

# Generic Semantic Security against a Kleptographic Adversary

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

September 16, 2018

## Abstract

Notable recent security incidents have generated intense interest in adversaries which attempt to subvert—perhaps covertly—cryptographic algorithms. In this paper we develop (IND-CPA) Semantically Secure encryption in this challenging setting.

This fundamental encryption primitive has been previously studied in the “kleptographic setting,” though existing results must relax the model by introducing trusted components or otherwise constraining the subversion power of the adversary: designing a Public Key System that is kletographically semantically secure (with minimal trust) has remained elusive to date.

In this work, we finally achieve such systems, even when all relevant cryptographic algorithms are subject to adversarial (kleptographic) subversion. To this end we exploit novel inter-component randomized cryptographic checking techniques (with an offline checking component), combined with common and simple software engineering modular programming techniques (applied to the system’s black box specification level). Moreover, our methodology yields a strong generic technique for the preservation of any semantically secure cryptosystem when incorporated into the strong kleptographic adversary setting.

## 1 Introduction

Modern cryptography has been spectacularly successful, leading to cryptographic tools with remarkable functionalities and security guarantees. Despite these advances, applying cryptographic tools to provide robust security *in practice* is a notorious challenge. In particular, practical settings often introduce threats that are not adequately reflected by conventional cryptographic security modeling. In this article, we focus on one such disparity between conventional modeling and practical application: the possibility of adversarial instantiation, subversion, or substitution of the cryptographic algorithms themselves.

One implicit assumption in typical cryptographic security modeling is that the deployed implementations of cryptographic algorithms actually realize their “official specifications.” In practice, cryptographic implementations may diverge from their specifications for a variety of reasons, including unintentional programming bugs or malicious tampering; in the *kleptographic* setting, one considers the pathological possibility of *fully adversarial implementations* of cryptographic algorithms. The goal of such an adversary is to produce implementations of cryptographic algorithms

---

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

<sup>†</sup>New Jersey Institute of Technology, qiang@njit.edu

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

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

which compromise security while appearing to be correct even in the face of fairly intensive testing. (Formal models are discussed below.) We remark that the possibility of such kleptographic attacks arise whenever a “third-party” software library or hardware device is relied upon for cryptographic purposes.

The consequences of such attacks are rather surprising: It turns out that—in wide generality—*adversarial implementations of randomized algorithms may leak private information while producing output that is indistinguishable from the specification*. The possibility of such threats was showcased over two decades ago by Young and Yung [YY96, YY97]. Recently, starting with Bellare, Paterson, and Rogaway [BPR14], the topic has received renewed formal attention [DGG<sup>+</sup>15, BH15, MS15, DFP15, BJK15, AMV15, RTYZ15, DMSD15] motivated by startling evidence from the Snowden revelations of past deployment of kleptographic attacks. The topic was also notably highlighted by Rogaway’s 2015 IACR Distinguished Lecture [Rog15] calling for community-wide efforts to explore defending mechanisms. One of the most striking recent discoveries [BPR14, BJK15] establishes that a steganographic channel can be embedded in the output of a subverted randomized (encryption) algorithm so that secret information can be exclusively leaked to the adversary. Such steganographic attacks can even be applied in settings where the subverted algorithms are stateless [BJK15, RTYZ15].

As mentioned above, the kleptographic setting features an adversary who may substitute malicious algorithms in place of intended cryptographic algorithms. The adversary’s attack, however, is constrained by introducing a trusted party (we call it the “watchdog”) who may *test* the (adversarially-provided) implementations against the specification—this testing can provide a measure of safety to the final users of the implementation. The major question is whether a combination of careful specification and testing can preserve security despite such a powerful adversary. The broad challenge is to rework the classical framework of cryptographic primitives and constructions so that they provide security in this new environment. Recent efforts have partially explored the landscape and clearly identified the critical role played by randomized algorithms in the kleptographic setting. Indeed, existing defending strategies roughly fall into two categories:

- “Abandon” randomized algorithms and turn to deterministic counterparts, e.g., use deterministic encryption with a unique ciphertext property as suggested in [BPR14, BH15, DFP15, BJK15] for encryption schemes. Of course, key generation is a sticking point in this approach, as it must be randomized; this has been handled by simply placing key generation outside the subversion model (i.e., assuming honest key generation).
- Use a trusted reverse firewall to re-randomize all incoming/outgoing communication generated by the randomized algorithms [MS15, DMSD15].

While insisting on deterministic algorithms has favorable properties, it necessarily places many central notions of security entirely out of reach: in particular, IND- security for public-key encryptions is unattainable. Furthermore, as mentioned above, the process of key generation is inherently stochastic; this difficulty has been avoided by placing key generation outside the subversion model (or, equivalently, assuming it to be trusted). On the other hand, the firewall model can provide general feasibility results but requires the assumption of an active trusted party (in particular, it requires a source of trusted randomness). As the large-scale goal of the study of kleptography is to reduce the need for trust in the underlying components, it is attractive to understand to what extent we can eliminate such trusted randomness, algorithmic elements, etc.

In this paper, we address the following central questions:

1. *Is it possible to generically annihilate subliminal channels in subverted randomized algorithms without a trusted party?*
2. *Is it possible to achieve IND-CPA secure public-key encryption without trust in any of the component algorithms.*

We provide affirmative answers to both questions.

Our principal technique (completely realizable in systems) involves a certain decomposition of algorithms into a few functional “components,” which are tested by the watchdog and may be independently run and “re-assembled” by the user.

**The kleptographic model, in brief.** The kleptographic model is meant to capture a situation where an adversary (or “big brother” as we shall occasionally say) has the opportunity to implement (and, indeed, “mis-implement” or subvert) our basic cryptographic tools. On the other hand, the model also introduces a “watchdog” who will attempt to check, via black-box testing, that the cryptographic tools have been faithfully implemented by the adversary. We imagine the adversary to be “proud but malicious”: the adversary wishes to interfere with security, but does not wish to be exposed as a fraud by the watchdog. The model, in brief:

1. **Specification (see also ★ below.)** The cryptographic primitive is *specified* as a tuple  $\Pi_{\text{SPEC}} = (F_{\text{SPEC}}^1, \dots, F_{\text{SPEC}}^k)$  of “functionalities.” Each  $F_{\text{SPEC}}^i$  is either a function  $F : \{0, 1\}^* \rightarrow \{0, 1\}^*$  (which specifies a deterministic algorithm) or a family of probability distributions, so that  $F(x)$  is a probability distribution over  $\{0, 1\}^*$  for each  $x$  (which specifies a randomized algorithm).
  2. **Subversion.** The adversary provides us with all algorithmic and cryptographic building blocks; that is, the adversary provides us with an “implementation”  $(F_{\text{IMPL}}^1, \dots, F_{\text{IMPL}}^k)$  for each of the functionalities  $F_{\text{IMPL}}^i$ . Observe that the adversary may provide even the algorithms that generate randomness and random objects such as keys. Of course, in general the implementations may disagree with the specification, which can provide the adversary a novel avenue to attack the primitive.
  3. **Testing; the watchdog.** The algorithmic and cryptographic building blocks are then sent to trusted testing facility, the *watchdog*. The watchdog is aware of the official specification, and may query the adversary’s implementations (treating them as block-boxes) in an attempt to detect disagreements between the implementations and the specifications.
  4. **The security game.** Assuming that the watchdog is satisfied, the implementations are pressed into service, at which point their security is modeled by a conventional security game.
- ★ **Remark: Decomposition and Amalgamation.** We permit the designer of the cryptographic primitive (who determines its specification) an extra dimension of freedom which can assist the watchdog in his verification task: We permit the designer to functionally *decompose* the primitives into a fixed number of “pieces.” For example, rather than specifying a function of interest  $f : X \rightarrow Y$ , the designer may instead specify two functions  $h : X \rightarrow W$  and  $g : W \rightarrow Y$  with the property that  $f = g \circ h$ . (Thus  $h$  and  $g$  together implicitly specify  $f$ .) An important example in our setting is specifying a randomized algorithm  $G(x)$  as a composition  $dG(x, \text{RG}(1^k))$ , where  $dG$  is a deterministic algorithm and  $\text{RG}$  is an algorithm which, given the input  $1^k$ , produces  $k$  uniformly random bits as output. In general, the decomposition may be arbitrary, but may only involve  $O(1)$  pieces and cannot depend on input length.

**Our contributions.** In this paper, we develop general techniques that eliminate subliminal channels introduced by adversarial implementations. We apply this general technique to construct symmetric-key and public-key encryption schemes preserving the IND-CPA security when all the algorithms are subject to adversarial subversion. In particular,

1. We first define a property of “stego-freeness” by adjusting the previous models of [BPR14, RTYZ15] to characterize whether an implementation of an algorithm can be considered to be following its specification in the kleptographic setting (where there is a watchdog). The model includes several variants depending on the algorithm inputs and choices about the decomposition. We also define a notion of “subversion resistance” for encryption schemes by considering the kleptographic IND-CPA security game where the challengers use adversarial implementations and the implementations are interrogated by an “offline” watchdog. (See below for formal definitions.)
2. We then consider how to defend against steganographic channel attacks by the simple non-black-box technique of decomposition-and-amalgamation. We first extend the attacks of [BPR14, BJK15, DGG<sup>+</sup>15] to the setting where a public-key encryption algorithm is split into two components: randomness generation (even when the immunization function is modeled as a random oracle) and a deterministic component. We then demonstrate a “double-splitting” strategy in which randomness generation is carried out by two independent components  $RG_0, RG_1$ . We prove that when  $r_0, r_1$  are sampled independently from  $RG_0$  and  $RG_1$ , mixing them with an immunization function  $\Phi$  can indeed destroy subliminal channels in the implementation of a wide class of randomized algorithms in the random oracle model.<sup>1</sup> We also consider how to achieve such results in the standard model (without a random oracle); see Sec. D in the appendix.
3. We further apply this general technique to immunize each algorithm of a symmetric-key (single-bit) encryption scheme, including key generation and encryption. Our construction preserves IND-CPA security of the underlying primitive. We then focus on constructing symmetric-key encryption schemes for large message spaces. To defend against input-trigger-style attacks [DFP15], we allow the user to carry out one single trusted addition. We also consider correctness in the kleptographic setting and draw connections to the theory of self-correcting programs. These techniques can be applied directly to immunize public-key encryption, and our construction gives the first IND-CPA secure scheme in the kleptographic setting without relying on a trusted party. Finally, we discuss some further applications, which include bypassing the impossibility results for publicly immunized *outputs* of a backdoored PRG; see Sec. E in the appendix.

