcal{W}^{\Pi_{\text{IMPL}}}(1^\lambda, pp)$  to reflect this. By replacing  $\mathcal{W}^{\Pi_{\text{IMPL}}}(1^\lambda)$  in Definition 2.2 with  $\mathcal{W}^{\Pi_{\text{IMPL}}}(1^\lambda, pp)$ , and suitably changing the security game so that the challenger does not generate  $pp$ , we can obtain the *adversarially chosen parameter model*. This model was used for studying pseudorandom generator under subversion in [DGG<sup>+</sup>15], we choose to present it as a general model that would be interesting to consider for any cryptographic primitive.

It is clear that if a primitive is secure in the adversarially chosen parameter model, then it is secure according to Definition 2.2. (Observe that the adversary is always free to generate  $pp$  according to the algorithm provided to the challenger.) We record this below.

**Lemma 2.6** *If  $\Pi$  is secure in cryptographic definition in the adversarially chosen parameter model, then  $\Pi$  is secure according to Definition 2.2.*

**Schemes with split-program.** Randomized algorithms play a distinguished role in our kleptographic setting. One technique we propose for immunization may also rely on the decomposition of a randomized generation algorithm  $y \leftarrow \text{Gen}(1^\lambda)$  into two algorithms: a random string generation algorithm RG responsible for producing a uniform  $\text{poly}(\lambda)$ -bit random string  $r$ , and a deterministic output generation algorithm dKG that transforms the randomness  $r$  into an output  $y$ . Note that dKG is deterministic and is always applied to a public input distribution. In light of Lemma 2.3, we may assume that the maliciously implemented  $\text{dKG}_{\text{IMPL}}$  is consistent with the honest implementation  $\text{dKG}_{\text{SPEC}}$  with overwhelming probability. See results in this model in Section 3.3, and definition in supporting material A.

We remark that this perspective only requires a change in the specification of  $\Pi_{\text{SPEC}}$ . When we apply Definition 2.2 with a specification that has been altered to reflect this split-program convention, we say that a primitive is proven secure in the *split-program model*.

*Discussions.* Intuitively, exposing the random coins offers the watchdog more flexibility. It seems a solution to the original scheme  $\Pi$  will immediately give us a solution to the modified scheme  $\Pi^{\text{SP}}$ . However, this is not necessary to be true. The adversary may also take advantage of the “exposed internal state” to break the scheme. The split-program model is quite general and can be applied to most practical algorithms. To see this, the user gets the source code of the implementation, which makes calls to some API for generating randomness (e.g., `rand()`) whenever necessary. The user can hook up the interface with the calls to the API with the separate program RG provided by the big brother. In principle, one can always augment a randomized KG to output the function index  $i$  together with the coin  $r$  in this way.

### 3 Subversion-Resistant One-Way Permutations

In this section, we study one-way permutations (OWP) in our cliptographic framework. As mentioned before, this is motivated by the problem of defending against subverted key generation algorithm. In particular, we propose general constructions for subversion-resistant OWPs that require only the weakest (offline) watchdog with adversarially chosen parameters. Our “immunizing strategy” consists of coupling the function generation algorithm with a hash function that is applied to the function index—intuitively, this makes it challenging for an adversary to mean-

ingly embed a backdoor in the permutation or its index.<sup>3</sup> We prove that if the specification of the hash function is modeled as a random oracle, then randomizing the permutation index using the (adversarially implemented) hash function destroys any potential backdoor structure. We emphasize that the permutation evaluation algorithm, the name generation algorithm, and the hash function may all be subverted by the adversary.

The cliptographic model introduces a number of new perspectives on the (basic) notion of security for one-way permutations. We actually consider three different notions below, each of which correspond to distinct practical settings: the first corresponds to the classical notion, where the challenger chooses the index of the function (using subverted code provided by the adversary)—we call this  $\text{OWP}^C$ ; the second corresponds to a setting where the adversary may choose the index—we call this  $\text{OWP}^A$ ; the last corresponds to our “split program model,” discussed above—we call this  $\text{OWP}^{\text{SP}}$ .

In many cases of practical interest, however, the permutation index may have special algebraic structure, e.g., RSA or DLP. In such cases, it would appear that the public hash function would require some further “structure preserving” property (so that it carries the space of indices to the space of indices). Alternatively, one can assume that the space of indices can be “uniformized,” that is, placed in one-to-one correspondence with strings of a particular length. In order to apply our approach to broader practical settings, we propose a natural “split-program” model that provides such uniformization by insisting that the function generation algorithm is necessarily composed of two parts: a random string generation algorithm  $\text{RG}$  that outputs random bits  $r$ , and a deterministic function index generation algorithm  $\text{dKG}$  which uses  $r$  to generate the index. See Section 3.3 for details.

### 3.1 Defining subversion resistant OWP/TDOWP

In this subsection, following our general definitional framework, we define the security of one-way permutations and trapdoor one-way permutations. We first recall the conventional definitions.

*One-way permutation (OWP).* A family of permutations  $\mathcal{F} = \{f_i : X_i \rightarrow X_i\}_{i \in I}$  is *one-way* if there are  $\text{PPT}$  algorithms  $(\text{KG}, \text{Eval})$  so that (i)  $\text{KG}$ , given a security parameter  $\lambda$ , outputs a function index  $i$  from  $I_\lambda = I \cap \{0, 1\}^\lambda$ ; (ii) for  $x \in X_i$ ,  $\text{Eval}(i, x) = f_i(x)$ ; (iii)  $\mathcal{F}$  is one-way; that is, for any  $\text{PPT}$  algorithm  $\mathcal{A}$ , it holds that  $\Pr[\mathcal{A}(i, y) \in f_i^{-1}(y) \mid i \leftarrow \text{KG}(\lambda); x \leftarrow X_i; y := f_i(x)] \leq \text{negl}(\lambda)$ .

*Trapdoor one-way permutation (TDOWP).* A family of permutations  $\mathcal{F} = \{f_i : X_i \rightarrow X_i\}_{i \in I}$  is *trapdoor one-way* if there are  $\text{PPT}$  algorithms  $(\text{KG}, \text{Eval}, \text{Inv})$  such that (i)  $\text{KG}$ , given a security parameter  $\lambda$ , outputs a function index and the corresponding trapdoor pair  $(i, t_i)$  from  $I_\lambda \times T$ , where  $I_\lambda = I \cap \{0, 1\}^\lambda$ , and  $T$  is the space of trapdoors; (ii)  $\text{Eval}(i, x) = f_i(x)$  for  $x \in X_i$ ; (iii)  $\mathcal{F}$  is one-way; and (iv) it holds that  $\Pr[\text{Inv}(t_i, i, y) = x \mid i \leftarrow \text{KG}(\lambda); x \leftarrow X_i; y := f_i(x)] \geq 1 - \text{negl}(\lambda)$ .

Sometimes, we simply write  $f_i(x)$  rather than  $\text{Eval}(i, x)$ .

**Subversion-resistant<sup>C</sup> one-way permutations:  $\text{OWP}^C$ .** As described in Section 2, we assume a “laboratory specification” of the OWP,  $(\text{KG}_{\text{SPEC}}, \text{Eval}_{\text{SPEC}})$ , which has been rigorously analyzed and certified (e.g., by the experts in the cryptography community). The adversary provides an alternate (perhaps subverted) implementation  $(\text{KG}_{\text{IMPL}}, \text{Eval}_{\text{IMPL}})$ . We study  $\text{OWP}/\text{TDOWP}$  in the offline

<sup>3</sup>In concrete constructions, the hash function becomes a component of e.g., the evaluation function, so that the syntax of the primitive is still the same.

watchdog model; while the implementations may contain arbitrary backdoors or other malicious features, they can not maintain any state.

Intuitively, the goal of the adversary is to privately maintain some “backdoor information”  $z$  so that the subverted implementation  $\text{KG}_{\text{IMPL}}$  will output functions that can be inverted using  $z$ . In addition, to avoid detection by the watchdog, the adversary must ensure that implementations  $(\text{KG}_{\text{IMPL}}(z), \text{Eval}_{\text{IMPL}}(z))$  are computationally indistinguishable from the specification  $(\text{KG}_{\text{SPEC}}, \text{Eval}_{\text{SPEC}})$  given only oracle access. Formally,

**Definition 3.1** A one-way permutation family  $\mathcal{F} = \{f_i : X_i \rightarrow X_i\}_{i \in I}$  with the specification  $\mathcal{F}_{\text{SPEC}} = (\text{KG}_{\text{SPEC}}, \text{Eval}_{\text{SPEC}})$ , is **subversion-resistant**<sup>C</sup> in the offline watchdog model if there exists a PPT watchdog  $\mathcal{W}$ , s.t.: for any PPT adversary  $\mathcal{A}$  playing with the challenger  $\mathcal{C}$  in the following game, (Fig. 1), either the detection probability  $\text{Det}_{\mathcal{W}, \mathcal{A}}$  is non-negligible, or the advantage  $\text{Adv}_{\mathcal{A}}$  is negligible. Here the detection probability of the watchdog  $\mathcal{W}$  with respect to  $\mathcal{A}$  is defined as

$$\text{Det}_{\mathcal{W}, \mathcal{A}}(1^\lambda) = \left| \Pr[\mathcal{W}^{\text{KG}_{\text{IMPL}}, \text{Eval}_{\text{IMPL}}}(1^\lambda) = 1] - \Pr[\mathcal{W}^{\text{KG}_{\text{SPEC}}, \text{Eval}_{\text{SPEC}}}(1^\lambda) = 1] \right|,$$

and the advantage of the adversary  $\mathcal{A}$  is defined as

$$\text{Adv}_{\mathcal{A}}(1^\lambda) = \Pr \left[ (\mathcal{A}(1^\lambda) \stackrel{\leftarrow}{\rightsquigarrow} \mathcal{C}^{\text{KG}_{\text{IMPL}}, \text{Eval}_{\text{IMPL}}}(1^\lambda)) = 1 \right].$$

For convenience, we also say that such  $\mathcal{F}_{\text{SPEC}}$  is a OWP<sup>C</sup> in the offline watchdog model. On the other hand, we say that an OWP is subvertible if:

$\text{Det}_{\mathcal{W}, \mathcal{A}}(1^\lambda)$  is negligible for all PPT  $\mathcal{W}$ , and  $\text{Adv}_{\mathcal{A}}(1^\lambda)$  is non-negligible.<sup>4</sup>

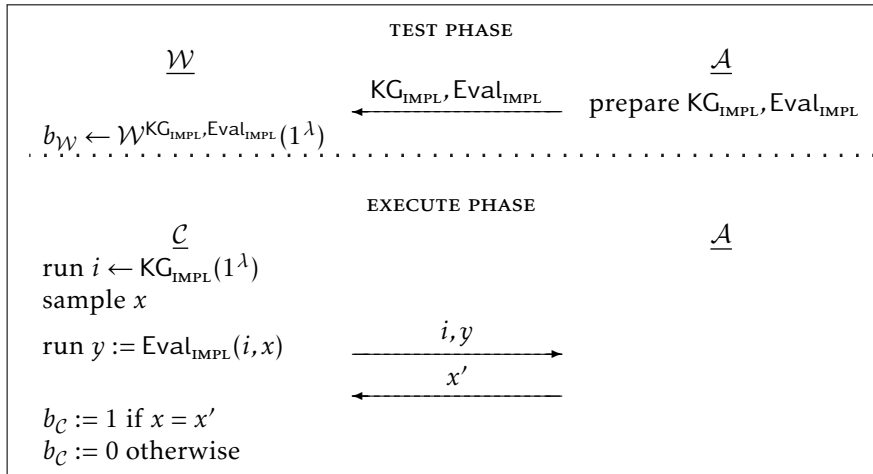


Figure 1: Subversion-resistant<sup>C</sup> security game: OWP<sup>C</sup>.

**Subvertible OWPs.** Next we observe that it is easy for an adversary to break the security of a conventional OWP in the kleptographic setting. In particular, the following lemma shows that one can construct a *subvertible* OWP (so the subverted implementation can evade detection by all watchdogs and the adversary can invert) using a conventional trapdoor OWP. In particular, if we

<sup>4</sup>We choose a stronger definition for subvertibility (swap the quantifiers of  $\mathcal{A}, \mathcal{W}$ ) to describe attacks instead of directly negating the definition of OWP<sup>C</sup>.

wish to use public-key cryptography in a kleptographic setting, nontrivial effort is required to maintain the security of even the most fundamental cryptographic primitives.

Our construction of a subvertible OWP substantiates the folklore knowledge that sufficient random padding can render cryptosystems vulnerable to backdoor attacks, e.g., [YY96, YY97]. Specifically, the random padding in the malicious implementation can be generated so that it encrypts the corresponding trapdoor using the backdoor as a key. For detailed proofs, we defer to supporting material C.1.

**Lemma 3.2** *One can construct a subvertible OWP from any TDOWP. In particular, given a TDOWP, we can construct a OWP that is not a OWP<sup>C</sup>.*

We defer the question of the existence of a OWP<sup>C</sup> to the next section, where we will construct permutations that satisfy a stronger property.

**Subversion-resistant OWPs with adversarially chosen indices: OWP<sup>A</sup>.** The notion of OWP<sup>C</sup> formulated above defends against kleptographic attacks when the adversary provides a subverted implementation of the defining algorithms. In many cases, however, it is interesting to consider a more challenging setting where the adversary may directly provide the public parameters, including the function index. Indeed, this is the case in many real-world deployment settings, where a “trusted” agency sets up (or recommends) the public parameters. One notorious example in a different setting is the Dual\_EC PRG [CNE<sup>+</sup>14]. Note that this notion is not very suitable for asymmetric key primitives, e.g. TDOWP, since allowing the adversary to set up the public key gives him the chance to generate the trapdoor. We will focus on OWP<sup>A</sup>.

**Definition 3.3** *A one-way permutation family  $\mathcal{F} = \{f_i : X_i \rightarrow X_i\}_{i \in I}$  with the specification  $\mathcal{F}_{\text{SPEC}} = (\text{KG}_{\text{SPEC}}, \text{Eval}_{\text{SPEC}})$ , is **subversion-resistant<sup>A</sup>** in the offline watchdog model, if there is a PPT watchdog  $\mathcal{W}$ , such that: for any PPT adversary  $\mathcal{A}$  playing the following game (Fig. 2) with the challenger  $\mathcal{C}$ , either the detection probability  $\text{Det}_{\mathcal{W}, \mathcal{A}}$  is non-negligible, or the advantage  $\text{Adv}_{\mathcal{A}}$  is negligible. Here the detection probability of the watchdog  $\mathcal{W}$  with respect to  $\mathcal{A}$  is defined as:*

$$\text{Det}_{\mathcal{W}, \mathcal{A}}(1^\lambda) = \left| \Pr[\mathcal{W}^{\text{Eval}_{\text{IMPL}}}(1^\lambda, i_\bullet) = 1] - \Pr[\mathcal{W}^{\text{Eval}_{\text{SPEC}}}(1^\lambda, i) = 1] \right|,$$

and the advantage of the adversary  $\mathcal{A}$  is defined as

$$\text{Adv}_{\mathcal{A}}(1^\lambda) = \Pr\left[(\mathcal{A}(1^\lambda) \overset{\text{c}}{\rightsquigarrow} \mathcal{C}^{\text{Eval}_{\text{IMPL}}}(1^\lambda, i_\bullet)) = 1\right].$$

where  $i \leftarrow \text{KG}_{\text{SPEC}}(1^\lambda)$ , and  $i_\bullet$  is chosen by the adversary.

We also say that such  $\mathcal{F}_{\text{SPEC}}$  is a OWP<sup>A</sup> in the offline watchdog model.

**Relating OWP<sup>C</sup> and OWP<sup>A</sup>.** Following Lemma 2.6, an adversary that successfully breaks the OWP<sup>C</sup> game can be easily transformed into an adversary that breaks the OWP<sup>A</sup> game; We can claim that *any OWP<sup>A</sup> is also a OWP<sup>C</sup>*. As far as existence is concerned, then, it suffices to construct a OWP<sup>A</sup> (which satisfies the stronger subversion-resistant<sup>A</sup> condition).

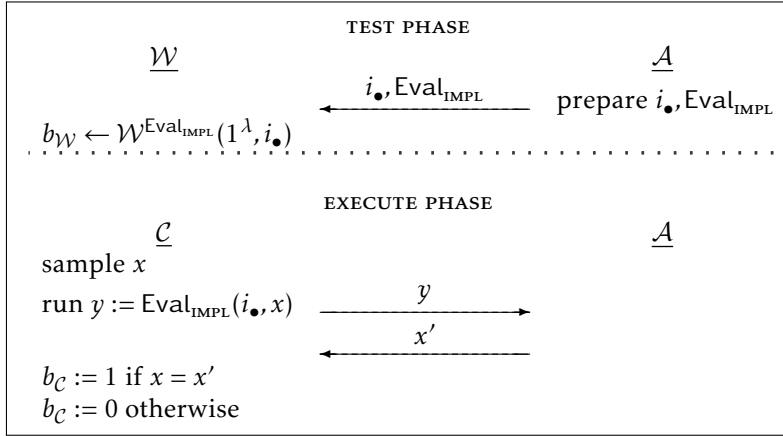


Figure 2: Subversion-resistant<sup>A</sup> security game: OWP<sup>A</sup>.

### 3.2 Constructing subversion-resistant<sup>A</sup> OWP

In this section, we discuss methods for safeguarding OWP against kleptographic attacks. We first present a general approach that transforms any OWP to a OWP<sup>A</sup> under the assumption that a suitable hash function can be defined on the index space. Specifically, we prove that randomizing the function index (via hashing, say) is sufficient to eliminate potential backdoor information. These results assume only the weakest (offline) watchdog. More importantly, we permit the hash function —like the other relevant cryptographic elements—to be implemented and potentially subverted by the adversary.

Note that we treat only the specification of the hash function in the random oracle model, assuming that the adversary may arbitrarily subvert the (randomly specified) hash function; thus the watchdog is provided both the adversary’s arbitrarily subverted “implementation” and the correct (random) hash function “specification.”<sup>5</sup> Despite the adversary’s control over the OWP and the hash function (which is partially constrained by the watchdog), it is difficult for the adversary to arrange a backdoor that works for a large enough target subset of function indices that these can be reliably “hit” by the hash function.

One more difficulty left is to keep the syntax intact, we propose to treat the hash function only as a component of (jumping ahead) the evaluation algorithm (see Fig. 3). The big brother only implements the evaluation algorithm as a whole with the hash function built in. In this case, the hash implementation (and specification) are not explicitly given to the watchdog anymore. However, we still manage to show the security by exploring the fact that both hash and the evaluation algorithm are deterministic algorithms with public input distribution, so that the offline watchdog can force the implementation of  $\text{Eval}_{\text{IMPL}}$  to agree with the specification  $\text{Eval}_{\text{SPEC}}$  with overwhelming probability when inputs are sampled according to the input generation distribution.

**General feasibility results for OWP<sup>A</sup>.** Let  $\mathcal{F}$  be any OWP family with specification  $\mathcal{F}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}^{\mathcal{F}}, \text{Eval}_{\text{SPEC}}^{\mathcal{F}})$ ; while we assume, of course,

<sup>5</sup>Note that we place no a priori constraints on the subverted hash function provided by the adversary. The watchdog, of course, can ensure that the subverted function and the specification (which is just a random function, in this case) agree with high probability on slices of the space, or possess other common statistical properties.

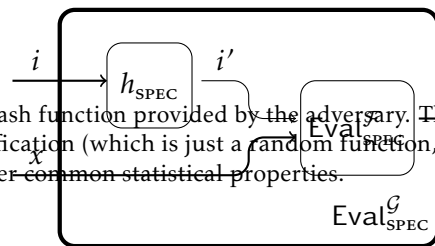


Figure 3: New specification  $\text{Eval}_{\text{SPEC}}^{\mathcal{G}}$ .



that it is one-way secure (in the classical sense), it may be subvertible. We also assume that  $\text{KG}_{\text{SPEC}}^{\mathcal{F}}(\lambda)$  outputs uniform  $i$  from the index set  $I_\lambda$  and that we have a public hash function with specification  $h_{\text{SPEC}} : I_\lambda \rightarrow I_\lambda$ , acting on this set. Then we construct a subversion-resistant<sup>A</sup> OWP family  $\mathcal{G}$  with specification  $\mathcal{G}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}^{\mathcal{G}}, \text{Eval}_{\text{SPEC}}^{\mathcal{G}})$  defined as follows:

- Function index generation  $i \leftarrow \text{KG}_{\text{SPEC}}^{\mathcal{G}}$ , where  $\text{KG}_{\text{SPEC}}^{\mathcal{G}}$  is given by: Sample  $i \leftarrow \text{KG}_{\text{SPEC}}^{\mathcal{F}}(\lambda)$ ; output  $i$ .
- Function evaluation  $y \leftarrow \text{Eval}_{\text{SPEC}}^{\mathcal{G}}(i, x)$ , where  $\text{Eval}_{\text{SPEC}}^{\mathcal{G}}$  is given by: Upon receiving inputs  $(i, x)$ , compute  $i' = h_{\text{SPEC}}(i)$  and compute  $y = \text{Eval}_{\text{SPEC}}^{\mathcal{F}}(i', x)$ ; output  $y$ . See also the pictorial illustration for  $\text{Eval}_{\text{SPEC}}^{\mathcal{G}}$  in Fig. 3.

**Remark 3.4** Note that the specification of the hash function is “part of” of the specification of the evaluation function. In fact, an interesting property of the construction above is that it is secure even if the (subverted) hash function is not separately provided to watchdog.<sup>6</sup>

**Security analysis.** Roughly, the proof relies on the following two arguments: (1.) any particular adversary can only invert a sparse subset of the one-way permutations; otherwise, such an adversary could successfully attack the (classical) security of the specification OWP. Thus, randomizing the function index will map it to a “safe” index, and destroy the possible correlation with any particular backdoor. (2.) The  $\text{Eval}_{\text{IMPL}}$  (having the hash function  $h_{\text{IMPL}}$  built in) is a *deterministic* function that is only called on fixed public input distributions  $(\mathcal{I} \times \mathcal{U})$ , where  $\mathcal{I}$  is the output distribution of  $\text{KG}_{\text{IMPL}}$  and  $\mathcal{U}$  is the uniform distribution over the input space, and both of them are known to the watchdog). Following Lemma 2.3 of Section 2, the watchdog can ensure that  $\text{Eval}_{\text{IMPL}}^{\mathcal{G}}$  is consistent with its specification an overwhelming probability when inputs are generated according to  $\mathcal{I} \times \mathcal{U}$ . We remark that on all inputs for which the hash implementation (running inside  $\text{Eval}_{\text{IMPL}}^{\mathcal{G}}$ ) is consistent with  $h_{\text{SPEC}}$ , random oracle queries have to be made.

**Theorem 3.5** Assume  $h_{\text{SPEC}}$  is random oracle, and  $\mathcal{F}$  with specification  $\mathcal{F}_{\text{SPEC}}$  is a OWP. Then  $\mathcal{G}$  with specification  $\mathcal{G}_{\text{SPEC}}$  defined above is a OWP<sup>A</sup> in the offline watchdog model.

**Proof 1** Suppose that  $\mathcal{G}$  is not subversion-resistant<sup>A</sup>. Then for any watchdog  $\mathcal{W}$ , there is a PPT adversary  $\mathcal{A}_{\mathcal{G}}$  so that the detection probability  $\text{Det}_{\mathcal{W}, \mathcal{A}_{\mathcal{G}}}$  is negligible and the advantage  $\text{Adv}_{\mathcal{A}_{\mathcal{G}}}$  is non-negligible, say  $\delta$ . We will construct an adversary  $\mathcal{A}_{\mathcal{F}}$  which will break the one-way security of  $\mathcal{F}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}^{\mathcal{F}}, \text{Eval}_{\text{SPEC}}^{\mathcal{F}})$  with non-negligible probability. In particular, we define a simple watchdog algorithm that samples a uniform input  $x$ , and compares whether  $\text{Eval}_{\text{SPEC}}^{\mathcal{G}}(i_\bullet, x) = \text{Eval}_{\text{IMPL}}^{\mathcal{G}}(i_\bullet, x)$ , where  $i_\bullet$  is the public parameter chosen by the adversary  $\mathcal{A}_{\mathcal{G}}$  (Note that the evaluation of  $\text{Eval}_{\text{SPEC}}^{\mathcal{G}}$  may involve querying random oracle).

Construction of  $\mathcal{A}_{\mathcal{F}}$ . Suppose  $(i^*, y^*)$  are the challenges that  $\mathcal{A}_{\mathcal{F}}$  receives from the challenger  $\mathcal{C}_{\mathcal{F}}$  (the challenger for one way security of  $\mathcal{F}_{\text{SPEC}}$ ), where  $y^* = \text{Eval}_{\text{SPEC}}^{\mathcal{F}}(i^*, x^*)$  for a randomly selected  $x^*$ .  $\mathcal{A}_{\mathcal{F}}$  simulates a copy of  $\mathcal{A}_{\mathcal{G}}$ . In addition  $\mathcal{A}_{\mathcal{F}}$  simulates the subversion-resistant<sup>A</sup> OWP game with  $\mathcal{A}_{\mathcal{G}}$ .

<sup>6</sup>In general, development of secure primitives in the complete subversion model would presumably be easier if the watchdog can separately “check” the implementation of  $h$  even though we do not need this for the above construction.

Before receiving the function index  $i_\bullet$  and the implementation  $\text{Eval}_{\text{IMPL}}^{\mathcal{G}}$  from  $\mathcal{A}_G$ , the adversary  $\mathcal{A}_F$  (also acting as the challenger in the  $\text{OWP}^A$  game playing with  $\mathcal{A}_G$ ) operates as follows: First, note that  $h_{\text{SPEC}}$  is random oracle; whenever  $\mathcal{A}_G$  wants to evaluate  $h_{\text{SPEC}}$  on some points (or implementing the component for  $\text{Eval}_{\text{IMPL}}^{\mathcal{G}}$  that is consistent with  $h_{\text{SPEC}}$  for those points),  $\mathcal{A}_G$  has to make random oracle queries. Without loss of generality, assume  $\mathcal{A}_G$  asks  $q$  number of random oracle queries on  $i_1, \dots, i_q$  where  $q = \text{poly}(\lambda)$ . Here  $\mathcal{A}_F$  randomly chooses a bit  $b$  to decide whether to embed  $i^*$  to the answers of random oracle queries in this stage. If  $b = 0$ ,  $\mathcal{A}_F$  randomly selects an index  $t \in \{1, \dots, q\}$ , and sets  $i^*$  as the answer for  $h_{\text{SPEC}}(i_t)$ ; for all others  $j \in \{1, \dots, q\} \setminus \{t\}$ ,  $\mathcal{A}_F$  uniformly samples  $i'_j$  from the index set  $I$  and sets  $h_{\text{SPEC}}(i_j) = i'_j$ . If  $b = 1$ , for all  $j \in \{1, \dots, q\}$ , the adversary  $\mathcal{A}_F$  uniformly samples  $i'_j$  from the index set  $I$  and sets  $h_{\text{SPEC}}(i_j) = i'_j$ .

After receiving  $i_\bullet, \text{Eval}_{\text{IMPL}}^{\mathcal{G}}$  from  $\mathcal{A}_G$ , if  $b = 1$  the adversary  $\mathcal{A}_F$  sets  $i^*$  as  $h_{\text{SPEC}}(i_\bullet)$ ; otherwise, it chooses a random value and sets that as  $h_{\text{SPEC}}(i_\bullet)$ . Next,  $\mathcal{A}_F$  gives  $y^*$  to  $\mathcal{A}_G$  as the challenge and receives an answer  $x'$  from  $\mathcal{A}_G$ . Note that in this phase, whenever  $\mathcal{A}_G$  makes random oracle queries on  $i$ , if  $i \in \{i_1, \dots, i_q\} \cup \{i_\bullet\}$ , then returns the previous response as answer; otherwise, randomly choose  $i'$  in the index set  $I$ , and return  $i'$  as the answer.

Last,  $\mathcal{A}_F$  checks whether  $b = 0 \wedge i_\bullet \neq i_t$ , or  $b = 1 \wedge i_\bullet \in \{i_1, \dots, i_q\}$  (in those cases,  $\mathcal{A}_F$  fails to embed  $i^*$  into the right value); if yes,  $\mathcal{A}_F$  aborts; otherwise,  $\mathcal{A}_F$  submits  $x'$  to challenger  $\mathcal{C}_F$  as his answer.