**Related works.** Kleptography, as noted, was introduced by Young and Yung [YY96, YY97]; they primarily highlighted the possibility of subverting key generation, and left open the problem of defending against such subversion. Recent work [RTYZ15] has made initial progress on the problem of protecting key generation for specific cryptographic algorithms (trapdoor one-way permutations, pseudorandom generators, and digital signature scheme). However, these techniques are highly tuned to specific algorithms and do not remove arbitrary steganographic channels, which is one of our main goals.

---

<sup>1</sup>Same as [RTYZ15], we only assume the specification  $RG_{\text{SPEC}}$  to be an random oracle, while the implementation  $RG_{\text{IMPL}}$  can be arbitrarily subverted.

Several research threads have studied the kleptographic setting, developing both new attacks and defending mechanisms. In particular, Bellare, Paterson, and Rogaway [BPR14] studied subverted randomized encryption algorithms, building a steganographic channel that leaks secrets bit by bit; they also developed defending mechanisms in the setting where key generation is honest. Such subliminal channel attacks turn out to be the major obstacle in this area, and have been further explored by Ateniese et al. [AMV15], Bellare et al. [BH15, BJK15], Degabriele et al. [DFP15], and Dodis et al. [DGG<sup>+</sup>15]. A common feature of these works [BPR14, BH15, BJK15, DFP15] is to adopt deterministic algorithms and to assume honest key generation. Additionally, these works do not rely merely on testing: Most, in fact, require an a priori “decryptability” condition which demands that every message encrypted using the implementation should be decrypted correctly using the specification. A notable exception is [DFP15]; however, they rely on a watchdog that possesses access to the actual challenger–adversary communication transcript (including the internal state of the challenger).

Other works [MS15, DMSD15] considered defending mechanisms with a “reverse firewall” that is trusted to generate good randomness and can “re-randomize” incoming and outgoing communication. This model is attractive as it may permit quite general feasibility results; on the other hand, it introduces another trusted party (and source of trusted randomness).

In contrast to previous work, our goal is to develop CPA-secure encryption in a much stricter model that does not require strong watchdogs, clean keys, trusted randomness, or decryptability assumptions.

## 2 Definitions and Models

The adversary in kleptographic settings is “*proud-but-malicious*”: The adversary prefers subversion to be “under the radar” of any possible detection; on the other hand, she still wishes exploit her power of subversion to violate security. As explained in the introduction, our central focus will be the challenge of *generically destroying* subliminal channels which may have been adversarially embedded in a subverted algorithm. We briefly recall the notion of a subliminal channel to set the stage for the basic definitions below.

Consider an (honest) randomized algorithm  $A$  which takes an input  $x$  and has additional access to a “secret” bit  $s \in \{0, 1\}$ . The algorithm produces a random output  $y$ , which we assume leaks no information about  $s$ . A fundamental result in steganography [Sim83, Sim86, HLv02] asserts that it is possible to construct a subverted algorithm  $\tilde{A}_z$ , whose behavior is determined by a hidden random string  $z$ , so that

- for all inputs  $x$  and  $s$ , the distribution produced by  $\tilde{A}_z(x, s)$  (including the random selection of  $z$ ) is *identical* to the distribution produced by  $A(x, s)$ , and hence leaks no information about  $s$ ; but,
- with knowledge of  $z$ , the output of  $\tilde{A}_z$  is highly correlated with  $s$ . In particular, an adversary with access to  $z$  can use the output of  $\tilde{A}_z$  to infer  $s$  with high probability.

See Figs. 1a and 1b.

As mentioned above, subliminal channels are a major security obstacle in the kleptographic setting and our main result is a method for destroying them. Intuitively, to achieve this goal, we demand that the adversary cannot learn any extra information from the output of a subverted

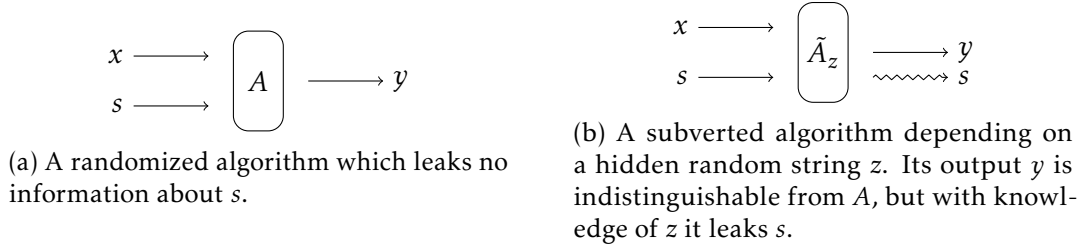


Figure 1: Embedding a subliminal channel in a randomized algorithm  $A$ .

implementation of an algorithm so long as it has passed the checking of the watchdog. We adopt the general kleptographic definitional framework of [RTYZ15], and generalize the notion defined by the surveillance game in [BPR14] to formulate our definition of “destroying a steganographic channel”—this defines a new notion that we call “stego-freeness.” Stego-free specifications for algorithms will be the stepping stones to our final construction of cryptographic primitives (e.g., public-key encryption) with subversion resistance.

In the definitional framework of [RTYZ15], the adversary  $\mathcal{A}$  provides subverted implementations of the cryptographic algorithms for a particular primitive; the challenger  $\mathcal{C}$  must then play the (standard) cryptographic security game for the primitive with  $\mathcal{A}$ . Of course, the challenger uses the subverted implementations during the security game. However, to protect the challenger, there is a *watchdog*  $\mathcal{W}$  who tests the subverted implementations, by comparing them with the specification of the algorithms. The adversary “wins” if she can manufacture implementations so that she can win the security game and—at the same time—evade the detection of  $\mathcal{W}$ . One can arrive at a variety of different definitions based on the order of quantification for  $\mathcal{W}$  and  $\mathcal{A}$ , how exactly  $\mathcal{W}$  is permitted to test the implementations, and whether  $\mathcal{W}$  is given any further information (such as a transcript of the security game). We refer to [RTYZ15] for detailed discussion.

In this paper we will adopt the strongest of the definitions of [RTYZ15] (which gives the watchdog the least power): in their terminology, we will consider a *universal* and *offline* watchdog. In such a definition, the watchdog only tests the implementation once with only oracle access. In particular,  $\mathcal{W}$  has no access to the actual communications during the security game. Moreover, the description of the watchdog is quantified before the adversary.<sup>2</sup> (Thus, security for a particular primitive requires that there is a single watchdog that can protect against all possible adversaries.)

To formalize the notion that the adversary cannot learn unintended information from an implementation, we adapt the surveillance game from [BPR14] (which was defined for symmetric key encryption): specifically, we compare the information leaked by the implementation with that leaked by the specification (or, equivalently, an honest implementation).

**Defining stego-freeness.** We now formally define stego-freeness for any (randomized) algorithm  $G$  under subversion. Following the basic kleptographic models described above, the adversary  $\mathcal{A}$  prepares a (potentially subverted) implementation  $G_{\text{IMPL}}$  of the algorithm  $G$ ; we let  $G_{\text{SPEC}}$  denote the specification of the algorithm. The goal of the adversary is to utilize  $G_{\text{IMPL}}$  to leak secret information *exclusively* to her via the outputs that  $G_{\text{IMPL}}$  produces (as in the discussion above). Stego-freeness means **either** the adversary  $\mathcal{A}$  cannot learn any extra information from the outputs of  $G_{\text{IMPL}}$  (in comparison with that of  $G_{\text{SPEC}}$ ), **or** the subversion can be detected by the watchdog  $\mathcal{W}$  (using oracle

<sup>2</sup>This is stronger than most of the definitions in the literature. The closest one is [DFP15]; however, their watchdog has to take the transcript between  $\mathcal{C}$  and  $\mathcal{A}$  as inputs which implicitly implies the dependence of the running time on  $\mathcal{A}$ .

access to  $G_{\text{IMPL}}$ )—this is characterized by the detection advantage  $\mathbf{Det}_{\mathcal{W},\mathcal{A}}$  below. Depending on how communication is generated and whether the randomized algorithm can take rich inputs, we have a variety of definitions; we begin with the following elementary version for randomized algorithms—such as key generation—that rely only on a length parameter rather than inputs drawn from a large space.