Probabilistic analysis. Now we bound the success probability of  $\mathcal{A}_F$ . Suppose  $x^*$  is the random input chosen by  $\mathcal{C}_F$ ; Let  $W$  denote the event that  $\mathcal{A}_F$  aborts,  $W_1$  the event that  $b = 0 \wedge i \neq i_t$ , and  $W_2$  the event that  $b = 1 \wedge i \in \{i_1, \dots, i_q\}$ . Recall that  $\Pr[x' = x^*] = \Pr[x' = x^* | \overline{W}] \Pr[\overline{W}]$ . Let  $Q = \{i_1, \dots, i_q\}$ .

We first bound  $\Pr[\overline{W}]$ . Note that  $\Pr[\overline{W}] = 1 - \Pr[W]$ , and  $\Pr[W] = \Pr[W_1 \vee W_2] \leq \Pr[W_1] + \Pr[W_2]$ . Assuming  $\Pr[i_\bullet \in Q] = \eta$ , it follows that:

$$\begin{aligned} \Pr[W_1] &= \Pr[b = 0 \wedge i_\bullet \neq i_t] = \Pr[b = 0] \cdot \Pr[i_\bullet \neq i_t] \\ &= \frac{1}{2} (\Pr[i_\bullet \neq i_t | i_\bullet \in Q] \Pr[i_\bullet \in Q] + \Pr[i_\bullet \neq i_t | i_\bullet \notin Q] \Pr[i_\bullet \notin Q]) \\ &= \frac{1}{2} [(1 - (1/q)) \cdot \eta + (1 - \eta)] = \frac{1}{2} (1 - \eta/q) \end{aligned}$$

While  $\Pr[W_2] = \Pr[b = 1] \cdot \Pr[i \in Q] = \eta/2$ , we have:  $\Pr[W] \leq \frac{1}{2} (1 - \frac{\eta}{q} + \eta) = \frac{1}{2} (1 + \eta(1 - \frac{1}{q})) \leq \frac{1}{2} (1 + 1 \cdot (1 - \frac{1}{q})) = 1 - \frac{1}{2q}$ . Thus we can derive that  $\Pr[\overline{W}] \geq 1/(2q)$ .

Next, we bound  $\Pr[x' = x^* | \overline{W}]$ . From the assumption that  $\mathcal{A}_G$  breaks the security of  $\mathcal{G}$ , we have the following two conditions: (1) the detection probability  $\mathbf{Det}_{\mathcal{W}, \mathcal{A}_G}$  is negligible; (2) the advantage  $\mathbf{Adv}_{\mathcal{A}_G}$  is non-negligible  $\delta$ .

From condition (1), we claim:  $\Pr[\text{Eval}_{\text{IMPL}}^{\mathcal{G}}(i_\bullet, x) = \text{Eval}_{\text{SPEC}}^{\mathcal{G}}(i_\bullet, x)] \geq 1 - \text{negl}(\lambda)$ . The probability is over the choices of  $x$  from uniform distribution over the input space. If not, the portion of inputs that  $\text{Eval}_{\text{IMPL}}^{\mathcal{G}}$  deviates from its specification is non-negligible (say,  $\delta_1$ ) in the whole domain. The watchdog  $\mathcal{W}$  we defined (that simply samples an  $x$  and tests if the values  $\text{Eval}_{\text{IMPL}}^{\mathcal{G}}(i_\bullet, x)$  and  $\text{Eval}_{\text{SPEC}}^{\mathcal{G}}(i_\bullet, x)$  are equal) satisfies that  $\Pr[\mathcal{W}^{\text{Eval}_{\text{IMPL}}^{\mathcal{G}}}(1^\lambda, i_\bullet) = 1] = 1 - \delta_1$ . On the other hand,  $\Pr[\mathcal{W}^{\text{Eval}_{\text{SPEC}}^{\mathcal{G}}}(1^\lambda, i) = 1] = 1$ . This implies that  $\mathbf{Det}_{\mathcal{W}, \mathcal{A}_G}$  is  $\delta_1$ , which contradicts to condition (1). Conditioned on  $\overline{W}$ , the equalities:

$$y^* = \text{Eval}_{\text{SPEC}}^{\mathcal{F}}(i^*, x^*) = \text{Eval}_{\text{SPEC}}^{\mathcal{F}}(h_{\text{SPEC}}(i_\bullet), x^*) = \text{Eval}_{\text{SPEC}}^{\mathcal{G}}(i_\bullet, x^*) = \text{Eval}_{\text{IMPL}}^{\mathcal{G}}(i_\bullet, x^*),$$

hold with an overwhelming probability. That said, conditioned on  $\overline{W}$ , from  $\mathcal{A}_G$ 's view, the distribution of  $y^*$  is identical to what she expects as a challenge in the subversion-resistant<sup>A</sup> game.

Recall now from condition (2) the advantage  $\mathbf{Adv}_{\mathcal{A}_G}$  is non-negligible  $\delta$ ; this means  $\mathcal{A}_G$  inverts challenge  $y^* = \text{Eval}_{\text{IMPL}}^G(i_\bullet, x^*)$  and returns a correct  $x' = x^*$  with probability  $\delta$ . Combining the above, we can conclude that:

$$\Pr[x' = x^*] \geq \delta(1 - \text{negl}(\lambda)) \frac{1}{2q} = \frac{\delta}{2q} - \text{negl}(\lambda)$$

which is non-negligible; note that  $q = \text{poly}(\lambda)$ . Thus  $\mathcal{A}_F$  breaks the security of  $\mathcal{F}_{\text{SPEC}}$ , which leads to a contradiction.

### 3.3 Constructing subversion-resistant<sup>SP</sup> OWP/TDOWP

We can define the notion of subversion-resistant<sup>C</sup> TDOWP similar as  $\text{OWP}^C$ . (Note that a *subvertible* TDOWP means that the adversary can invert the TDOWP using a backdoor which may have no relation to the regular trapdoor.) We defer the formal definition to supporting material A.

Indices (names) of a OWP family may have structure. For example, for OWPs based on discrete logarithm,  $f_{g,p}(x) = g^x \bmod p$ , the function index consists of an algebraically meaningful pair  $(p, g)$ , where  $p$  is a prime and  $g$  a random generator. Applying the immunization strategy above would then require a hash function that respects this algebraic structure, mapping meaningful pairs  $(g, p)$  to meaningful pairs  $(g', p')$ . Furthermore, for a TDOWP, we must assume there is a public algorithm that can map between (public key, trapdoor) pairs.

To address these difficulties, we propose a practical *split-program* model in which every function generation algorithm (and, in general, any randomized algorithm) is composed of two components: a “random string generation algorithm”  $\text{RG}$  that outputs a uniform  $\ell$ -bit string  $r$ , and a deterministic function index generation algorithm  $\text{dKG}$  that transforms the randomness  $r$  into a function index  $i$ . In this model,  $\text{dKG}$  is deterministic and is coupled with a known public input distribution (the output distribution of  $\text{RG}$ ). Following Lemma 2.3 and the elaboration in Section 3.1, a watchdog can ensure that the implementation  $\text{dKG}_{\text{IMPL}}$  is “almost consistent” with  $\text{dKG}_{\text{SPEC}}$  (the specification) over this input distribution, i.e.,  $\Pr[\text{dKG}_{\text{IMPL}}(r) = \text{dKG}_{\text{SPEC}}(r) : r \leftarrow \text{RG}_{\text{IMPL}}] \approx 1$ . Morally, this forces the adversary to concentrate his malicious efforts on subverting  $\text{RG}$ .

Since we already demonstrated how to analyze the immunizing strategy for OWP, in this section we present results for  $\text{TDOWP}^{\text{SP}}$ . It is straightforward to adapt the construction and analysis to  $\text{OWP}^{\text{SP}}$ . The standard TDOWP definitions can be easily adapted in the split-program model, where the challenge index is generated by running  $\text{dKG}_{\text{SPEC}}$  on a uniform string  $r$  generated by  $\text{RG}_{\text{SPEC}}$ . It is easy to see that a standard TDOWP is also a TDOWP in the split program model. For detailed definition, we refer to Definition A.2 in supporting material A.

Next we define the notion of a subversion-resistant<sup>SP</sup> TDOWP in the split-program model by simply augmenting Definition 3.1, see Fig. 9 in supporting material A. It is easy to see the same method applies to  $\text{OWP}^{\text{SP}}$  as well. For detailed discussions of  $\text{OWP}^{\text{SP}}$ , we defer to the full version.

**Generic construction of  $\text{TDOWP}^{\text{SP}}$ .** Consider a TDOWP family  $\mathcal{F}$  with specification  $\mathcal{F}_{\text{SPEC}} := (\text{RG}_{\text{SPEC}}^{\mathcal{F}}, \text{dKG}_{\text{SPEC}}^{\mathcal{F}}, \text{Eval}_{\text{SPEC}}^{\mathcal{F}}, \text{Inv}_{\text{SPEC}}^{\mathcal{F}})$ , where  $\text{RG}_{\text{SPEC}}^{\mathcal{F}}$  outputs uniform bits. Assuming a public hash function with specification  $h_{\text{SPEC}} : \{0, 1\}^* \rightarrow \{0, 1\}^*$ , we construct a  $\text{TDOWP}^{\text{SP}}$  family  $\mathcal{G}$  with specification  $\mathcal{G}_{\text{SPEC}} = (\text{RG}_{\text{SPEC}}^{\mathcal{G}}, \text{dKG}_{\text{SPEC}}^{\mathcal{G}}, \text{Eval}_{\text{SPEC}}^{\mathcal{G}}, \text{Inv}_{\text{SPEC}}^{\mathcal{G}})$ , defined below:

- Randomness generation  $r \leftarrow \text{RG}_{\text{SPEC}}^{\mathcal{G}}$ :  $\text{RG}_{\text{SPEC}}^{\mathcal{G}}$  is the same as  $\text{RG}_{\text{SPEC}}^{\mathcal{F}}$ . That is,  $\text{RG}_{\text{SPEC}}^{\mathcal{G}}$  runs  $\text{RG}_{\text{SPEC}}^{\mathcal{F}}$  to get  $r$  and outputs  $r$ .
- Index/trapdoor generation algorithm  $(i, t_i) \leftarrow \text{dKG}_{\text{SPEC}}^{\mathcal{G}}(r)$ : Upon receiving inputs  $r$ , it computes  $\tilde{r} \leftarrow h_{\text{SPEC}}(r)$ , and outputs  $(i, t_i) \leftarrow \text{dKG}_{\text{SPEC}}^{\mathcal{F}}(\tilde{r})$ .

–  $\text{Eval}_{\text{SPEC}}^{\mathcal{G}}, \text{Inv}_{\text{SPEC}}^{\mathcal{G}}$  are the same as  $\text{Eval}_{\text{SPEC}}^{\mathcal{F}}, \text{Inv}_{\text{SPEC}}^{\mathcal{F}}$ .<sup>7</sup>

See also the pictorial description for  $\text{dKG}_{\text{SPEC}}^{\mathcal{G}}$  in Fig. 4:

**Security analysis.** The security of  $\text{OWP}^{\text{SP}}/\text{TDOWP}^{\text{SP}}$  is more subtle than it looks. Randomizing the function index directly indeed destroys any backdoor structure; however, simply randomizing the random coins for generating the function index might lead the adversary to another index/backdoor pair. It will be critical in the split-program model that, with an offline watchdog, the output of  $\text{RG}_{\text{IMPL}}$  is unpredictable even to the adversary who implements it.

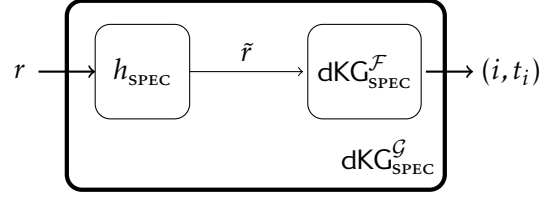


Figure 4: New specification  $\text{dKG}_{\text{SPEC}}^{\mathcal{G}}$ .

A few words about the security proof: Recall that in the  $\text{OWP}^{\text{A}}$  proof, the reduction tries to “program the random oracle” so that the challenge of the specification can be embedded into the challenge to the adversary. In the split-program model, however, the reduction can directly embed the challenge if outputs of  $\text{RG}$  are unpredictable to the adversary; in this case, from the view of the adversary, any random index as challenge is possible to appear in the  $\text{TDOWP}^{\text{SP}}$  game. Therefore, we here do not need to program the random oracle. We defer the full proof to supporting material C.2.

**Theorem 3.6** *Assume  $h_{\text{SPEC}}$  is random oracle, and  $\mathcal{F}$  with specification  $\mathcal{F}_{\text{SPEC}}$  is a TDOWP. Then  $\mathcal{G}$  with specification  $\mathcal{G}_{\text{SPEC}}$  defined above is a  $\text{TDOWP}^{\text{SP}}$  in the offline watchdog model.*

## 4 Subversion-Resistant Signatures

In this section, we consider how to design digital signatures against the kleptographic attacks. Previously results [AMV15, BPR14] suggest that a unique signature scheme [GO93, Lys02] is secure against subversion on the signing algorithm assuming it satisfies the *verifiability* condition that every message signed by the sabotaged  $\text{Sign}_{\text{IMPL}}$  should be verified via  $\text{Verify}_{\text{SPEC}}$ . The main question left is that in those constructions *the key generation and verification algorithms have to be faithfully implemented*, while in practice, all implementations normally come together. Note that all our results obtained in previous sections are actually in the *complete subversion* model. Here we will apply our subversion-resistant TDOWPs to bridge this gap, and construct the first subversion-resistant signature schemes in the more realistic complete subversion model.

We stress that bringing a reduction between two primitives in the classical crypto world to clipto world turns out to be highly non-trivial. We will see that the well-known reduction for full domain hash can not go through in clipto setting when we try to build a subversion-resistant<sup>C</sup> signature from a  $\text{TDOWP}^{\text{C}}$ . (see Remark 4.2 and the proof for Thm. 4.3 for more details).