**Definition 2.1** (stego-free, basic form). *Consider a (randomized) algorithm  $G$  with specification  $G_{\text{SPEC}}$ . We say such  $G_{\text{SPEC}}$  is stego-free in the offline watchdog model if there exists a PPT watchdog  $\mathcal{W}$  so that for any PPT adversary  $\mathcal{A}$  playing the following game (see Fig. 2), it satisfies that either*

$$\mathbf{Adv}_{\mathcal{A}} \text{ is negligible, or } \mathbf{Det}_{\mathcal{W},\mathcal{A}} \text{ is non-negligible}$$

where

$$\mathbf{Adv}_{\mathcal{A}}(1^\lambda) = |\Pr[b_C = 1] - 1/2| \quad \text{and} \quad \mathbf{Det}_{\mathcal{W},\mathcal{A}}(1^\lambda) = |\Pr[\mathcal{W}^{G_{\text{IMPL}}}(1^\lambda) = 1] - \Pr[\Pr[\mathcal{W}^{G_{\text{SPEC}}}(1^\lambda) = 1]]|.$$

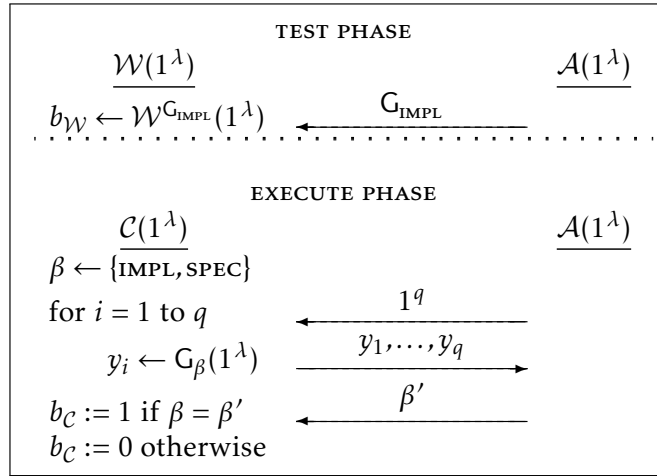


Figure 2: A game for stego-freeness.

Note that the definition requires only non-negligible detection probability on the part of the watchdog. Note that detection probabilities can be directly amplified by repetition.<sup>3</sup>

**Remark 2.2.** *Our constructions actually satisfy a stronger condition that directly relates the advantage to the detection probability. The definition above demands the asymptotic guarantee that either  $\mathbf{Adv}_{\mathcal{A}}$  is negligible or  $\mathbf{Det}_{\mathcal{W},\mathcal{A}}$  is non-negligible. This can be strengthened to demand that for every adversary  $\mathcal{A}$  there is a polynomial  $s(\lambda)$  (with no constant term, so that  $s(0) = 0$ ) so that  $s(\mathbf{Adv}_{\mathcal{A}}(1^\lambda)) \leq \mathbf{Det}_{\mathcal{W},\mathcal{A}}(1^\lambda)$ .*

Observe that an offline watchdog can ensure that the implementation of a *deterministic* algorithm disagrees with its specification with negligible probability when inputs are drawn from a public input distribution. Throughout, we use the term “public” distribution to refer to any efficiently sampleable source that the watchdog can construct, perhaps using  $F_{\text{SPEC}}^i$  and  $F_{\text{IMPL}}^i$ .

<sup>3</sup>Trivial amplification transforms a gap of  $\epsilon$  to  $1 - \delta$  with  $k = \epsilon^{-1} \log(\delta^{-1})$  repetitions. As the watchdog’s running time is fixed independent of the adversary, however, amplification cannot be adapted to a particular non-negligible function. If the watchdog is permitted a number of samples that depends on the adversary, then one can amplify non-negligible detection probability to  $1 - o(1)$  for an infinite set of inputs.



**Lemma 2.3** ([RTYZ15]). *Consider an adversarial implementation  $\Pi_{\text{IMPL}} := (F_{\text{IMPL}}^1, \dots, F_{\text{IMPL}}^k)$  of a specification  $\Pi_{\text{SPEC}} = (F_{\text{SPEC}}^1, \dots, F_{\text{SPEC}}^k)$ , 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  $\exists j \in [k], \Pr[F_{\text{IMPL}}^j(x) \neq F_{\text{SPEC}}^j(x) : x \leftarrow X_\lambda^j]$  is non-negligible, this can be detected by a PPT offline watchdog with non-negligible probability.*

**More general definitions of stego-freeness.** In the above game,  $G$  only takes as input a fixed security parameter (often ignored later in the paper); this definition can capture algorithms like randomness generation and key generation when we instantiate  $G$  to be the corresponding functionality. Besides the security parameter, we can consider algorithms which take richer inputs. Such extensions will be important for our applications, and can significantly complicate the task of destroying an embedded steganographic channel. One note is that for input taken from a small domain, (for example,  $\{0, 1\}$ ), we simply allow the adversary to query the evaluation on all inputs.

Beyond the previous cases, we may consider algorithms taking inputs from a large domain. The most straightforward adaptation permits the adversary to sample  $G_{\text{IMPL}}(1^\lambda, x_i)$  at inputs  $x_i$  of her choice. However, this model suffers from a crippling “input trigger” attack [DFP15] (where the adversary hides some secret information at a particular “trigger” location  $x$  which can be impossible for an offline watchdog to detect); we discuss this in detail later. However, there is a compromise setting that captures many cases of actual interest and permits strong feasibility results. In this setting we permit the adversary to determine inputs to a randomized algorithm  $G$  by specifying a randomized *input generator*  $\text{IG}$ : The input generator may be an arbitrary PPT algorithm with the condition that given  $1^\lambda$  it produces (typically random) outputs of length exactly  $\lambda$ . This implicitly defines the randomized algorithm  $G(1^\lambda, \text{IG}(1^\lambda))$ . In our setting, the watchdog is provided (oracle access) to  $\text{IG}$ , which it may use during its testing of  $G$ . Note that  $\text{IG}$  is not part of the specification of  $G$ , but rather chosen by the adversary during the security game; thus there is no reason to consider subversion of  $\text{IG}$ . Revisiting the security game in this new setting, challenges  $\{y_i\}$  are generated by first sampling  $m_i \leftarrow \text{IG}(1^\lambda)$ , and then obtaining  $y_i \leftarrow G_\beta(1^\lambda, m_i)$  by calling  $G_\beta$  using inputs  $1^\lambda$  and  $m_i$ . Note that the adversary could use  $\text{IG}$  to produce some specific input “triggers” where  $G_{\text{IMPL}}$  deviates from  $G_{\text{SPEC}}$ . This more general notion of stego-freeness (with a “public” input distribution) captures algorithms that take the output of other algorithms as input, which will be critical when we reason about amalgamation of algorithms. See Figure 3 below for a unified game, where the algorithm may take both types of inputs.

**Definition 2.4** (stego-free, general form). *We say that a randomized algorithm  $G$  is stego-free if it satisfies Definition 2.1 with the security game of Figure 3. Note that the PPT input generator  $\text{IG}$  may be determined by the adversary during the game.*

Which of the definitions (2.1 or 2.4) is appropriate for a given randomized algorithm can be determined from context, depending on whether an input generator is specified.

As mentioned above, an even stronger definition is obtained by permitting the adversary to simply choose the input  $m_i$  for each  $y_i$  directly. This notion reflects stego-freeness for algorithms with adversarially chosen inputs. Such a subverted implementation may have a hidden “trigger” that was randomly drawn during the (adversarial) manufacturing process and can permit the adversary to easily win the stego-freeness distinguishing game. In fact, such a trigger attack does not even require that  $G$  be randomized: for example, consider the algorithm  $G_{\text{SPEC}}(1^\lambda, x) := x$ ,



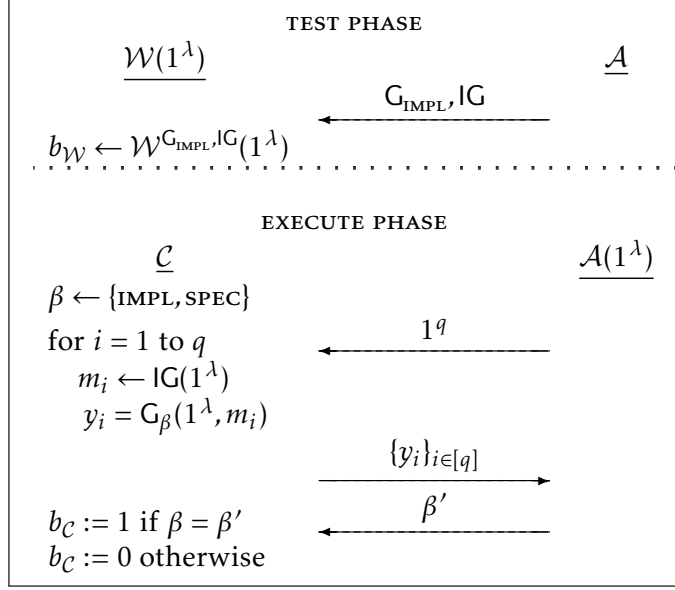


Figure 3: The stego-freeness game with input distribution  $\{1^\lambda\} \times \text{IG}$ .

defined for  $x \in \{0, 1\}^\lambda$ . The adversary then uniformly draws  $z \leftarrow \{0, 1\}^\lambda$  and defines

$$G_{\text{IMPL}}(1^\lambda, x) = \begin{cases} 0^\lambda & \text{if } x = z, \\ x & \text{otherwise.} \end{cases}$$

As the placement of the trigger ( $z$ ) is random, the watchdog cannot detect disagreement between  $G_{\text{IMPL}}$  and  $G_{\text{SPEC}}$ , while the adversary can distinguish these algorithms easily by querying  $z$ . In a practical setting, an algorithm with such an input trigger can leak arbitrary private data to an adversary in a way undetectable to an offline watchdog. This was formally demonstrated in [DFP15] and called an “input-triggered subversion attack.” Nevertheless, we will discuss in Section 4.1 a method for sidestepping this impossibility with only an offline watchdog by assuming some minimum trusted operations, such as “one trusted addition.”<sup>4</sup>

These definitions of stego-freeness are sufficient for capturing most of the interesting use cases we will require (e.g., reflecting key generation and encryption).

**Remark 2.5.** *Following Lemma 2.3, observe that if  $G_{\text{SPEC}}$  is deterministic, an offline watchdog can ensure that inconsistencies ( $G_{\text{SPEC}}(x) \neq G_{\text{IMPL}}(x)$ ) occur with only negligible probability when inputs are sampled from the input distribution  $\text{IG}$  (by drawing and testing a sample). In particular, deterministic algorithms with a public input distribution satisfy stego-freeness in a straightforward fashion.*

**Discussions about stego-freeness and steganography.** We emphasize two properties of these definitions. First, if a proposed specification satisfies such definitions, direct use of the implementation—rather than the specification—preserves the typical security guarantees originally possessed by the

<sup>4</sup>All previous works either simply assume it won’t happen (the decryptability assumption) or employ an omniscient watchdog who has access to the transcript between the challenger  $\mathcal{C}$  and the adversary  $\mathcal{A}$ , (and the secret key of  $\mathcal{C}$ ).

specification. This enables us to provide fairly modular security proofs by designing specifications for each algorithm with stego-freeness.