Following our general framework, it is easy to derive a definition for subversion-resistant signature scheme. As pointed out in [AMV15], it is impossible to achieve unforgeability without the verifiability condition. Using our terminology, it is impossible to construct a subversion-resistant

<sup>7</sup> We remark that in the split-program model, the hash function applies to the random bits, and the hash function is implemented by the adversary inside  $\text{Eval}_{\text{IMPL}}^{\mathcal{G}}$ . The specification of the hash function can be modeled as a random oracle so that replacing the random oracle with an explicit function like SHA256 may be heuristically justified.

signature scheme with just an offline watchdog, even if only the Sign algorithm is subverted. So we will work in the online watchdog model where the watchdog can check the transcripts generated during the game between  $\mathcal{C}$  and  $\mathcal{A}$ .<sup>8</sup> Next we define the security for digital signature scheme in the complete subversion model where all algorithms are implemented by the adversary, including the key generation algorithm.

**Definition 4.1** Given a signature scheme  $\Pi$ , a specification defined as  $\Pi_{\text{SPEC}} = (\text{KG}_{\text{SPEC}}, \text{Sign}_{\text{SPEC}}, \text{Verify}_{\text{SPEC}})$ , is **subversion-resistant<sup>C</sup>** in the online watchdog model, if there exists a PPT watchdog  $\mathcal{W}$ , such that: for any PPT adversary  $\mathcal{A}$  playing the following game (Fig 5) with the challenger  $\mathcal{C}$ , either the detection probability  $\text{Det}_{\mathcal{W}, \mathcal{A}}$  is non-negligible, or the advantage  $\text{Adv}_{\mathcal{A}}$  is negligible.

Here the detection probability of the watchdog  $\mathcal{W}$  with respect to  $\mathcal{A}$  is defined as

$$\text{Det}_{\mathcal{W}, \mathcal{A}}(1^\lambda) = \left| \Pr[\mathcal{W}^{\text{KG}_{\text{IMPL}}, \text{Sign}_{\text{IMPL}}, \text{Verify}_{\text{IMPL}}}(1^\lambda, \tau) = 1] - \Pr[\mathcal{W}^{\text{KG}_{\text{SPEC}}, \text{Sign}_{\text{SPEC}}, \text{Verify}_{\text{SPEC}}}(1^\lambda, \widehat{\tau}) = 1] \right|,$$

and the advantage of the adversary  $\mathcal{A}$  is defined as

$$\text{Adv}_{\mathcal{A}}(1^\lambda) = \Pr[(\mathcal{A}(1^\lambda) \leftrightarrow \mathcal{C}^{\text{KG}_{\text{IMPL}}, \text{Sign}_{\text{IMPL}}, \text{Verify}_{\text{IMPL}}}(1^\lambda)) = 1].$$

where  $\tau$  is the transcript that generated when the challenger uses  $\Pi_{\text{IMPL}}$  and  $\widehat{\tau}$  is the transcript generated when the challenger uses  $\Pi_{\text{SPEC}}$ .

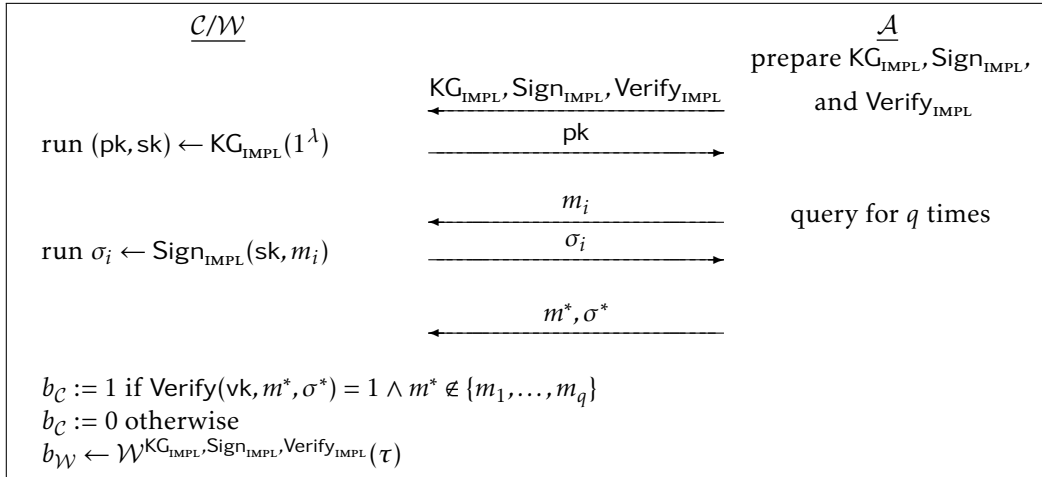


Figure 5: Subversion-resistant<sup>C</sup> Signature Game, where  $\tau := (\text{pk}, \{m_i, \sigma_i\}_{i \in [q]}, m^*, \sigma^*)$

**Discussion.** To upgrade the previous results to the complete subversion model, the main challenging part is to protect the randomized key generation algorithm against subversion attacks. While the attacks shown by Bellare et al [BPR14] are so devastating that the adversary can use a sabotaged randomized implementation to leak any secret bit by bit, we observe that the key generation algorithm will be run only once (as in the security definition). It leaves us the opportunity that if we suggest a better designed specification, the leaked information via the implementations is

<sup>8</sup> Note that for digital signature schemes in practice, it is much more reasonable to assume an online watchdog than an omniscient watchdog as in [BPR14, DFP15]. Indeed, due to the nature of signature schemes, the transcripts consist of message-signature pairs, can be *publicly* verified, and an online watchdog is sufficient.

very limited, so that unforgeability of the signature scheme could still be preserved. Thus we can mitigate the subliminal channel in this way.

Next, we will prove that a variant of the widely deployed full domain hash scheme [Cor00, BR96] can achieve the security in the complete subversion model. More concretely, in this variant, the signing algorithm needs to hash  $m$  together with  $\text{pk}$ ; we remark that this modification is critical for the security reduction (see Remark 4.2). When instantiating its key generation with our subversion-resistant TDOWP<sup>C</sup>, this variant gives a subversion-resistant signature scheme.

**Constructing signature schemes with an online watchdog.** Given a TDOWP<sup>C</sup>, with specification  $\mathcal{F}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}^{\mathcal{F}}, \text{Eval}_{\text{SPEC}}^{\mathcal{F}}, \text{Inv}_{\text{SPEC}}^{\mathcal{F}})$ , and a public hash function with specification  $h_{\text{SPEC}} : \mathcal{PK} \times \mathcal{M} \rightarrow \mathcal{M}$ , where  $\mathcal{PK}$  is the public key space,  $\mathcal{M}$  is the message space, we construct a subversion-resistant<sup>C</sup> signature scheme  $\mathcal{SS}$  with specification  $\mathcal{SS}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}^{\mathcal{SS}}, \text{Sign}_{\text{SPEC}}^{\mathcal{SS}}, \text{Verify}_{\text{SPEC}}^{\mathcal{SS}})$  as follows:

- Key generation  $(\text{pk}, \text{sk}) \leftarrow \text{KG}_{\text{SPEC}}^{\mathcal{SS}}(\lambda)$ , where  $\text{KG}_{\text{SPEC}}^{\mathcal{SS}}$  is given by:  
Compute  $(i, t_i) \leftarrow \text{KG}_{\text{SPEC}}^{\mathcal{F}}(\lambda)$ , and set  $\text{pk} := i$  and  $\text{sk} := t_i$ ;
- Signature generation  $\sigma \leftarrow \text{Sign}_{\text{SPEC}}^{\mathcal{SS}}(\text{pk}, \text{sk}, m)$ , where  $\text{Sign}_{\text{SPEC}}^{\mathcal{SS}}$  is given by:  
Upon receiving message  $m$ , compute  $\tilde{m} = h_{\text{SPEC}}(\text{pk}, m)$ , and  $\sigma = \text{Inv}_{\text{SPEC}}^{\mathcal{F}}(\text{pk}, \text{sk}, \tilde{m})$ , where  $\text{pk} = i, \text{sk} = t_i$ .
- Verification algorithm  $b \leftarrow \text{Verify}_{\text{SPEC}}^{\mathcal{SS}}(\text{pk}, m, \sigma)$ , where  $\text{Verify}_{\text{SPEC}}^{\mathcal{SS}}$  is given by:  
Upon receiving message-signature pair  $(m, \sigma)$ , if  $\text{Eval}_{\text{SPEC}}^{\mathcal{F}}(\text{pk}, \sigma) = h_{\text{SPEC}}(\text{pk}, m)$  then set  $b := 1$ , otherwise set  $b := 0$ ; here  $\text{pk} = i$ .

**Remark 4.2** We emphasize that using the full domain hash directly without inputting  $\text{pk}$  into the hash, it is possible that the random oracle query for  $m^*$  is asked before the implementations are prepared. In this case, the simulator has not yet received  $y^*$  from the TDOWP challenger, thus the simulator has no way to embed  $y^*$  into the target. Including the  $\text{pk}$  in the hash essentially makes the random oracle queries made before the implementation are provided essentially useless (unrelated to any of the signature), since the adversary cannot predict the actual value of  $\text{pk}$ . We defer the detailed proof to supporting material C.3.

**Theorem 4.3** Assume  $h_{\text{SPEC}}$  is random oracle, and  $\mathcal{F}$  with specification  $\mathcal{F}_{\text{SPEC}}$  is a TDOWP<sup>C</sup> in the offline watchdog model. Then the signature scheme  $\mathcal{SS}$  with specification  $\mathcal{SS}_{\text{SPEC}}$  constructed above is subversion-resistant<sup>C</sup> in the online watchdog model.

## 5 Subversion-Resistant Pseudorandom Generators

Having studied the classical fundamental building blocks (OWPs and TDOWPs) in the complete subversion model, we now attempt to mimic the classical program of constructing richer cryptographic primitives from OWP/TDOWPs. We proceed in two different ways. The first is to carry “black-box” construction, the second is to generalize the immunizing strategy to broader settings. We remark that typical “black-box” constructions and reductions may not survive in the cryptographic model (indeed, even such basic features as the presence of multiple calls to a randomized algorithm can significantly affect security [BPR14].) We begin by focusing on pseudorandom generators (PRG).

We first review the basic notions of PRG under subversion attacks and then provide a specific construction that mimics the classical Blum-Micali PRG construction in this cryptographic

context. We then examine how to extend the applicability of our general sanitizing strategy for OWP/TDOWPs to more general settings, demonstrating a strategy of public immunization for PRGs. Note that an impossibility result was shown in [DGG<sup>+</sup>15] that no public immunizing strategy is possible if it is applied to the output of the backdoored PRG, so a solution involves some trusted randomness is proposed. We also remark that all algorithms in our backdoor-free PRG construction—including the sanitizing function (which can be part of the KG algorithm in the specification)—can be subverted. Thus we provide the first PRG constructions secure in the complete subversion model.

We remark that since we follow the formalization of [DGG<sup>+</sup>15], the stretching algorithm is deterministic and stateful. In this case, the input distribution is evolving and not fixed, a universal watchdog cannot exhaust all those distributions. Fortunately, in the case of PRG stretching algorithm, we can still establish such a result, see security analysis of theorem 5.3.

### 5.1 Preliminaries: The definition of a subversion-resistant<sup>A</sup> PRG

We adopt the definition from [DGG<sup>+</sup>15]: a pseudorandom generator consists of a pair of algorithms (KG, PRG), where KG outputs a public parameter  $pk$  and  $\text{PRG} : \{0, 1\}^* \times \{0, 1\}^\ell \rightarrow \{0, 1\}^\ell \times \{0, 1\}^{\ell'}$  takes the public parameter  $pk$  and an  $\ell$ -bit random seed  $s$  as input; it returns a state  $s_1 \in \{0, 1\}^\ell$  and an output string  $r_1 \in \{0, 1\}^{\ell'}$ . PRG may be iteratively executed; in the  $i$ -th iteration, it takes the state from the previous iteration  $s_{i-1}$  as the seed and generates the current state  $s_i$  and output  $r_i$ . We use  $q$ -PRG to denote the result of  $q$  iterations of PRG with outputs  $r_1, \dots, r_q$  (each  $r_i \in \{0, 1\}^{\ell'}$ ).

They then define the notion of a *backdoored PRG* [DGG<sup>+</sup>15]: the adversary sets up a public parameter  $pk$  and may keep the corresponding backdoor  $sk$ . The output distribution  $\text{PRG}(pk, \mathcal{U})$  must still look pseudorandom to all algorithms that do not hold the backdoor  $sk$  (e.g., it fools the watchdog), where  $\mathcal{U}$  is the uniform distribution; however, with  $sk$ , the adversary is able to distinguish the output from a uniform string, breaking the PRG.

The definition of a backdoored-PRG [DGG<sup>+</sup>15] is closely related to the subversion-resistant<sup>A</sup> definition in our definitional framework, as the adversary is empowered to choose the “index”  $pk$ . Although there are several variants that all appear meaningful and interesting for PRG in the cliptographic settings, we will initially focus on the subversion-resistant<sup>A</sup> PRG as the striking real world example of Dual\_EC subversion is indeed in this model. Additionally, from Lemma 2.6, we remark that any  $\text{PRG}^A$  is a  $\text{PRG}^C$ .