The second, and more critical, issue pertains to the feasibility of achieving these notions of stego-freeness: in particular, at first glance they appear hopeless. It is known that general steganography is always possible over a channel with sufficient entropy [Sim83, Sim86, HLv02]. This implies that the subverted algorithm  $G_{\text{IMPL}}$  can always produce a sequence of messages that enable the adversary to retrieve secret data from the (public) outputs  $y_1, \dots, y_q$ . In particular, as shown by Bellare, Paterson, Rogaway in the seminal result [BPR14], a subverted randomized encryption algorithm can generate ciphertexts so that the adversary can recover the secret key bit-by-bit from the sequence of ciphertexts. Moreover, the distribution of these subverted ciphertexts is statistically close to the natural, unsubverted ciphertext distribution. To make matters worse, such attacks can be launched even if the subverted implementations are stateless [BJK15]. As a simple example of such subversion in our setting, consider the algorithm  $G_{\text{SPEC}}(1^\lambda)$  which outputs a uniformly random element of  $\{0, 1\}^\lambda$ . Consider then the subverted implementation  $G_{\text{IMPL}}^z(1^\lambda)$  whose behavior is determined by a uniformly random string  $z \leftarrow \{0, 1\}^\lambda$  chosen by the adversary: the algorithm  $G_{\text{IMPL}}^z(1^\lambda)$  outputs a uniformly random element of the set

$$H = \{w \in \{0, 1\}^\lambda \mid \text{lsb}(F_z(w)) = 0\},$$

where  $\text{lsb}(x)$  denotes the least-significant bit of  $x$  and  $F_z(\cdot)$  denotes a pseudorandom function (PRF) with key  $z$ . (Note that elements of  $H$  can be drawn by rejection sampling.) Of course, it is easy for the adversary to distinguish  $G_{\text{IMPL}}$  from  $G_{\text{SPEC}}$  (as  $G_{\text{IMPL}}$  only outputs strings with a particular property that is easily testable by the adversary who has  $z$ ). On the other hand, no watchdog can distinguish these algorithms without breaking the PRF. This suggests that if the user makes only *black-box* use of the subverted implementation of randomized algorithms, it is hopeless to achieve stego-freeness. This motivates the following *non-black-box* model.

**The split-program methodology and trusted amalgamation.** To overcome the steganographic attacks discussed above, we propose a slightly modified model which permits the specification of an algorithm to be split into several components. In this *split-program model*, each component of the implementation is exposed to the watchdog to check, while the challenger will amalgamate the components to yield the fully functional implementation. Of course, the implementation of each component is still presented by the adversary. We permit decomposition into only a constant number of components (independent of input length), with the demand that the desired algorithm can furthermore be expressed as the composition of a constant number of the components. Note that such a “split-program” presentation of an algorithm rules out any gate-by-gate treatment, as the model permits only a constant number of compositions. Intuitively, this simple non-black-box presentation of a randomized algorithm not only provides more opportunity to enforce the malicious implementation to follow a certain pattern, but also enables the watchdog to do more delicate checking on the inner structure.

One example of this framework is the simple split-program method proposed in [RTYZ15] to study certain randomized algorithms: they begin by specifying a (general) randomized algorithm  $G$  as a pair  $(\text{RG}, \text{dG})$  where  $\text{RG}$  is the *randomness generation* algorithm, responsible for generating a uniform random string of appropriate length, and  $\text{dG}$  is a deterministic algorithm that takes the input to the original randomized algorithm  $G$  and the random coins produced by  $\text{RG}$  to produce the final output. They then add to this specification a third deterministic algorithm  $\Phi$  which acts as a kind of “immunization function” for the random bits generated by  $\text{RG}$ . Specifically, given the

implementations  $(\text{RG}_{\text{IMPL}}, \text{dG}_{\text{IMPL}}, \Phi_{\text{IMPL}})$ , the challenger *amalgamates* them by first querying  $r_0 \leftarrow \text{RG}_{\text{IMPL}}$ , “sanitizing” this randomness by passing it into  $\Phi_{\text{IMPL}}$  to receive  $r \leftarrow \Phi_{\text{IMPL}}(r_0)$  and, finally, running  $y \leftarrow \text{dG}_{\text{IMPL}}(r)$ . They show that in several contexts such an “immunization” can preserve security even under subversion. We remark that a simple decomposition and amalgamation of this form cannot destroy steganography in general, and we show an explicit attack in this model; see Sec. 3.1.

To reflect such trusted amalgamation in our security model, we permit the challenger to carry out (a constant number of) compositions without molestation; that is, the notion of “composition” is protected from adversarial subversion. This can be inferred from the definition of the specification, and it is implicit in the security games defined, e.g., in Figure 11, and in Figure 12, in the appendix.

As mentioned above, the split-program method proposed in [RTYZ15] permitted them to establish security (in the kleptographic model) for specific cryptographic primitives. In this paper, we will show that this general methodology has remarkable power against subversion: by further decomposition and amalgamation, we show it is possible to generically destroy steganographic channels. This provides us a family of tools for developing kleptographically-secure cryptography without abandoning randomized algorithms.

**Stateful algorithms.** As discussed above, steganographic attacks can be launched even if the implementation is stateless. To simplify our presentation, most of our discussion adopts this stateless assumption. However, adaptations of our techniques can provide security even for stateful implementations (in the sense that each functionality maintains internal state). These amplified results require a slightly stronger watchdog (whose running time can depend on the adversary) and more detailed control of the subverted algorithm to ensure that they receive inputs from public distributions. See Remark 3.2 in Sec. 3.2 for more discussion.

### 3 Eliminating Subliminal Channels in Randomized Algorithms

In this section, we will present our main result: provable destruction of any subliminal channels in subverted implementations of randomized algorithms.

First, we motivate our new constructions by showing that the steganographic attacks of [BPR14, BJK15] can still be carried out in the simple split-program model introduced by [RTYZ15]. The attack succeeds even if the associated “immunizing function”  $\Phi$  is a trusted hash function modeled as a random oracle. This indicates that some stronger form of immunization is necessary for destroying steganography.

Next, we present our main result in Section 3.2: a generic transformation that destroys steganography in randomized algorithms. The basic technique relies on “double-splitting” the randomness generating procedure coupled with a public immunizing function: specifically, randomness generation is expressed as two algorithms  $\text{RG}_0, \text{RG}_1$  in conjunction with a public immunizing function  $\Phi$ . We prove that with this “extra” split, when the outputs  $r_0 \leftarrow \text{RG}_0, r_1 \leftarrow \text{RG}_1$  are independently sampled from  $\text{RG}_i$  and mixed with  $\Phi$ , the new specification of the randomness generation algorithm is indeed stego-free in the random oracle model. Note that all components, including the immunizing function  $\Phi$ , are subject to subversion. With randomness generation cleaned in this way, we can further destroy the subliminal channel in a large class of randomized algorithms. These results enable us to overcome the major obstacle in designing cryptographic specifications that satisfy subversion resistance.

**Transition to the standard model.** Finally, we consider how to achieve stego-freeness without a random oracle. The main observation is that the watchdog can guarantee that each copy of a particular  $\text{RG}_{\text{IMPL}}^i$  provides at least  $\log n$  bits of (min-)entropy. If we are willing to have more components for randomness generation, we can accumulate entropy using a simple immunizing function and stretch the result using a PRG. See Sec. D in the appendix.

### 3.1 Impossibility of publicly immunizing a single random source

Previous works [BPR14, BJK15] demonstrated that if a (subverted) randomized algorithm is used in a black-box fashion, a subliminal channel can always be embedded in its output distribution. Here we point out that a similar attack exists even if we adopt the techniques described in [RTYZ15] which (i.) split the algorithm  $G$  into  $\text{RG}$  (which generates randomness) and  $\text{dG}$  (which is deterministic), (ii.) introduce  $\Phi$ —an “immunizing” function—and, ultimately, (iii.) generate output via the composition  $\text{dG} \leftarrow \Phi \leftarrow \text{RG}$ . Here  $\Phi$  is responsible for “cleaning” the randomness produced by the possibly subverted randomness generator  $\text{RG}$ . In fact, this approach can fail even in the most generous setting when  $\Phi$  is given by a random oracle and the adversary only subverts  $\text{RG}$ .

The attack is a straightforward adaptation of the techniques from [BPR14, BJK15]: the subverted implementation  $\text{RG}_{\text{IMPL}}$  can evaluate  $\text{dG}$  and appropriately query the random oracle  $\Phi$  during the procedure of rejection sampling. It is easy, then, to arrange for the the output of  $\text{RG}_{\text{IMPL}}$  to be biased in a way only detectable by the adversary. While a generic attack is possible, for concreteness we present an attack on a subverted public-key cryptosystem which permits the adversary to effortlessly determine the (plaintext) message bit. This indicates that more advanced non-black-box techniques are necessary to remove steganographic channels in general.

A detailed description appears in Sec. B.

### 3.2 Purifying randomness via double splitting

The attacks described above (and by [BPR14, BJK15]) demonstrate a core obstacle to defending randomized algorithms against subversion: the random coins drawn by the subverted implementation can be biased—for example, by rejection sampling—even if they are then “immunized” by a random oracle. On the other hand, the split-program model intuitively offers the watchdog an opportunity for fine-grained testing: Is the situation really no better than the purely black-box model? Let us take a closer look from the security analysis point of view.