We first reformulate the definition of [DGG<sup>+</sup>15] in the subversion-resistant<sup>A</sup> cliptographic framework: There exist “specification” versions of the algorithms and an offline watchdog. The parameter generation algorithm  $\text{KG}_{\text{SPEC}}$  has the requirement that the distribution of the adversarially generated public parameter must be indistinguishable from the output distribution of  $\text{KG}_{\text{SPEC}}$ . Additionally, as the PRG algorithm is deterministic, and its input distribution is public, a offline watchdog can ensure that it is consistent with its specification  $\text{PRG}_{\text{SPEC}}$  on an overwhelming fraction of the inputs. The formal definitions are as follows:

**Definition 5.1** *We say that a PRG (with the specification  $(\text{KG}_{\text{SPEC}}, \text{PRG}_{\text{SPEC}})$ ) is  $q$ -subversion-resistant<sup>A</sup> in the offline watchdog model if, there exists a PPT watchdog  $\mathcal{W}$  such that: for any PPT adversary  $\mathcal{A}$  playing the following game (Fig. 6) with the challenger  $\mathcal{C}$ , either the detection probability  $\text{Det}_{\mathcal{W}, \mathcal{A}}$  is non-negligible, or the advantage  $\text{Adv}_{\mathcal{A}}$  is negligible.*

*Here the detection probability of the watchdog  $\mathcal{W}$  with respect to  $\mathcal{A}$  is defined as*

$$\text{Det}_{\mathcal{W}, \mathcal{A}}(1^\lambda) = \left| \Pr[\mathcal{W}^{\text{PRG}_{\text{IMPL}}}(1^\lambda, pk_\bullet) = 1] - \Pr[\mathcal{W}^{\text{PRG}_{\text{SPEC}}}(1^\lambda, pk) = 1] \right|,$$

and the advantage of the adversary  $\mathcal{A}$  is defined as

$$\mathbf{Adv}_{\mathcal{A}}(1^\lambda) = \left| \Pr \left[ (\mathcal{A}(1^\lambda) \leftrightarrow \mathcal{C}^{\text{PRG}_{\text{IMPL}}}(1^\lambda, pk_\bullet)) = 1 \right] - \frac{1}{2} \right|.$$

where  $pk \leftarrow \text{KG}_{\text{SPEC}}(1^\lambda)$ , and  $\text{PRG}_{\text{IMPL}}, pk_\bullet$  are chosen by the adversary.

We say that such PRG is a  $\text{PRG}^{\text{A}}$  to stress that the public parameters are generated by the adversary.

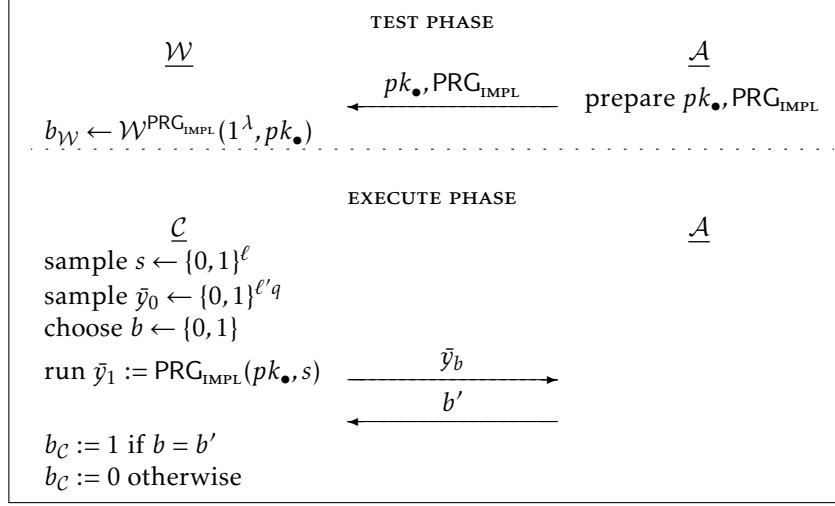


Figure 6: Subversion-resistant<sup>A</sup> PRG Game

## 5.2 Constructing $q$ -PRG<sup>A</sup> from a OWP<sup>A</sup>

In this section, we provide constructions of a PRG<sup>A</sup> based on a OWP<sup>A</sup>. We start with a construction based on a (simplified) Blum-Micali PRG, and then extend it to a full-fledged solution. We remark that a similar reduction can be used to construct a subversion-resistant<sup>C</sup> PRG from a subversion-resistant<sup>C</sup> OWP (where the challenger queries the KG implementation to choose a public parameter).

Before describing the details of our construction, we recall the classic generic construction of Goldreich-Levin (GL), yielding a hardcore predicate [GL89] for any OWF  $f$ . We suppose the input  $x$  of  $f$  is divided into two halves  $x = (x_1, x_2)$  and define the bit  $B(x) = \langle x_1, x_2 \rangle$ ;  $B(x)$  is hard to predict given  $x_1, f(x_2)$ , assuming that  $f$  is one-way. Moreover, if there is a PPT algorithm that predicts  $B(x)$  with significant advantage  $\delta$  given  $x_1, f(x_2)$ , then there is a PPT algorithm  $I$  that inverts  $f$  with probability  $\text{poly}(\delta)$ .

**Basic construction.** We will show that given a subversion-resistant<sup>A</sup> one-way permutation (OWP) family  $\mathcal{F}$  with specifications and implementations  $\mathcal{F}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}^{\mathcal{F}}, \text{Eval}_{\text{SPEC}}^{\mathcal{F}})$  and  $(\text{KG}_{\mathcal{F}}, \text{Eval}_{\mathcal{F}})$  respectively, the classic Blum-Micali PRG [BM82] (using the GL hardcore predicate) is 1-subversion-resistant<sup>A</sup>. Our basic construction  $\mathcal{G}$  with the specification  $\mathcal{G}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}^{\mathcal{G}}, \text{PRG}_{\text{SPEC}}^{\mathcal{G}})$  is as follows:

- Parameter generation algorithm  $pk \leftarrow \text{KG}_{\text{SPEC}}^{\mathcal{G}}(\lambda)$ : compute  $i \leftarrow \text{KG}_{\text{SPEC}}^{\mathcal{F}}(\lambda)$  and set  $pk := i$ ;
- Bit string generation algorithm  $(s', b) \leftarrow \text{PRG}(pk, s)$ : upon receiving  $s$  and  $pk$ , where  $pk = i$ ,  $s = s_1 \| s_2$  and  $|s_1| = |s_2| = \ell$ , compute the following:  $s'_1 := s_1$ ,  $s'_2 := \text{Eval}_{\text{SPEC}}^{\mathcal{F}}(i, s_2)$ , and  $s' = s'_1 \| s'_2$ ,  $b := \langle s_1, s_2 \rangle$ .



**Security analysis.** We can show in the lemma below that, with a specification designed as above, the basic construction above is a 1-subversion-resistant PRG. The intuition is that in the (simplified) Blum-Micali PRG, a distinguisher can be transformed into an OWP inverter (following the GL proof); thus an adversary who can build a backdoor for this PRG violates the subversion-resistance of  $\mathcal{F}$ . We present the lemma for its security, while due to lack of space, we refer the detailed proof to supporting materials C.4.

**Lemma 5.2** *If  $\mathcal{F}$  with specification  $\mathcal{F}_{\text{SPEC}}$  is a OWP<sup>A</sup> in the offline watchdog model, then  $\mathcal{G}$  with specification  $\mathcal{G}_{\text{SPEC}}$  constructed above is a 1-subversion-resistant<sup>A</sup> PRG in the offline watchdog model.*

**Full-fledged PRG<sup>A</sup>.** We can easily adapt our basic construction to the full-fledged PRG<sup>A</sup> construction via the iteration as the BM-PRG and argue the security following the classic hybrid lemma. We refer the details of the construction and analysis to supporting material B.1.

### 5.3 A general public immunization strategy for PRG<sup>A</sup>

An impossibility result concerning public immunization of a PRG (to yield a PRG<sup>A</sup>) was presented in [DGG<sup>+</sup>15]. However, we observe that this impossibility result only applies to an immunization procedure that operates on the *output* of the PRG<sup>A</sup>. The general construction of OWP<sup>A</sup> shown above inspires us to consider an alternate general immunizing strategy for (potentially subvertible) PRGs. We establish that—similar to the procedure above for eliminating backdoors in OWPs—one can randomize the public parameter to sanitize a PRG.<sup>9</sup>

The intuition for this strategy to be effective in the setting of PRG is similar: considering a specification  $\text{KG}_{\text{SPEC}}$  that outputs a uniform  $pk$  from its domain, no single backdoor can be used to break the security for a large fraction of public parameter space; otherwise, one can use this trapdoor to break the PRG security of the specification. As above, while the adversary can subvert the hash function, an offline watchdog can ensure the hash function is faithful enough to render it difficult for the adversary arrange for the result of the hashed parameter to be amenable to any particular backdoor.

Consider a (potentially subvertible) PRG with specification  $\mathcal{F}_{\text{SPEC}} = (\text{KG}_{\text{SPEC}}^{\mathcal{F}}, \text{PRG}_{\text{SPEC}}^{\mathcal{F}})$ ; we assume that  $\text{KG}_{\text{SPEC}}^{\mathcal{F}}$  outputs a uniform element of its range  $PP$ . Consider hash function with specification  $h_{\text{SPEC}} : PP \rightarrow PP$ . Then we construct a PRG<sup>A</sup>  $\mathcal{G}$  with its specification  $\mathcal{G}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}^{\mathcal{G}}, \text{PRG}_{\text{SPEC}}^{\mathcal{G}})$ :

- Parameter generation algorithm  $pk \leftarrow \text{KG}_{\text{SPEC}}^{\mathcal{G}}$ : Compute  $\text{KG}_{\text{SPEC}}^{\mathcal{F}}$ , resulting in the output  $pk$ ;
- Bit string stretch algorithm  $(s', r) \leftarrow \text{PRG}_{\text{SPEC}}^{\mathcal{G}}(pk, s)$  which is given by: Upon receiving a random seed  $s$  and public keys  $pk$  as inputs, it computes  $\tilde{pk} = h_{\text{SPEC}}(pk)$  and it computes  $\text{PRG}_{\text{SPEC}}^{\mathcal{F}}(\tilde{pk}, s)$  and obtains  $s', r$  as outputs, where  $r$  would be the actual output, while  $s'$  would be used as the seed for next iteration. See also the pictorial illustration for  $\text{PRG}_{\text{SPEC}}^{\mathcal{G}}$  in Fig 7.

**Security analysis.** If the above PRG only iterates once, the security analysis would be very similar to that of Theorem 3.5; since any potential backdoor embedded in the public parameter is now destroyed, and the stretch algorithm is a deterministic algorithm with a public input distribution, thus an offline watchdog can already ensure it to be (essentially) consistent with its specification.

Things become trickier when the PRG may be iterated with arbitrary number of times. For example, suppose the watchdog checks only for  $t$  iterations,  $\text{PRG}_{\text{IMPL}}$  might deviate from the

<sup>9</sup>To interpret this results, the solution of [DGG<sup>+</sup>15] is in a semi-private model which requires a trusted seed/key generation, thus part of the PRG algorithms can not be subverted. It follows that the construction of PRG in the complete subversion model was still open until our solution. In contrast, our sanitizing strategy does not require any secret, and even the deterministic hash function can be implemented by the adversary as part of the KG algorithm.

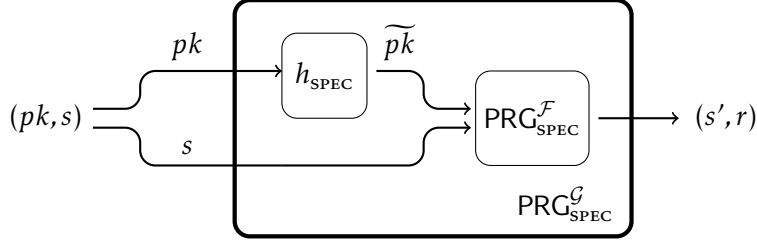


Figure 7: Public immunization strategy for PRG.

$t + 1$ -th iteration. This might be indeed problematic for general deterministic algorithms. Fortunately, for this particular example of PRG, the watchdog simply checks for one uniform input and compares the output with that generated by the specification is enough to ensure almost-everywhere consistency. To see this, the adversary can create a subset of inputs  $S = \{s\}$ , such that:  $\text{PRG}_{\text{IMPL}}(pk, s) \neq \text{PRG}_{\text{SPEC}}(pk, s)$ , where  $pk$  is the adversarially generated public parameter. Observe that the probability that a randomly chosen input  $s$  falls in  $S$  would be negligible. Otherwise the watchdog can detect with a non-negligible probability. While the difference with a stateful stretching algorithm is that it offers the adversary more chances to hit the bad set  $S$  because of the iterations. Note that when  $\text{PRG}_{\text{IMPL}}(pk, s) = \text{PRG}_{\text{SPEC}}(pk, s)$  for some randomly chosen  $s$ , then the output  $s'$  would also be pseudorandom; iterating on this input, the stretching algorithm yields a polynomially many pseudorandom strings, thus the probability of any of those hit the bad set  $S$  would be still negligible. With this observation, we can still claim that with an overwhelming probability,  $\text{PRG}_{\text{IMPL}}$  will be consistent with  $\text{PRG}_{\text{SPEC}}$  even after arbitrary number of iterations (polynomially bounded). We defer the proof to the full version.

**Theorem 5.3** *Assume  $h_{\text{SPEC}}$  is random oracle, and  $\mathcal{F}$  with specification  $\mathcal{F}_{\text{SPEC}} = (\text{KG}_{\text{SPEC}}^{\mathcal{F}}, \text{PRG}_{\text{SPEC}}^{\mathcal{F}})$  is a pseudorandom generator, where  $\text{KG}_{\text{SPEC}}^{\mathcal{F}}$  outputs  $pk$  randomly from its range. Then  $\mathcal{G}$  with specification  $\mathcal{G}_{\text{SPEC}}$  in the above construction yields a  $q$ -subversion-resistant  $\text{PRG}^{\mathcal{A}}$  for any polynomially large  $q$ .*

**Remark 5.4** *If the public parameter contains only random group elements, e.g., the Dual\_EC PRG, we may simply encode them into bits and use a regular hash function like SHA-256, and convert the resulting bits back to a group element;*

## References

- [AMV15] Giuseppe Ateniese, Bernardo Magri, and Daniele Venturi. Subversion-resilient signature schemes. In Indrajit Ray, Ninghui Li, and Christopher Kruegel, editors, *ACM CCS 15*, pages 364–375. ACM Press, October 2015.
- [BH15] Mihir Bellare and Viet Tung Hoang. Resisting randomness subversion: Fast deterministic and hedged public-key encryption in the standard model. In Elisabeth Oswald and Marc Fischlin, editors, *EUROCRYPT 2015, Part II*, volume 9057 of *LNCS*, pages 627–656. Springer, Heidelberg, April 2015.
- [BJK15] Mihir Bellare, Joseph Jaeger, and Daniel Kane. Mass-surveillance without the state: Strongly undetectable algorithm-substitution attacks. In Indrajit Ray, Ninghui Li, and Christopher Kruegel, editors, *ACM CCS 15*, pages 1431–1440. ACM Press, October 2015.
- [BM82] Manuel Blum and Silvio Micali. How to generate cryptographically strong sequences of pseudo random bits. In *23rd FOCS*, pages 112–117. IEEE Computer Society Press, November 1982.
- [BPR14] Mihir Bellare, Kenneth G. Paterson, and Phillip Rogaway. Security of symmetric encryption against mass surveillance. In Juan A. Garay and Rosario Gennaro, editors, *CRYPTO 2014, Part I*, volume 8616 of *LNCS*, pages 1–19. Springer, Heidelberg, August 2014.