Intuitively, specifications of the form  $(\text{RG}_{\text{SPEC}}, \text{dG}_{\text{SPEC}}, \Phi_{\text{SPEC}})$  can provide security in a kleptographic setting if the immunization function  $\Phi$  can suitably interfere with generation of biased output by the implementation of  $\text{RG}_{\text{SPEC}}$ . To simplify our presentation, we assume throughout that  $\text{RG}_{\text{SPEC}}$  produces at least  $\lambda$  bits of randomness; this does not affect the generality of the results. (Our techniques can be adapted to a low-entropy setting with some changes to running time of the watchdog.<sup>5</sup>) An important feature in this setting is that an offline watchdog  $\mathcal{W}$  can at least guarantee that the output  $r_0$  of  $\text{RG}_{\text{IMPL}}$  is *unpredictable* to the adversary  $\mathcal{A}$ . Otherwise, the distribution given by  $\text{RG}_{\text{IMPL}}$  would have significant (non-negligible) collision probability,<sup>6</sup> which can be easily

<sup>5</sup>Observe that if  $\text{RG}_{\text{SPEC}}$  produces only  $O(\log n)$  random coins then an offline watchdog, by a suitable regimen of repeated sampling, can empirically approximate (with high probability) the distribution of  $\text{RG}_{\text{IMPL}}$  with high accuracy. This can be directly compared with  $\text{RG}_{\text{SPEC}}$  using distance in total variation. Note that such a watchdog requires a number of samples polynomial in the resulting error.

<sup>6</sup>Observe that if  $D$  is a probability distribution on a set  $X$ , the optimal strategy for predicting the result of drawing

tested by  $\mathcal{W}$  who simply draws two samples and rejects if it observes a collision. (As with the other tests we discuss, the success of this test can be amplified by repetition.) On the other hand, the collision probability of  $\text{RG}_{\text{SPEC}}$  is negligible. This suggests the intuition that  $\Phi_{\text{SPEC}}(r_0)$  appears to  $\mathcal{A}$  to be a randomly drawn value if  $\Phi_{\text{SPEC}}$  is a random oracle. Unfortunately,  $\mathcal{A}$  also holds the backdoor  $z$  which may contain information about the final output  $r = \Phi(r_0)$  generated by the sampling and “cleaning” process. In particular, as shown in the attack, the subverted implementation has full access to  $\Phi_{\text{SPEC}}$  and may thus bias the output  $r = \Phi(r_0)$  as a function of  $z$ , which can be noticed by  $\mathcal{A}$ .

To circumvent the above obstacle, we introduce a new technique that further splits randomness generation into *two* random algorithms,  $\text{RG}_{\text{SPEC}}^0$  and  $\text{RG}_{\text{SPEC}}^1$ , and combines their outputs using an immunization function  $\Phi_{\text{SPEC}}$ ; we shall see that this mechanism can destroy subliminal channels. In general, in the trusted amalgamation model, (see Definition A.1 in appendix A.1) the user runs  $\text{RG}_{\text{IMPL}}^0$  and  $\text{RG}_{\text{IMPL}}^1$  independently and passes the joint outputs to  $\Phi_{\text{IMPL}}$ ; the final output will have the form  $r = \Phi_{\text{IMPL}}(r_0 \circ r_1)$  (where  $\circ$  denotes concatenation). The main idea behind this strategy is that it frustrates attempts by  $\text{RG}_{\text{IMPL}}^0$  and  $\text{RG}_{\text{IMPL}}^1$  to launch sophisticated rejection-sampling because the final output is not fully determined by either output. (In particular, neither can evaluate  $\Phi_{\text{IMPL}}(r_0 \circ r_1)$  during the generation of  $r_0$  or  $r_1$ .) In this way, if  $\Phi_{\text{SPEC}}$  is modeled as a random oracle, the final output  $\Phi_{\text{SPEC}}(r_0 \circ r_1)$  will be uncorrelated with  $\mathcal{A}$ 's state (which includes both  $\mathcal{A}$ 's random oracle queries and  $z$ ). Now we can safely claim that  $r$  looks uniform even to  $\mathcal{A}$ .<sup>7</sup>

We remark on a similarity between the setting above and the topic of randomness extractors [NZ96]. Recall that an extractor  $E$  is a deterministic function which takes as input a number of imperfect random sources  $X_1, X_2, \dots, X_k$  and produces a nearly uniform output  $E(X_1, \dots, X_k)$ . It is a fact that extraction is not possible when  $k = 1$ —no fixed deterministic function  $E$  can produce clean randomness from a source with bounded min-entropy. On the other hand, it is possible when  $k \geq 2$  [CG88, Bou05, CZ15].

With this approach, we demonstrate a qualitative advantage of the split-program methodology. We first describe a stego-free specification of randomness generation in Fig. 4, and then proceed to give an immunization strategy for arbitrary randomized algorithms so long as they have a public input distribution (or a small input domain). We apply these tools in next section to construct an IND-CPA public key encryption scheme that retains security under subversion; note that this is impossible to achieve if randomness generation is used as a black-box (even if it is separated as an individual component).

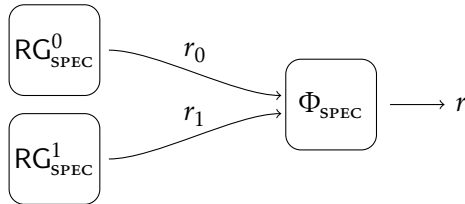


Figure 4: A stego-free specification for randomness generation.

---

an element of  $X$  according to  $D$  is simply to guess  $\max_x D(x)$ . If this maximum probability is  $\epsilon$ , then the collision probability of  $D$ , equal to  $\sum_x D(x)^2$ , is at least  $\epsilon^2$ .

<sup>7</sup>We remark that  $\Phi_{\text{IMPL}}$  can be subverted, but it is a deterministic function with a public input distribution, the inconsistency can be ensured to be at only a negligible fraction of places due to the watchdog, see Lemma 2.3.

**Theorem 3.1.** Consider a randomness generation algorithm  $\text{RG}$  with specification  $(\text{RG}_{\text{SPEC}}^0, \text{RG}_{\text{SPEC}}^1, \Phi_{\text{SPEC}})$  as described in Fig. 4:

- $\text{RG}_{\text{SPEC}}^0$  and  $\text{RG}_{\text{SPEC}}^1$ , given  $1^\lambda$ , output uniformly random strings of length  $\lambda$ ;
- $\Phi_{\text{SPEC}}$  is a hash function so that  $\Phi_{\text{SPEC}}(w)$  has length  $\lceil |w|/2 \rceil$ ; and
- the specification for  $\text{RG}(1^\lambda)$  is the amalgamation  $\Phi_{\text{SPEC}}(\text{RG}_{\text{SPEC}}^0(1^\lambda), \text{RG}_{\text{SPEC}}^1(1^\lambda))$ .

Then  $\text{RG}_{\text{SPEC}}$  is stego-free in the trusted amalgamation model (see Def. A.1 in Sec. A.1 in the appendix) if  $\Phi_{\text{SPEC}}$  is modeled as a random oracle.

*Proof.* If the specification of Fig. 4 is not stego-free, then for any offline watchdog  $\mathcal{W}$  there is an adversary  $\mathcal{A}$  that can prepare an implementation  $\text{RG}_{\text{IMPL}} := (\text{RG}_{\text{IMPL}}^0, \text{RG}_{\text{IMPL}}^1, \Phi_{\text{IMPL}})$  satisfying the following: (1.)  $\mathcal{W}$  cannot distinguish  $\text{RG}_{\text{IMPL}}$  from  $\text{RG}_{\text{SPEC}}$  via oracle access; (2.) The adversary  $\mathcal{A}$  can distinguish output of  $\text{RG}_{\text{IMPL}}$  from  $\text{RG}_{\text{SPEC}}$ , i.e., she can win the game defined in Figure 11 in supporting material, Sec. A.1. We will then define an offline watchdog such that these two conditions cannot hold simultaneously for any adversary.

An offline watchdog. The watchdog  $\mathcal{W}$ 's strategy is as follows:  $\mathcal{W}$  first draws a constant number of samples from  $\text{RG}_{\text{IMPL}}^0$  and  $\text{RG}_{\text{IMPL}}^1$ ; if  $\mathcal{W}$  observes a collision in either distribution, it rejects the implementation outright (as collisions are negligible in SPEC). Next,  $\mathcal{W}$  draws pairs of samples (again) from  $\text{RG}_{\text{IMPL}}^0$  and  $\text{RG}_{\text{IMPL}}^1$  and evaluates  $\Phi_{\text{IMPL}}$  on (the concatenation of) each pair to ensure that the result is consistent with  $\Phi_{\text{SPEC}}$ . (As usual, this testing involves only  $O(1)$  samples and can be trivially amplified by repetition.)

Next, we will show, for any PPT adversary  $\mathcal{A}$ , if the detection probability  $\text{Det}_{\mathcal{W}, \mathcal{A}}$  is negligible, then the advantage  $\text{Adv}_{\mathcal{A}}$  will also be negligible, thus the two conditions cannot hold simultaneously.

Game transitions. We will go through the security game part of Def. A.1 step by step. Without loss of generality, we assume the challenge  $r$  contains only one element (i.e.,  $q = 1$  in Def. A.1).

In Game-0, the adversary  $\mathcal{A}$  prepares subverted implementations  $\text{RG}_{\text{IMPL}} := (\text{RG}_{\text{IMPL}}^0, \text{RG}_{\text{IMPL}}^1, \Phi_{\text{IMPL}})$ ; we let  $Q$  be the set of random oracle queries  $\mathcal{A}$  made during preparation of  $\text{RG}_{\text{IMPL}}$ .