- [BR96] Mihir Bellare and Phillip Rogaway. The exact security of digital signatures: How to sign with RSA and Rabin. In Ueli M. Maurer, editor, *EUROCRYPT'96*, volume 1070 of *LNCS*, pages 399–416. Springer, Heidelberg, May 1996.
- [CNE<sup>+</sup>14] Stephen Checkoway, Ruben Niederhagen, Adam Everspaugh, Matthew Green, Tanja Lange, Thomas Ristenpart, Daniel J. Bernstein, Jake Maskiewicz, Hovav Shacham, and Matthew Fredrikson. On the practical exploitability of dual EC in TLS implementations. In *Proceedings of the 23rd USENIX Security Symposium, San Diego, CA, USA, August 20-22, 2014.*, pages 319–335, 2014.
- [Cor00] Jean-Sébastien Coron. On the exact security of full domain hash. In Mihir Bellare, editor, *CRYPTO 2000*, volume 1880 of *LNCS*, pages 229–235. Springer, Heidelberg, August 2000.
- [DFP15] Jean Paul Degabriele, Pooya Farshim, and Bertram Poettering. A more cautious approach to security against mass surveillance. In Gregor Leander, editor, *FSE 2015*, volume 9054 of *LNCS*, pages 579–598. Springer, Heidelberg, March 2015.
- [DGG<sup>+</sup>15] Yevgeniy Dodis, Chaya Ganesh, Alexander Golovnev, Ari Juels, and Thomas Ristenpart. A formal treatment of backdoored pseudorandom generators. In Elisabeth Oswald and Marc Fischlin, editors, *EUROCRYPT 2015, Part I*, volume 9056 of *LNCS*, pages 101–126. Springer, Heidelberg, April 2015.
- [DMSD15] Yevgeniy Dodis, Ilya Mironov, and Noah Stephens-Davidowitz. Message transmission with reverse firewalls—secure communication on corrupted machines. *Cryptology ePrint Archive*, Report 2015/548, 2015. <http://eprint.iacr.org/2015/548>.
- [GL89] Oded Goldreich and Leonid A. Levin. A hard-core predicate for all one-way functions. In *21st ACM STOC*, pages 25–32. ACM Press, May 1989.
- [GO93] Shafi Goldwasser and Rafail Ostrovsky. Invariant signatures and non-interactive zero-knowledge proofs are equivalent (extended abstract). In Ernest F. Brickell, editor, *CRYPTO'92*, volume 740 of *LNCS*, pages 228–245. Springer, Heidelberg, August 1993.
- [HH09] Iftach Haitner and Thomas Holenstein. On the (im)possibility of key dependent encryption. In Omer Reingold, editor, *TCC 2009*, volume 5444 of *LNCS*, pages 202–219. Springer, Heidelberg, March 2009.
- [HLv02] Nicholas J. Hopper, John Langford, and Luis von Ahn. Provably secure steganography. In Moti Yung, editor, *CRYPTO 2002*, volume 2442 of *LNCS*, pages 77–92. Springer, Heidelberg, August 2002.
- [JG02] Ari Juels and Jorge Guajardo. RSA key generation with verifiable randomness. In David Naccache and Pascal Paillier, editors, *PKC 2002*, volume 2274 of *LNCS*, pages 357–374. Springer, Heidelberg, February 2002.
- [Lys02] Anna Lysyanskaya. Unique signatures and verifiable random functions from the DH-DDH separation. In Moti Yung, editor, *CRYPTO 2002*, volume 2442 of *LNCS*, pages 597–612. Springer, Heidelberg, August 2002.
- [MS15] Ilya Mironov and Noah Stephens-Davidowitz. Cryptographic reverse firewalls. In Elisabeth Oswald and Marc Fischlin, editors, *EUROCRYPT 2015, Part II*, volume 9057 of *LNCS*, pages 657–686. Springer, Heidelberg, April 2015.
- [NIS12] NIST. Special publication 800-90: Recommendation for random number generation using deterministic random bit generators. National Institute of Standards and Technology, 2012. <http://csrc.nist.gov/publications/PubsSPs.html>.
- [PLS13] Nicole Perlroth, Jeff Larson, and Scott Shane. N.S.A. able to foil basic safeguards of privacy on web. *The New York Times*, 2013. <http://www.nytimes.com/2013/09/06/us/nsa-foils-much-internet-encryption.html>.
- [Rog15] Phillip Rogaway. The moral character of cryptographic work. *Cryptology ePrint Archive*, Report 2015/1162, 2015. <http://eprint.iacr.org/2015/1162>.
- [Sim83] Gustavus J. Simmons. The prisoners' problem and the subliminal channel. In David Chaum, editor, *CRYPTO'83*, pages 51–67. Plenum Press, New York, USA, 1983.
- [Sim86] Gustavus J. Simmons. A secure subliminal channel (?). In Hugh C. Williams, editor, *CRYPTO'85*, volume 218 of *LNCS*, pages 33–41. Springer, Heidelberg, August 1986.
- [Wik] Wikipedia. Nothing up my sleeve. [https://en.wikipedia.org/wiki/Nothing\\_up\\_my\\_sleeve\\_number](https://en.wikipedia.org/wiki/Nothing_up_my_sleeve_number).
- [YY96] Adam Young and Moti Yung. The dark side of “black-box” cryptography, or: Should we trust capstone? In Neal Koblitz, editor, *CRYPTO'96*, volume 1109 of *LNCS*, pages 89–103. Springer, Heidelberg, August 1996.
- [YY97] Adam Young and Moti Yung. Kleptography: Using cryptography against cryptography. In Walter Fumy, editor, *EUROCRYPT'97*, volume 1233 of *LNCS*, pages 62–74. Springer, Heidelberg, May 1997.

## A Supporting Material: Omitted Definitions

**Definition A.1** A trapdoor OWP family  $\mathcal{F} = \{f_i : X_i \rightarrow Y_i\}_{i \in I}$  with specification  $(\text{KG}_{\text{SPEC}}, \text{Eval}_{\text{SPEC}}, \text{Inv}_{\text{SPEC}})$ , is **subversion-resistant<sup>C</sup>** in the offline watchdog model, if there exists a watchdog  $\mathcal{W}$ , such that: for any PPT adversary  $\mathcal{A}$  playing with the challenger  $\mathcal{C}$  in the following game, (Fig 8), either the detection probability  $\text{Det}_{\mathcal{W}, \mathcal{A}}$  is non-negligible, or the advantage  $\text{Adv}_{\mathcal{A}}$  is negligible. Here the detection probability of the watchdog  $\mathcal{W}$  with respect to  $\mathcal{A}$  is defined as

$$\text{Det}_{\mathcal{W}, \mathcal{A}}(1^\lambda) = \left| \Pr[\mathcal{W}^{\text{KG}_{\text{IMPL}}, \text{Eval}_{\text{IMPL}}, \text{Inv}_{\text{IMPL}}}(1^\lambda) = 1] - \Pr[\mathcal{W}^{\text{KG}_{\text{SPEC}}, \text{Eval}_{\text{SPEC}}, \text{Inv}_{\text{SPEC}}}(1^\lambda) = 1] \right|,$$

and the advantage of the adversary  $\mathcal{A}$  is defined as

$$\text{Adv}_{\mathcal{A}}(1^\lambda) = \Pr\left[(\mathcal{A}(1^\lambda) \stackrel{\leftarrow}{\rightsquigarrow} \mathcal{C}^{\text{KG}_{\text{IMPL}}, \text{Eval}_{\text{IMPL}}, \text{Inv}_{\text{IMPL}}}(1^\lambda)) = 1\right].$$

For convenience, we also say that such  $\mathcal{F}_{\text{SPEC}}$  is a TDOWP<sup>C</sup> in the offline watchdog model.

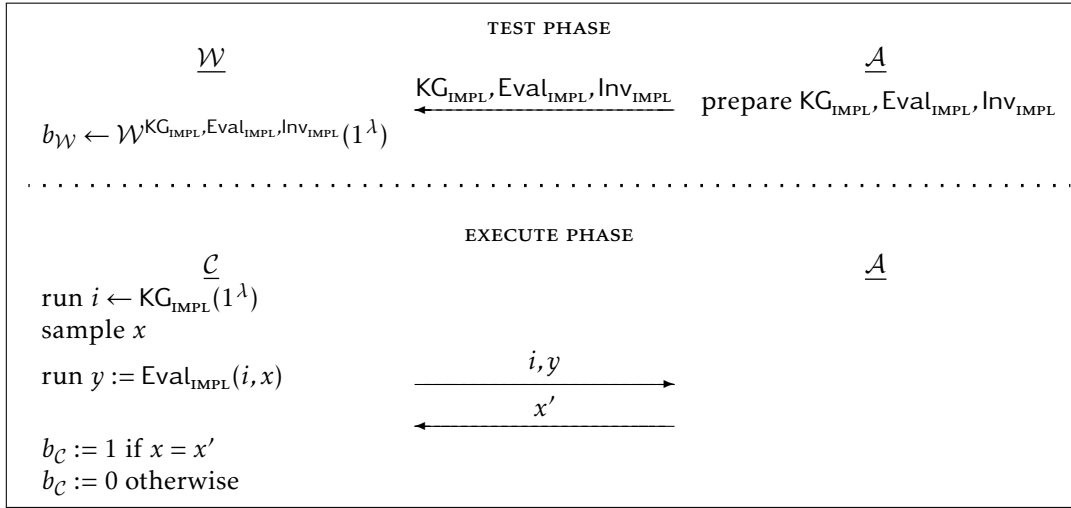


Figure 8: Subversion-resistant<sup>C</sup> TDOWP Game

**Definition A.2** A function family  $\mathcal{F}$  is trapdoor one way in the split-program model if there exist a pair of algorithms  $(\text{RG}, \text{dKG}, \text{Eval}, \text{Inv})$  where (i.)  $\text{RG}$ , given a security parameter  $\lambda$ , outputs a uniform  $\ell(\lambda)$ -bit string  $r$ ; (ii.)  $\text{dKG}$  is deterministic: given the randomness  $r$  it outputs a function index and trapdoor pair  $i, t$ ; and (iii.)  $\mathcal{F}$  is one-way under this procedure for generating  $(i, t) \leftarrow \text{dKG}(r) : r \leftarrow \text{RG}$ ; (4).  $\forall x, \text{Inv}(t, \text{Eval}(i, x)) = x$ .

## B Supporting Material: Omitted Constructions

### B.1 Full Fledged Subversion Resistant PRG<sup>A</sup>

We now extend our basic construction via iteration to show that the full-fledged Blum-Micali PRG construction, using a subversion-resistant<sup>A</sup> OWP, achieves a  $q$ -subversion-resistant<sup>A</sup> PRG for any

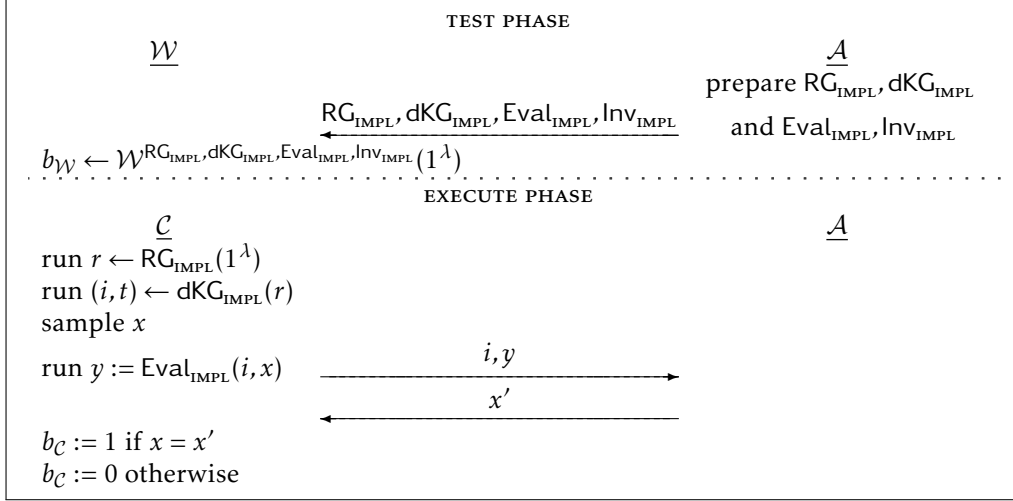


Figure 9: Subversion-resistant<sup>SP</sup> TDOWP Game

$q = \text{poly}(\lambda)$ . Given a subversion-resistant<sup>A</sup> OWP  $\mathcal{F}$  with its specification  $\mathcal{F}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}^{\mathcal{F}}, \text{Eval}_{\text{SPEC}}^{\mathcal{F}})$ , our full-fledged construction with its specification  $\mathcal{H}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}^{\mathcal{H}}, \text{PRG}_{\text{SPEC}}^{\mathcal{H}})$ <sup>10</sup> is as follows:

- Parameter generation algorithm  $pk \leftarrow \text{KG}_{\text{SPEC}}^{\mathcal{H}}(\lambda)$ :  
 compute  $i \leftarrow \text{KG}_{\text{SPEC}}^{\mathcal{F}}(\lambda)$  and set  $pk := i$ ;
- Bit string generation algorithm  $(s', r) \leftarrow \text{PRG}_{\text{SPEC}}^{\mathcal{H}}(pk, s)$ :  
 upon receiving  $s$  and  $pk$  where  $pk = i$ ,  $s = s_1 || s_2$ , and  $|s_1| = |s_2| = \ell$ , compute the following:
  - let  $s_1^0 := s_1$  and  $s_2^0 := s_2$ ;
  - for  $j = 1, \dots, \ell'$ ,
    - $b_j := \langle s_1^{j-1}, s_2^{j-1} \rangle$ ;
    - $s_1^j := s_1^{j-1}$ ;  $s_2^j := \text{Eval}_{\text{SPEC}}^{\mathcal{F}}(i, s_2^{j-1})$ ;  $s^j := s_1^j || s_2^j$ ;
  - $s' = s^{\ell'} = s_1 || s_2^{\ell'}$ ; and  $r = b_1 \dots b_{\ell'}$ .

Please see also Figure 10 for pictorial illustration for implementations.

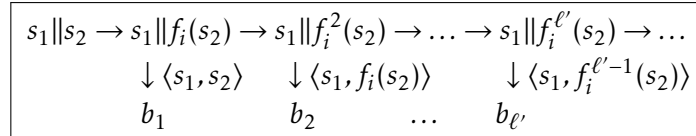


Figure 10: One iteration of BM-PRG, and  $f_i(x) := \text{Eval}_{\text{IMPL}}^{\mathcal{F}}(i, x)$

<sup>10</sup> $\text{PRG}_q$  can be defined in a straightforward manner that runs the above PRG for  $q$  iterations, each iteration outputs  $\ell'$  bits and updates the state for next iteration.

**Theorem B.1** *The full fledged construction above with specification  $\mathcal{H}_{\text{SPEC}}$  is  $q$ -subversion-resistant<sup>A</sup> (for any polynomially large  $q$ ), if  $\mathcal{F}$  with specification  $\mathcal{F}_{\text{SPEC}}$  is a subversion-resistant<sup>A</sup> OWP.*

**Proof 2 (Proof sketch)** *Following Lemma 5.2,  $s_1 \| f_i(s_2) \| b_1$  is pseudorandom, i.e., it is indistinguishable from “ $u_1, \dots, u_{2\ell}, v_1$ ”, even to an adversary  $\mathcal{A}$  who may set the public parameter  $i$ , where  $\{u_t\}_{t \in [\ell]}, v_1$  are all random bits, and  $f_i(s_2) = \text{Eval}_{\text{IMPL}}^{\mathcal{F}}(i, s_2)$ .*

*Observe that  $b_2$  can be computed from  $s_1, f_i(s_2)$ ; it follows that the adversary  $\mathcal{A}$  (who has the backdoor) can not predict  $b_2$  from  $b_1$ , otherwise she trivially distinguishes  $s_1 \| f_i^1(s_2) \| b_1$  from random, simply by computing  $b_2$  from  $s_1 \| f_i^1(s_2)$  and predicting using  $b_1$  to see whether these are consistent. Similarly,  $\mathcal{A}$  cannot predict  $b_3$  from  $b_1, b_2$ : To see this, first observe that  $\mathcal{A}$  can not predict  $b_3$  from  $b_2$ , thus  $\mathcal{A}$  can not predict  $b_3$  from  $v_1, b_2$  where  $v_1$  is a random bit. If  $b_3$  is predictable by  $\mathcal{A}$  from  $b_1, b_2$ , then starting from  $s_1 \| f_i(s_2)$ ,  $\mathcal{A}$  computes  $b_2, b_3$ , and simply uses  $b_1$  to test whether the predication is correct, so that she can distinguish  $s_1 \| f_i(s_2) \| b_1$  from  $s_1 \| f_i(s_2) \| v_1$ , and further from  $u_1 \dots, u_{2\ell}, v_1$ .*

*The above argument can be applied to any  $j \in [\ell']$ , so that  $\mathcal{A}$  can not predict  $b_{j+1}$  from  $b_1, \dots, b_j$ . Then—following the classic reduction from pseudorandomness to next-bit unpredictability—we can conclude that  $b_1 \dots b_{\ell'}$  is indistinguishable from uniform bits  $\{0, 1\}^{\ell'}$ , even to  $\mathcal{A}$ . (This can be shown via the standard hybrid argument.) Last, inductively, we can conclude that  $r_1, \dots, r_q$  are indistinguishable from  $\ell' \cdot q$  uniform bits.*

## C Supporting Material: Omitted Proofs

### C.1 Proof of Lemma 3.2

**Proof 3** *Consider a TDOWP  $\mathcal{F} = \{f_i\}$  with the associated specifications  $\mathcal{F}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}^{\mathcal{F}}, \text{Eval}_{\text{SPEC}}^{\mathcal{F}}, \text{Inv}_{\text{SPEC}}^{\mathcal{F}})$ . Assuming the trapdoors can be represented using  $\ell(\lambda)$  bits, we construct a subvertible OWP family  $\mathcal{G}$  with specification  $\mathcal{G}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}^{\mathcal{G}}, \text{Eval}_{\text{SPEC}}^{\mathcal{G}}, \text{Inv}_{\text{SPEC}}^{\mathcal{G}})$  as follows:*

- *Function generation  $(i, r) \leftarrow \text{KG}_{\text{SPEC}}^{\mathcal{G}}$ , where  $\text{KG}_{\text{SPEC}}^{\mathcal{G}}$  is given by:  
Run the  $\text{KG}_{\text{SPEC}}^{\mathcal{F}}$  algorithm and receive a function index/trapdoor pair  $(i, t_i)$ ; then discard  $t_i$ , and sample randomly  $r \leftarrow \{0, 1\}^{\ell(\lambda)}$ ; It outputs  $(i, r)$ .*
- *Function evaluation  $y \leftarrow \text{Eval}_{\text{SPEC}}^{\mathcal{G}}(i, r, x)$ , where  $\text{Eval}_{\text{SPEC}}^{\mathcal{G}}$  is given by:  
Upon inputting  $i, r, x$ , discard  $r$ , compute  $y \leftarrow \text{Eval}_{\text{SPEC}}^{\mathcal{F}}(i, x)$ ; it outputs  $y$ .*
- *Invert function  $\text{Inv}_{\text{SPEC}}^{\mathcal{G}}$  is the same as  $\text{Inv}_{\text{SPEC}}^{\mathcal{F}}$ .*

*It is easy to see that  $\mathcal{G}$  is one way because  $\mathcal{F}$  is one way (without  $t_i$ ) in the classical setting.*

*While in the cliptographic setting, an adversary can first chooses a random key  $k$  for a symmetric key encryption scheme  $\text{SE} = (\text{SE.Enc}, \text{SE.Dec})$ . Then she provides a implementations below that can evade the detection and breaks the one-way security of the functions generated by the implementation:*

- *Function generation implementation,  $(i, \tilde{r}) \leftarrow \text{KG}_{\text{IMPL}}^{\mathcal{G}}(k)$ , where  $\text{KG}_{\text{IMPL}}^{\mathcal{G}}(k)$  is given by:  
It first runs  $\text{KG}_{\text{SPEC}}^{\mathcal{F}}$ , and receives an index  $i$  together with the corresponding trapdoor  $t_i$ ;  $\tilde{r} = \text{SE.Enc}(k, t_i)$ , i.e.,  $\tilde{r}$  is generated by encrypting  $t_i$  using  $k$ . It outputs  $(i, \tilde{r})$ .*
- *$\text{Eval}_{\text{IMPL}}^{\mathcal{F}}$  and  $\text{Inv}_{\text{IMPL}}^{\mathcal{F}}$  are the same as the specifications.*

It is easy to see that when adversary obtains  $i, \tilde{r}$ , and  $y = \text{Eval}_{\text{SPEC}}^{\mathcal{F}}(i, x)$  she can simply decrypts  $\tilde{r}$  to get  $t_i = \text{SE.Dec}(k, \tilde{r})$ , and then inverts  $y$  to get  $x$ .

Furthermore, since  $\text{SE.Enc}$  is modeled as a which is assumed to be a pseudorandom permutation (PRP). PRP, the distributions of  $(i, r)$  returned by  $\text{KG}_{\text{SPEC}}^{\mathcal{G}}$  and  $(i, \tilde{r})$  returned by  $\text{KG}_{\text{IMPL}}^{\mathcal{G}}(k)$  for any  $k$  are computationally indistinguishable.

## C.2 Proof of Theorem 3.6

**Proof 4** Suppose there is a PPT adversary  $\mathcal{A}_{\mathcal{G}}$  that subverts  $\mathcal{G}$  in the game defined in Fig. 9. We will build an adversary  $\mathcal{A}_{\mathcal{F}}$  that breaks the one-way security of  $\mathcal{F}_{\text{SPEC}}$ .

Construction of  $\mathcal{A}_{\mathcal{F}}$ . Assume  $(i^*, y^*)$  are the challenges sent from the  $\mathcal{F}_{\text{SPEC}}$  challenger  $\mathcal{C}_{\mathcal{F}}$  (suppose  $x^*$  is the input chosen by  $\mathcal{C}_{\mathcal{F}}$  that generates this challenge).  $\mathcal{A}_{\mathcal{F}}$  will simulate the game with  $\mathcal{A}_{\mathcal{G}}$ . When  $\mathcal{A}_{\mathcal{G}}$  asks random oracle queries  $r_1, \dots, r_q$ ,  $\mathcal{A}_{\mathcal{F}}$  answers all those queries with uniform strings.

$\mathcal{A}_{\mathcal{G}}$  then provides the implementations  $(\text{RG}_{\text{IMPL}}^{\mathcal{G}}, \text{dKG}_{\text{IMPL}}^{\mathcal{G}}, \text{Eval}_{\text{IMPL}}^{\mathcal{G}}, \text{Inv}_{\text{IMPL}}^{\mathcal{G}})$ .  $\mathcal{A}_{\mathcal{F}}$  continues the simulation, querying  $\text{RG}_{\text{IMPL}}^{\mathcal{G}}$  to receive  $r$ . If  $r \in \{r_1, \dots, r_q\}$ ,  $\mathcal{A}_{\mathcal{F}}$  aborts; otherwise  $\mathcal{A}_{\mathcal{F}}$  sends  $(i^*, y^*)$  to  $\mathcal{A}_{\mathcal{G}}$ . Finally,  $\mathcal{A}_{\mathcal{F}}$  submits the answer  $x'$  from  $\mathcal{A}_{\mathcal{G}}$  as his answer to  $\mathcal{A}_{\mathcal{F}}$ .

Probabilistic analysis. Now we bound the success probability of  $\mathcal{A}_{\mathcal{F}}$ . Suppose  $x$  is the random input chosen by  $\mathcal{C}_{\mathcal{F}}$ , and let  $W$  denote the event that  $\mathcal{A}_{\mathcal{F}}$  aborts. It follows that:  $\Pr[x^* = x'] = \Pr[x^* = x' | \overline{W}] \Pr[\overline{W}]$ .

From the assumption that  $\mathcal{A}_{\mathcal{G}}$  breaks the security of  $\mathcal{G}$ , we can infer that the following two conditions: (1) the detection probability **Det** is negligible; (2) the advantage **Adv** is non-negligible  $\delta$ . From condition (1), we have

$$\Pr[r \notin \{r_1, \dots, r_q\} : r \leftarrow \text{RG}_{\text{IMPL}}^{\mathcal{G}}] \geq 1 - \text{negl}(\lambda).$$

Otherwise, there is a watchdog algorithm that simply sample two outputs to detect whether there is collision, in which case the implementation  $\text{RG}_{\text{IMPL}}^{\mathcal{G}}$  is rejected. On the other hand, if  $\text{RG}_{\text{SPEC}}^{\mathcal{G}}$  outputs uniform bits, the collision probability is negligible; note that here  $\text{RG}_{\text{SPEC}}^{\mathcal{G}} = \text{RG}_{\text{SPEC}}^{\mathcal{F}}$ . Thus  $\Pr[\overline{W}] \geq 1 - \text{negl}(\lambda)$ .

Next, we bound  $\Pr[x' = x^* | \overline{W}]$ . From condition (1) again, following the proof of Lemma 2.3, we claim:

$$\Pr[\text{dKG}_{\text{IMPL}}^{\mathcal{G}}(r) = \text{dKG}_{\text{SPEC}}^{\mathcal{G}}(r) : r \leftarrow \text{RG}_{\text{IMPL}}^{\mathcal{F}}] \geq 1 - \text{negl}(\lambda); \quad (1)$$

$$\Pr[\text{Eval}_{\text{IMPL}}^{\mathcal{G}}(i^*, x) = \text{Eval}_{\text{SPEC}}^{\mathcal{F}}(i^*, x)] \geq 1 - \text{negl}(\lambda); \quad (2)$$

otherwise there is a trivial watchdog that samples, tests equality, and rejects misbehaving implementations.

Now let us analyze the view of the adversary  $\mathcal{A}_{\mathcal{G}}$  in the subversion-resistant<sup>SP</sup> TDOWP game. From Ineq. (1), with overwhelming probability  $\mathcal{A}_{\mathcal{G}}$  is supposed to see an index  $i$  in  $(i, t) \leftarrow \text{dKG}_{\text{IMPL}}^{\mathcal{G}}(r) = \text{dKG}_{\text{SPEC}}^{\mathcal{G}}(r)$ , where  $r \leftarrow \text{RG}_{\text{IMPL}}^{\mathcal{G}}$ ; similarly,  $y \leftarrow \text{Eval}_{\text{IMPL}}^{\mathcal{G}}(i, x)$  for a randomly selected  $x$ . While  $i^*$  is generated by calling  $\text{dKG}_{\text{SPEC}}^{\mathcal{F}}(r^*)$  for an uniform  $r^*$ . Conditioned on  $\overline{W}$ , the distribution of  $\text{dKG}_{\text{SPEC}}^{\mathcal{G}}(r) = \text{dKG}_{\text{SPEC}}^{\mathcal{F}}(h_{\text{SPEC}}(r))$  is identical to  $\text{dKG}_{\text{SPEC}}^{\mathcal{F}}(r^*)$ . Also, from Ineq. (2), with overwhelming probability,  $y^* = \text{Eval}_{\text{SPEC}}^{\mathcal{F}}(i^*, x^*)$  would be consistent with  $\text{Eval}_{\text{IMPL}}^{\mathcal{G}}(i^*, x^*)$ . Thus we can claim that the distribution of  $i^*, y^*$  is indistinguishable from what  $\mathcal{A}_{\mathcal{G}}$  expects. Together with condition (2),  $\mathcal{A}_{\mathcal{G}}$  will return a correct  $x^*$  with a noticeable probability  $\delta$ . Combined all above, we can conclude that

$$\Pr[x = x^*] \geq \delta(1 - \text{negl}(\lambda))(1 - \text{negl}(\lambda)) \geq \delta - \text{negl}(\lambda).$$

Thus  $\mathcal{A}_{\mathcal{F}}$  breaks the security of  $\mathcal{F}_{\text{SPEC}}$ , a contradiction.