The challenger  $\mathcal{C}$  samples from  $\text{RG}_{\text{IMPL}}^0$  and  $\text{RG}_{\text{IMPL}}^1$  respectively and receives  $r_0$  and  $r_1$ ; then  $\mathcal{C}$  evaluates  $\Phi_{\text{IMPL}}$  at  $r_0 \circ r_1$  and sends the output  $r$  as the challenge to  $\mathcal{A}$ . Let  $Q_b$  (for  $b = 0, 1$ ) be the set of random oracle queries made by  $\text{RG}_{\text{IMPL}}^b$  before outputting  $r_b$ . All random oracle queries will be (consistently) answered with randomly chosen values.

Game-1 is identical to Game-0, except that  $\Phi_{\text{IMPL}}$  is replaced with  $\Phi_{\text{SPEC}}$ ; Game-2 is identical to Game-1, except that the challenger  $\mathcal{C}$  simply chooses a uniform  $r$  and directly sends it to  $\mathcal{A}$  as the challenge; Game-3 is identical to Game-2, except that  $\text{RG}_{\text{IMPL}}$  is completely replaced with  $\text{RG}_{\text{SPEC}}$ ; Game-4 is identical to Game-3, except that  $r$  is generated as in Game-0, but the challenger uses  $\text{RG}_{\text{SPEC}}$  instead.

Probabilistic analysis. We will analyze the gaps of each game transition conditioned on the event that  $\text{Det}_{\mathcal{W}, \mathcal{A}}$  is negligible, we denote this event as  $D_{\mathcal{W}}$ . All the probability gap would be under the condition  $D_{\mathcal{W}}$ . For brevity, we assume the condition  $D_{\mathcal{W}}$  without mention for the analysis of each probability gap.

First, since  $\Phi_{\text{SPEC}}$  is a deterministic function with a public input distribution (the output distribution of  $\text{RG}_{\text{IMPL}}^0 \times \text{RG}_{\text{IMPL}}^1$ ), following Lemma 2.3,

$$\Pr \left[ \Phi_{\text{IMPL}}(r_0 \circ r_1) \neq \Phi_{\text{SPEC}}(r_0 \circ r_1) : r_0 \leftarrow \text{RG}_{\text{IMPL}}^0, r_1 \leftarrow \text{RG}_{\text{IMPL}}^1 \right] \leq \text{negl}(\lambda).$$



Otherwise, the gap is non-negligible and the watchdog  $\mathcal{W}$  will notice the inconsistency (with non-negligible probability). It follows that replacing  $\Phi_{\text{IMPL}}$  with  $\Phi_{\text{SPEC}}$  would incur only a negligible difference (conditioned on  $D_{\mathcal{W}}$ ), thus:

$$|\Pr[b_C = 1 \text{ in Game-0}] - \Pr[b_C = 1 \text{ in Game-1}]| \leq \text{negl}(\lambda).$$

Second, we will argue that the probability that  $r_0 \circ r_1$  is ever queried (falling in  $Q \cup Q_0 \cup Q_1$ ) is negligible, and now we are in Game-1 using  $\Phi_{\text{SPEC}}$  which is a random oracle.

It is easy to see that  $\Pr[r_0 \circ r_1 \in Q] \leq \text{negl}(\lambda)$ ; otherwise, the watchdog will observe a collision in  $\text{RC}_{\text{IMPL}}^i$  with non-negligible probability. To see this, let  $R_0 = \{r_0 \mid \exists r_1, r_0 \circ r_1 \in Q\}$ , note that  $Q, R_0$  are only polynomially large. If the probability  $r_0 \circ r_1$  falls into  $Q$  (thus the probability that  $r_0$  falls into  $R_0$ ) is non-negligible, say  $\delta$ , that means  $r_0$  will be generated by  $\text{RC}_{\text{IMPL}}^0$  with probability at least  $\delta_0 = \delta/\text{poly}(\lambda)$ . It follows that the collision probability that  $\text{RC}_{\text{IMPL}}^0$  produces the same output  $r_0$  would be  $\delta_0^2$ . While on the other hand,  $\text{RC}_{\text{SPEC}}^0$  produces uniform bits, the collision probability (that  $\text{RC}_{\text{IMPL}}^0$  produces the same uniform output string  $r_0$ ) would be negligible. Thus the watchdog can easily distinguish  $\text{RC}_{\text{IMPL}}^0$  from  $\text{RC}_{\text{SPEC}}^0$  when drawing, say 2 samples.

Similarly, we bound the probability for  $Q_0, Q_1$ . Let  $R_{0,1} = \{r_1 \mid \exists r_0, r_0 \circ r_1 \in Q_0\}$ . Since  $\text{RC}_{\text{IMPL}}^0, \text{RC}_{\text{IMPL}}^1$  are independently run, the probability that  $r_1$  falls into the polynomially large set  $R_{0,1}$  would be negligible; otherwise,  $\text{RC}_{\text{IMPL}}^1$  outputs  $r_1$  with a non-negligible probability, then  $\mathcal{W}$  can notice the difference between implementations and the specifications by identifying collisions. Thus  $\Pr[r_0 \circ r_1 \in Q_0] \leq \Pr[r_1 \in R_{0,1}] \leq \text{negl}(\lambda)$ . The same holds for  $Q_1$ .

The adversary  $\mathcal{A}$  is holding the set of random oracle queries  $Q$ , and a backdoor  $z$ . The only way  $r$  may correlate with  $z$  is that  $r_0 \circ r_1$  is queried during the execution of  $\text{RC}_{\text{IMPL}}^0$ , or  $\text{RC}_{\text{IMPL}}^1$ . If  $r_0 \circ r_1 \notin Q \cup Q_0 \cup Q_1$ ,  $r_0 \circ r_1$  will be independent with  $\mathcal{A}$ 's view  $(Q, z)$ , thus  $r = \Phi_{\text{SPEC}}(r_0 \circ r_1) = \text{RO}(r_0 \circ r_1)$  looks uniform to  $\mathcal{A}$ . We can claim that:

$$|\Pr[b_C = 1 \text{ in Game-1}] - \Pr[b_C = 1 \text{ in Game-2}]| \leq \text{negl}(\lambda).$$

Next, it is easy to see that  $\Pr[b_C = 1 \text{ in Game-2}] = \Pr[b_C = 1 \text{ in Game-3}]$  since the adversary receives the identical challenge. Also,  $\Pr[b_C = 1 \text{ in Game-3}] = \Pr[b_C = 1 \text{ in Game-4}]$  since querying  $\text{RC}_{\text{SPEC}}$  yields a uniform output  $\text{RO}(u_0 \circ u_1)$ , where  $u_0, u_1$  are uniformly chosen.

To conclude, conditioned on  $D_{\mathcal{W}}$ , we have:

$$|\Pr[b_C = 1 \text{ in Game-0}] - \Pr[b_C = 1 \text{ in Game-4}]| \leq \text{negl}(\lambda).$$

Observe that Game-0 corresponds to the case that challenger flips a coin to be 0, i.e.,  $\mathcal{C}$  uses  $\text{RC}_{\text{IMPL}}$  to generate the challenge messages, while Game-4 corresponds to the case that  $b = 1$ , when  $\mathcal{C}$  uses  $\text{RC}_{\text{SPEC}}$ . It follows that:

$$\text{Adv}_{\mathcal{A}} = |\Pr[b_C = 1] - 1/2| \leq \text{negl}(\lambda).$$

Combine all above, we can conclude that the  $\text{RC}_{\text{SPEC}}$  defined in Figure 4 is stego-free.  $\square$

**Implementation considerations.** Practical deployment of such splitting and amalgamation—especially as it requires independence and (for the watchdog) copying internal state—clearly requires detailed consideration of the particular computational environment. In general, there are two natural approaches to achieve independence: The most modular approach relies on modern lightweight virtualization layers such as Docker [Doc13] to insulate individual copies of



the adversary’s code; we remark that this also permits state duplication (which may be necessary for the watchdog in the stateful case). More aggressive complete virtualization is also possible, but more cumbersome. An alternate approach relies on constraining the source code (or the compiler) to directly limit I/O and system calls; this has the advantage that the components can be run efficiently in the native environment. Finally, there may also be settings where it is possible to isolate the program in the architectural/hardware layer or physically separate various components (e.g., using Intel’s secure isolation gateway, or even move one  $RG^i$  outside the user’s computer, use a random beacon, etc.).

**Remark 3.2 (Stateful algorithms).** *The previous construction gives a concrete analysis of how a simple watchdog can secure a split specification. At the end of Section 2, we remarked that our results hold even if the implementations are stateful, in the sense that they may maintain local state that persists between executions.*

*To see this in the example above, we need to ensure that the implementation  $RG_{\text{IMPL}}^b$  produces unpredictable outputs (even to the adversary) for polynomially many invocations. Even if the implementation keeps some internal state, we can recover the same result by adopting a slightly more advanced watchdog which is permitted to have running time that depends on the adversary. In particular, if the adversary is permitted to run the implementation  $k$  times, the watchdog may, for each  $1 \leq k' \leq k$ , (1.) run the implementation  $k'$  times to generate a particular internal state  $s$ , and then (2.) run the implementation a constant number of times with state  $s$  to test for collisions. By this process, the watchdog can guarantee the entropy of a sample conditioned on previous samples. A straightforward analysis shows that if there is a particular  $k'$  for which the resulting distribution has non-negligible collision probability (with non-negligible probability in the state produced by the first  $k' - 1$  invocations), the watchdog can detect this (with non-negligible probability). This was also observed in [RTYZ15]. Note that the actual procedure carried out by the watchdog is still universal, but the running time (and number of samples) may be adapted to the adversary. Note also that no strictly universal watchdog (whose running time is independent of the adversary) can possibly detect an algorithm that suddenly becomes deterministic after some (polynomial) number of steps determined by the adversary.*

*We further remark that the above can be easily generalized to stateful algorithm having public input distributions.*

### 3.3 Stego-free specifications for randomized algorithms; a general transformation

Now we are ready to establish the general result which yields a stego-free specification for any randomized algorithm in the trusted amalgamation model (see Definition A.2 in Sec. A.2 in the appendix). We use a randomized algorithm with a public input distribution as a running example; this can be generalized directly to the setting with an extra small input domain. As discussed in Section 2, this already covers many of the interesting cases such as key generation and bit encryption.

**The transformation.** Consider a randomized algorithm  $G$  which uses  $\lambda = \lambda(n)$  random bits for inputs of length  $n$ . Let  $(dG, RG)$  denote the natural specification of  $G$  that isolates randomness generation, so that  $RG(1^\lambda)$  produces  $\lambda$  uniformly random bits and  $dG(r, x)$  is a deterministic algorithm so that for every  $x$ ,  $G(x)$  is equal to  $dG(RG(1^\lambda(n), x))$  for  $n = |x|$ . (“Equal” here simply means these have identical output distributions.) As described above, we consider the transformed specification for  $G$  of the form  $(RG_1, RG_2, \Phi, dG)$  where  $dG$  is as above, both  $RG_1(1^\lambda)$  and  $RG_2(1^\lambda)$

produce  $\lambda$  uniformly random bits, and  $\Phi$  is a hash function that carries strings of length  $2k$  to strings of length  $k$ . (See Fig. 5 below, which shows the transformation applied to an algorithm with a public input distribution generated by IG.) We will prove that when  $\Phi$  is given by a random oracle, this is a stego-free specification of  $G$ .

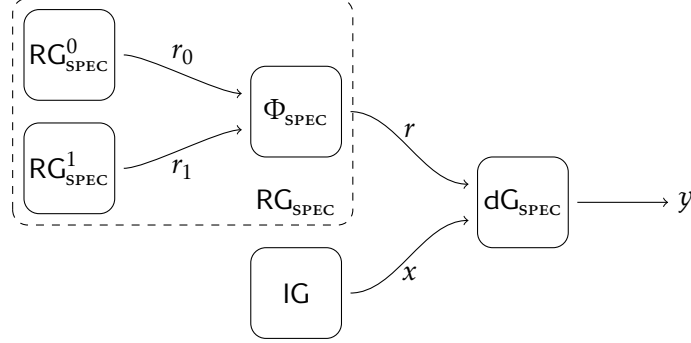


Figure 5: A stego-free specification for randomized algorithm  $G$ , where  $x$  is obtained from IG.

Security analysis. The intuition is as follows: We know that a subverted implementation  $dG_{\text{IMPL}}$  of a deterministic algorithm will be consistent with its specification with an overwhelming probability when the inputs are sampled from a public distribution (Lemma 2.3). In this case, it is generated by  $\text{IG} \times \text{RG}_{\text{IMPL}}$ . Thus  $dG_{\text{IMPL}}$  can be considered as good as  $dG_{\text{SPEC}}$  when restricted to this public input distribution. Furthermore,  $\text{RG}_{\text{SPEC}}$  is stego-free, as discussed above. Thus all implementations can be replaced with their specifications, and the security follows easily. We record this in the theorem below, and we defer the detailed proof to Appendix C.

**Theorem 3.3.** *For any randomized algorithm  $G$ , consider the specification  $G_{\text{SPEC}} := (\text{RG}_{\text{SPEC}}, dG_{\text{SPEC}})$ , where  $\text{RG}_{\text{SPEC}}$  generates  $\lambda = \lambda(n)$  bits of uniform randomness and  $dG_{\text{SPEC}}$  is deterministic. Let  $(\text{RG}_{\text{SPEC}}^0, \text{RG}_{\text{SPEC}}^1, \Phi_{\text{SPEC}})$  be the double-split specification of  $\text{RG}_{\text{SPEC}}$  as above. If (1)  $\text{RG}_{\text{SPEC}}^0(1^\lambda)$  and  $\text{RG}_{\text{SPEC}}^1(1^\lambda)$  output  $\lambda$  uniform bits; (2)  $\Phi_{\text{SPEC}}$  takes  $r_0 \circ r_1$  as input, and outputs  $r$  (so it maps strings of length  $2k$  to strings of length  $k$ ) (see Fig. 5), then  $G_{\text{SPEC}}$  is stego-free with a trusted amalgamation (according to Def. A.2 in appendix. A.2). Here  $\Phi_{\text{SPEC}}$  is modeled as a random oracle, and  $\Phi_{\text{IMPL}}^0, \Phi_{\text{IMPL}}^1$  are executed independently.*

It is straightforward to generalize the result above to support algorithms with an extra small size (polynomially large  $q$ ) input as they essentially expand the input distribution to be  $q$  distributions, which can all be checked by the watchdog.

**Corollary 3.4.** *For any randomized algorithm  $G$ , consider a specification  $G_{\text{SPEC}} := (\text{RG}_{\text{SPEC}}, dG_{\text{SPEC}})$ , where  $\text{RG}_{\text{SPEC}}$  is defined as  $(\text{RG}_{\text{SPEC}}^0, \text{RG}_{\text{SPEC}}^1, \Phi_{\text{SPEC}})$  (as in Fig. 4): If  $dG_{\text{SPEC}}$  takes  $r, x$  and  $m$  as input, where  $x$  is generated by a sampler IG, and  $m$  is taken from a polynomial size public domain  $D$ . Then the specification  $G_{\text{SPEC}}$  is stego-free with a trusted amalgamation (according to Def. A.2 in appendix Sec. A.2) if  $\Phi_{\text{SPEC}}$  is modeled as a random oracle.*

## 4 Subversion-Resistant Encryption

Encryption is a fundamental cryptographic tool, and is well-understood in conventional settings. Unfortunately, encryption in the kleptographic setting has been a particular challenge; existing

solutions do not offer satisfactory CPA security or protection from subversion of key generation. The main motivation for this paper is to develop subversion resistant encryption schemes (both symmetric and asymmetric) that meet these elementary criteria.

To briefly survey the state-of-the-art: Previous results [BPR14, DFP15, BH15], in order to protect against subliminal channel attacks, adopt a special kind of *deterministic* encryption scheme. Note, however, that IND-CPA security is impossible for deterministic public-key encryption (or deterministic, stateless symmetric-key encryption). (An important exception is [DMSD15], which assumes that a trusted, subversion-free, party re-randomizes all communication.) These results permit only the encryption algorithm to be subverted: key generation and decryption are assumed to be honest. A further “decryptability” assumption was adopted in [BPR14, BH15] to ensure security; it assumes that every ciphertext generated by the subverted implementation can be correctly decrypted using the specification. While this helps to achieve security, it seems difficult to justify; see the criticism in [DFP15]. In general, then, correctness is not placed on the same footing as security (which is established via the specification in tandem with watchdog testing), but is rather provided by fiat.

In this section, we will address all these concerns.

#### 4.1 Subversion-resistant symmetric encryption

We first construct subversion-resistant symmetric-key encryption in the case when all algorithms are subject to subversion. We then discuss correctness and how to remove the “decryptability” assumption.

**Defining subversion resistance and correctness.** We follow the definitional framework of [RTYZ15] to define a subversion-resistant symmetric-key encryption scheme. The definition formalizes the intuition that an encryption scheme is secure under subversion if IND-CPA security can be preserved in the “subverted security game” conditioned on the implementation evading detection. Recall that in the “subverted security game” the challenger uses the subverted implementations of all the algorithms. Initially, we focus on *stateless* encryption algorithms, and discuss how to relax this assumption when we discuss large message spaces.

**Definition 4.1.** A (stateless) symmetric-key (bit) encryption scheme with specification  $\mathbf{E}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}, \text{Enc}_{\text{SPEC}}, \text{Dec}_{\text{SPEC}})$  is subversion-resistant in the offline watchdog model if there exists a PPT watchdog  $\mathcal{W}$  so that, for any PPT adversary  $\mathcal{A}$  playing the game described in Figure 6, either

$\mathbf{Adv}_{\mathcal{A}}$  is negligible, or  $\mathbf{Det}_{\mathcal{W}, \mathcal{A}}$  is non-negligible.

where  $\mathbf{Adv}_{\mathcal{A}}(1^\lambda) = |\Pr[b_C = 1] - 1/2|$ , and,  $\mathbf{Det}_{\mathcal{W}, \mathcal{A}}(1^\lambda) = |\Pr[\mathcal{W}^{\mathbf{E}_{\text{IMPL}}}(1^\lambda) = 1] - \Pr[\Pr[\mathcal{W}^{\mathbf{E}_{\text{SPEC}}}(1^\lambda) = 1]]|$ .

Correctness under subversion is an overlooked, but fundamental property. For example, one can imagine that the adversary “hates” a certain message  $m$  that is unknown to the watchdog (e.g., “your cryptosystem is subverted”). The subverted implementation can then check whether the plaintext matches  $m$  and, if so,  $\text{Dec}_{\text{IMPL}}$  outputs an arbitrary value other than  $m$ . This can be used by the adversary to effectively implement censorship. We say a symmetric-key encryption scheme is *correct* under subversion if the following holds:

$$\forall M, \Pr \left[ \begin{array}{l} \text{Dec}_{\text{IMPL}}(C) \neq M : \\ C \leftarrow \text{Enc}_{\text{IMPL}}(K, M), K \leftarrow \text{KG}_{\text{IMPL}}(1^\lambda) \end{array} \right] \leq \text{negl}(\lambda),$$

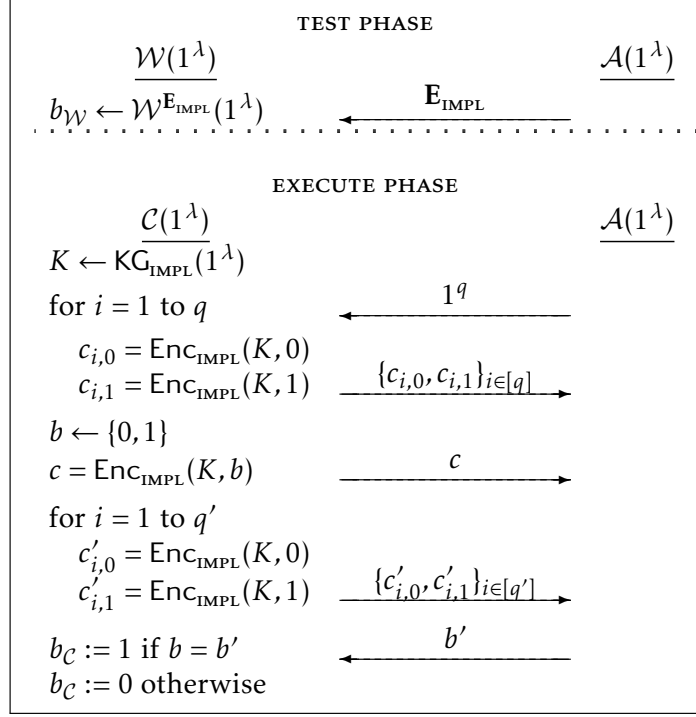


Figure 6: Subversion-resistant symmetric-key bit encryption. (The stateless case.)

where the probability is over the choice of  $K$  and the coins used in  $\text{Enc}_{\text{IMPL}}$ .