### C.3 Proof of Theorem 4.3

**Proof 5** Assume that  $\mathcal{SS}$  is not subversion-resistant<sup>C</sup>, then for any watchdog  $\mathcal{W}_{\mathcal{SS}}$ , there exists a PPT adversary  $\mathcal{A}_{\mathcal{SS}}$ , such that: (1) the detection probability  $\text{Det}_{\mathcal{W}_{\mathcal{SS}}, \mathcal{A}_{\mathcal{SS}}}$  is negligible; (2) the advantage  $\text{Adv}_{\mathcal{A}_{\mathcal{SS}}}$  is non-negligible. This means the implementations provided by  $\mathcal{A}_{\mathcal{SS}}$  will be accepted by  $\mathcal{W}_{\mathcal{SS}}$  who sees the transcripts generated in the game; in addition  $\mathcal{A}_{\mathcal{SS}}$  outputs a valid new forgery with non-negligible probability, say  $\delta$ . We will construct an adversary  $\mathcal{A}_{\mathcal{F}}$  that can break the security of  $\mathcal{F}_{\text{SPEC}}$ .

*Construction of  $\mathcal{A}_{\mathcal{F}}$ .* It is easy to see that for any efficient watchdog  $\mathcal{W}_{\mathcal{F}}$  for a TDOWP<sup>C</sup>, we can define a corresponding watchdog  $\mathcal{W}_{\mathcal{SS}}$  for the signature scheme. There will be a PPT adversary of the signature scheme  $\mathcal{A}_{\mathcal{SS}}$  that satisfies the above two conditions, and we will build another adversary  $\mathcal{A}_{\mathcal{F}}$ . In particular, we define  $\mathcal{W}_{\mathcal{SS}}$  as follows: it first runs  $\mathcal{W}_{\mathcal{F}}$ ; it runs  $\text{KG}_{\text{IMPL}}$  twice to see whether there is collision; it also runs the  $\text{Verify}_{\text{IMPL}}^{\text{SS}}$  algorithm on the transcripts and compares with the results of running the specification.

Suppose  $(*, m_{1,1}), \dots, (*, m_{q_1,1})$  are the random oracle queries that  $\mathcal{A}_{\mathcal{SS}}$  makes before outputting the implementations,  $\mathcal{A}_{\mathcal{F}}$  answers all those using uniform strings.

Upon receiving the implementations of  $(\text{KG}_{\text{IMPL}}^{\text{SS}}, \text{Sign}_{\text{IMPL}}^{\text{SS}}, \text{Verify}_{\text{IMPL}}^{\text{SS}})$  from  $\mathcal{A}_{\mathcal{SS}}$ ,  $\mathcal{A}_{\mathcal{F}}$  prepares implementations of  $\text{Eval}_{\text{IMPL}}^{\mathcal{F}}, \text{Inv}_{\text{IMPL}}^{\mathcal{F}}$  using  $\text{Verify}_{\text{IMPL}}^{\text{SS}}, \text{Sign}_{\text{IMPL}}^{\text{SS}}$  and sends them together with  $\text{KG}_{\text{IMPL}}^{\text{SS}}$  to the  $\mathcal{F}_{\text{SPEC}}$  watchdog  $\mathcal{W}_{\mathcal{F}}$ . It is easy to see that  $\mathcal{W}_{\mathcal{F}}$  will not complain, otherwise,  $\mathcal{W}_{\mathcal{SS}}$  will complain and thus it violates condition (1).

Upon receiving  $i^*, y^*$  from the  $\mathcal{F}_{\text{SPEC}}$  challenger  $\mathcal{C}_{\mathcal{F}}$ ,  $\mathcal{A}_{\mathcal{F}}$  directly forwards  $i^*$  as the public key  $\text{pk}$  to  $\mathcal{A}_{\mathcal{SS}}$ . Note that  $i^* = \text{pk}$  comes exactly from  $\text{KG}_{\text{IMPL}}^{\mathcal{F}} = \text{KG}_{\text{IMPL}}^{\text{SS}}$ , and  $y^* = \text{Eval}_{\text{IMPL}}^{\mathcal{F}}(x^*)$  (essentially  $\text{Verify}_{\text{IMPL}}^{\text{SS}}$ ) on a uniformly chosen  $x^*$ .

Now if  $\mathcal{A}_{\mathcal{SS}}$  makes another set of random oracle queries for  $(\text{pk}, m_{1,2}), \dots, (\text{pk}, m_{q_2,2})$ ,  $\mathcal{A}_{\mathcal{F}}$  chooses a random index  $t \in [q_2]$  and answers  $y^*$  as  $h_{\text{SPEC}}(\text{pk}, m_{t,2})$ , and randomly chooses  $\{\sigma_j\}_{j \neq t}$  and returns  $\{\text{Eval}_{\text{IMPL}}^{\mathcal{F}}(\text{pk}, \sigma_j)\}_{j \neq t}$  as answers for all other queries  $\{(\text{pk}, m_{j,2})\}_{j \neq t}$ .  $\mathcal{A}_{\mathcal{F}}$  keeps a list for the values of all those queries. (We remark that without loss of generality here we assume all the prefixes are  $\text{pk}$ , if not,  $\mathcal{A}_{\mathcal{F}}$  simply returns a random string as an answer.)

If  $\mathcal{A}_{\mathcal{SS}}$  then makes signing queries on  $m_{1,3}, \dots, m_{q_3,3}$ ,  $\mathcal{A}_{\mathcal{F}}$  can simply find the corresponding  $\{\sigma_{i,3}\}$  from the list and return them as the signatures. If  $\mathcal{A}_{\mathcal{SS}}$  outputs a forgery  $m^*, \sigma^*$  and  $m_{t,2} = m^*$ ,  $\mathcal{A}_{\mathcal{F}}$  returns  $\sigma^*$  to  $\mathcal{C}_{\mathcal{F}}$  as the pre-image for  $y^*$ .

*Probabilistic analysis.* Now we bound the success probability of  $\mathcal{A}_{\mathcal{F}}$ . First, if  $\text{pk}, m^*$  appeared in the random oracle queries before outputting the implementations, there would be inconsistency of answers to the random oracle queries. However, this can only happen with a negligible probability. To see this, the watchdog will sample  $\text{KG}_{\text{IMPL}}$  (which is resettable) twice to observe collision, if  $\mathcal{A}_{\mathcal{SS}}$  can predict the output of  $\text{KG}_{\text{IMPL}}$  with a non-negligible advantage, then the collision probability will also be non-negligible. Thus it violates condition (1), it has to be the case that  $\mathcal{A}_{\mathcal{SS}}$  can only predict  $\text{pk}$  with a negligible probability.

Conditioned on that the above event never happens. Let  $\delta' = \Pr[\sigma^* = x^*]$ , where  $x^*$  is the input chosen by  $\mathcal{C}_{\mathcal{F}}$ . It is easy to see that

$$\begin{aligned} \delta' &= \Pr[\sigma^* = x^* | m_{t,2} = m^*] \Pr[m_{t,2} = m^*] + \Pr[\sigma^* = x^* | m_{t,2} \neq m^*] \Pr[m_{t,2} \neq m^*] \\ &\geq \Pr[\sigma^* = x^* | m_{t,2} = m^*] \Pr[m_{t,2} = m^*]. \end{aligned}$$

Since  $\mathcal{W}_{\mathcal{SS}}$  checks  $(m^*, \sigma^*)$  using  $\text{Verify}_{\text{SPEC}}^{\text{SS}}$ , following condition (1), we can claim that  $\mathcal{A}_{\mathcal{SS}}$  made a random oracle query for  $\text{pk} || m^*$ ; if not, the probability that  $\text{Eval}_{\text{SPEC}}^{\mathcal{F}}(\text{pk}, \sigma^*)$  hits a completely random value is negligible.



With the above claim,  $(pk, m^*)$  will appear in the adversaries' random oracle queries. Observe that  $pk$  is output by  $\text{KG}_{\text{IMPL}}^{\text{SS}}$ , and we can see that with overwhelming probability,  $pk||m^*$  will appear in the random oracle queries after  $pk$  is provided. It follows that:  $\Pr[m_{t,2} = m^*] \geq \frac{1}{q_2}(1 - \text{negl}(\lambda))$ .

Conditioned on  $m_{t,2} = m^*$ , from  $\mathcal{A}_{\text{SS}}$ 's view, in the subversion-resistant<sup>C</sup> game, she is supposed to receive  $\sigma_1, \dots, \sigma_q$  that are output by  $\text{Sign}_{\text{IMPL}}^{\text{SS}}$  on inputting the corresponding  $sk$  and messages. With the online watchdog, and the property of the full domain hash, that those signatures will pass  $\text{Verify}_{\text{SPEC}}^{\text{SS}}$ , i.e.,  $\text{Eval}_{\text{SPEC}}^{\mathcal{F}}(\sigma_i) = h_{\text{SPEC}}(m_i)$ . This implies that the distribution of  $\sigma_i$  is the same from that comes from  $\text{Sign}_{\text{SPEC}}^{\text{SS}}$  by calling  $\text{Inv}_{\text{SPEC}}^{\mathcal{F}}(h_{\text{SPEC}}(m_i))$ , which is identically distributed as how  $\mathcal{A}_{\mathcal{F}}$  simulates the answers for the signing queries. Now together with condition (2),  $\mathcal{A}_{\text{SS}}$  will output a valid forgery w.r.t  $\text{Verify}_{\text{IMPL}}^{\text{SS}}$ . With the offline watchdog,  $\text{Verify}_{\text{IMPL}}^{\text{SS}}(m^*, \sigma^*) = \text{Verify}_{\text{SPEC}}^{\text{SS}}$ , i.e.,  $\text{Eval}_{\text{SPEC}}^{\mathcal{F}}(\sigma^*) = h_{\text{SPEC}}(m^*) = y$ , thus  $\Pr[\sigma^* = x^* | m^* = m_{t,2}] \geq \delta(1 - \text{negl}(\lambda))$ .

Combing all above, we see that  $\delta' \geq \delta/q_2(1 - \text{negl}(\lambda))$  which is non-negligible. We can see that  $\mathcal{A}_{\mathcal{F}}$  breaks the security of  $\mathcal{F}_{\text{SPEC}}$ , which completes the proof.

#### C.4 Proof of Lemma 5.2

**Proof 6** The specification  $\text{KG}_{\text{SPEC}}^{\mathcal{G}}$  of the simplified Blum-Micali PRG outputs a random function index from the index set (by simply running  $\text{KG}_{\text{SPEC}}^{\mathcal{F}}$ ).

It is easy to see that the OWP function family  $\overline{\mathcal{F}}$  given by  $\text{Eval}_{\text{SPEC}}^{\overline{\mathcal{F}}}(i, x_1 || x_2) := x_1 || \text{Eval}_{\text{SPEC}}^{\mathcal{F}}(x_2)$ , is subversion-resistant<sup>A</sup> if  $\mathcal{F}$  is subversion-resistant<sup>A</sup>. If the above basic construction is not subversion-resistant<sup>A</sup>, then there exists a PPT adversary  $\mathcal{A}$  such that (i.) all watchdogs will accept the public parameter  $pk$  and the implementation  $\text{PRG}_{\text{IMPL}}^{\mathcal{G}}$ , and (ii.)  $\mathcal{A}$  distinguishes the PRG output from a random  $2\ell + 1$ -bit string with some non-negligible probability  $\delta$ . To expand condition (ii.),  $\mathcal{A}$  can distinguish  $s_1 || f_i(s_2) || B(s)$  (where  $B(s) = \langle s_1, s_2 \rangle$ ) from a uniform  $(2\ell + 1)$ -bit string for an uniform  $s$ .

Now from  $\mathcal{A}$ , we construct an adversary  $\mathcal{A}_{\mathcal{F}}$  that breaks the subversion-resistance of  $\mathcal{F}_{\text{SPEC}}$ , as an OWP<sup>A</sup>:  $\mathcal{A}_{\mathcal{F}}$  runs  $\mathcal{A}$  to get  $pk, \text{PRG}_{\text{IMPL}}^{\mathcal{G}}$ .  $\mathcal{A}_{\mathcal{F}}$  then implements  $\text{Eval}_{\text{IMPL}}^{\mathcal{F}}$  as follows: when evaluating on inputs  $pk, x$ ,  $\text{Eval}_{\text{IMPL}}^{\mathcal{F}}$  calls  $\text{PRG}_{\text{IMPL}}^{\mathcal{G}}$  using  $pk, \bar{0} || x$  and receives  $\bar{0} || y || 0$ ,  $\text{Eval}_{\text{IMPL}}^{\mathcal{G}}$  then discards the first  $\ell$  bits and the last bit, and returns  $y$ . We can see that  $\text{Eval}_{\text{IMPL}}^{\mathcal{F}}$  and  $\text{PRG}_{\text{IMPL}}^{\mathcal{G}}$  have identical input/output behavior; thus for any watchdog, if it accepts  $pk, \text{PRG}_{\text{IMPL}}^{\mathcal{G}}$ , it also accepts  $pk, \text{Eval}_{\text{IMPL}}^{\mathcal{F}}$ .  $\mathcal{A}_{\mathcal{F}}$  then continues the simulation: when receiving a challenge  $y$ ,  $\mathcal{A}_{\mathcal{F}}$  randomly samples  $r$ , and a bit  $b$  and sends  $r || y || b$  for  $\mathcal{A}$  to distinguish. If  $b = B(s)$ ,  $\mathcal{A}$  will output 1 with probability  $1/2 + \delta$ . It is easy to see that such adversary can predict the GL hardcore predicate  $B$  with advantage  $\delta/2$ . Following the GL proof [GL89], there exists another algorithm  $I^{\mathcal{A}}$  that can invert  $y$  with probability  $\delta' = \text{poly}(\delta/2)$ .

$\mathcal{A}_{\mathcal{F}}$  runs  $I^{\mathcal{A}}$  to invert  $y$ , and will output a correct pre-image with probability  $\delta'$ . Combing both claims above,  $\mathcal{A}_{\mathcal{F}}$  subverts  $\mathcal{F}_{\text{SPEC}}$  with a non-negligible probability. This leads to a contradiction.