**Constructing subversion resistant symmetric-key bit encryption.** We proceed to design a specification for any secure symmetric-key encryption scheme so that the subliminal channels in all of the algorithms can be destroyed. With the general tool we developed in Section 3.2, we will “immunize” the algorithms one by one.

First, the (symmetric-key) key generation algorithm simply takes the security parameter as input and produces a random element in the key space. Following Thm. 3.1 directly, the specification  $\text{KG}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}^0, \text{KG}_{\text{SPEC}}^1, \Phi_{\text{SPEC}}^{\text{KG}})$  (see Fig. 7) is stego-free; here  $\text{KG}_{\text{SPEC}}^0, \text{KG}_{\text{SPEC}}^1$  both output random elements in the key space  $\mathcal{K}$ , and  $\Phi_{\text{SPEC}}^{\text{KG}} : \mathcal{K} \times \mathcal{K} \rightarrow \mathcal{K}$ , is modeled as a random oracle. As discussed above, the random oracle assumption can be removed if we allow randomness generation to be further subdivided; see appendix. D.

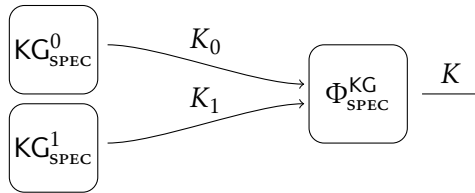


Figure 7: Stego-free specification for key generation algorithm in a symmetric-key bit encryption scheme.

Next, consider the encryption algorithm; we focus only on bit encryption. In this case, the

encryption algorithm takes three inputs: a security parameter, a small domain input (a bit), and a pair given by the random coin and the key, which come from a public input distribution  $\text{RG}_{\text{IMPL}} \times \text{KG}_{\text{IMPL}}$ . From Corollary 3.4, the specification  $\text{Enc}_{\text{SPEC}} := (\text{RG}_{\text{SPEC}}, \text{dEnc}_{\text{SPEC}})$  (as described in Fig. 8), where  $\text{RG}_{\text{SPEC}} := (\text{RG}_{\text{SPEC}}^0, \text{RG}_{\text{SPEC}}^1, \Phi_{\text{SPEC}}^{\text{RG}})$  defined as in Fig. 4, is stego-free if  $\Phi_{\text{SPEC}}^{\text{RG}}$  is assumed to be a random oracle.

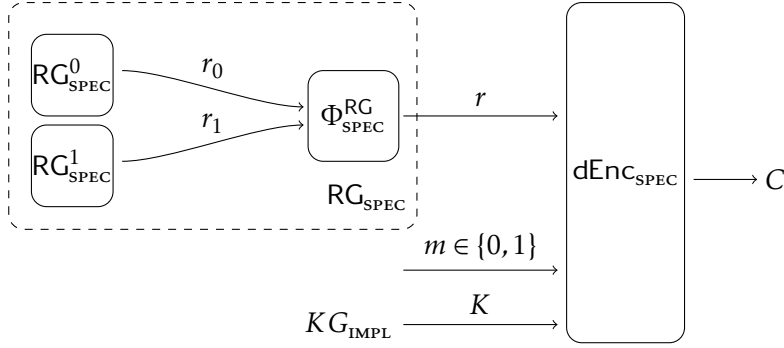


Figure 8: Stego-free specification for encryption algorithm in a symmetric-key bit encryption scheme.

Finally, consider the decryption algorithm. This algorithm does not influence CPA security, but directly influences correctness. Observe that the decryption algorithm is deterministic and can be analyzed with a public input distribution as we focus on bit encryption. To see this, the decryption algorithm will input a key  $K$  (generated by  $\text{KG}_{\text{IMPL}}$ ) and a ciphertext  $C$  (generated by encrypting 0 or 1 using  $K$ ). The watchdog can sample from the input distribution of  $\text{Dec}_{\text{IMPL}}$  to check the correctness, i.e., the consistency with  $\text{Dec}_{\text{SPEC}}$ .

With all of the algorithms handled individually as above, we present the first general immunizing result for randomized encryption algorithms.

**Theorem 4.2.** *Given any stateless IND-CPA secure symmetric bit encryption scheme, the specification described above is subversion resistant and correct according to Def. 4.1.*

(*sketch*). The specification is described above (with the randomized algorithm split into randomness generation and a deterministic part, and the randomness generation split into two components together with an immunizing function as in Figure 4.).

Security. The watchdog is the combination of those described above: it guarantees that  $\text{RG}_{\text{SPEC}}$  is stego-free (making samples to observe collisions as in the proof of Theorem 3.1) and guarantees that  $\text{dEnc}_{\text{IMPL}}$  is consistent with the specification on inputs sampled from  $\text{KG}_{\text{IMPL}} \times \text{RG}_{\text{IMPL}}$  (with 0 and 1); cf. Lemma 2.3. Here we only sketch the game changes and explain the negligible differences arising during the game transitions conditioned the watchdog's result.

Game-0 is the original game as described in Figure 6 with a trusted amalgamation; Game-1 is the same as Game-0 except  $\text{KG}_{\text{IMPL}}$  is replaced with  $\text{KG}_{\text{SPEC}}$ ; Game-2 is the same as Game-1 except  $\text{Enc}_{\text{IMPL}}$  is replaced with  $\text{Enc}_{\text{SPEC}}$ .

The adversary's advantage in Game-0 and Game-1 are negligibly close because of the stego-freeness of  $\text{KG}_{\text{SPEC}}$ . To see this, a simulator can simulate the rest of the games in Game-0 (or Game-1) after receiving  $K$  sent by the challenger in the game defining stego-freeness for  $\text{KG}_{\text{SPEC}}$ . If one can distinguish Game-0 from Game-1, then the simulator can easily tell apart  $\text{KG}_{\text{IMPL}}$  from  $\text{KG}_{\text{SPEC}}$

(even if the challenger does the amalgamation). Similarly, we can argue Game-1 and Game-2 are indistinguishable because  $\text{Enc}_{\text{SPEC}}$  is stego-free. The fact that  $\text{Enc}_{\text{SPEC}}$  is stego-free follows from Corollary 3.4.

Of course, in Game-2 every algorithm used in the game is faithfully implemented; thus the IND-CPA security of the underlying encryption scheme would hold.

Regarding correctness. As described above, the decryption algorithm for a bit encryption scheme has a public input distribution. The watchdog can sample random keys and check whether  $\text{Dec}_{\text{IMPL}}$  works properly for bit encryption. Following Lemma 2.3, we conclude that:  $\forall b \in \{0, 1\}$ ,

$$\Pr[\text{Dec}_{\text{IMPL}}(K, \text{Enc}_{\text{SPEC}}(K, b)) \neq b : K \leftarrow \text{KG}_{\text{IMPL}}] \leq \text{negl}(\lambda),$$

while for the encryption algorithm,  $\text{Enc}_{\text{IMPL}}$  can be used interchangeably with  $\text{Enc}_{\text{SPEC}}$  assuming trusted amalgamation. Thus,  $\forall b \in \{0, 1\}$ ,

$$\Pr[\text{Dec}_{\text{IMPL}}(\text{Enc}_{\text{IMPL}}(K, b)) \neq b : K \leftarrow \text{KG}_{\text{IMPL}}] \leq \text{negl}(\lambda).$$

Combining these yields the statement of the theorem. □

**Subversion resistant symmetric encryption with large message spaces. (Also handling state).**

For large message spaces, the security game must be adapted to allow the adversary to query; see Figure 13 in appendix. A. As mentioned in Section 2, this immediately invites an input-trigger attack. Specifically, for a particular query  $m_i$  (chosen randomly by  $A$  during subversion), the subverted encryption algorithm may directly output the secret key; the same threat exists for stateful *bit* encryption, where a particular sequence of encryption bits may act as a trigger. Note that this simply cannot be detected by a PPT watchdog (making polynomially many samples). Furthermore, the same attack can be launched on  $\text{Dec}_{\text{IMPL}}$  to ruin correctness (as  $\text{Dec}_{\text{IMPL}}$  can output a different value). This suggests the principle that a subverted implementation should never be given a “naked” input generated by the adversary, e.g., the queried message. (Such an input can act as a direct trigger when drawn from a large space, or can be remembered by stateful subversions to effectively simulate a large space.)

However, observe that if the input message can be forced to come from a known distribution  $U$  (e.g., the uniform distribution) the watchdog can check consistency between  $\text{dEnc}_{\text{IMPL}}$  and  $\text{dEnc}_{\text{SPEC}}$  on  $U$  (we ignore the other inputs here for simplicity, since they are either fixed or from a public input distribution). Indeed, the watchdog can guarantee that with overwhelming probability  $\text{dEnc}_{\text{IMPL}}$  is as good as  $\text{dEnc}_{\text{SPEC}}$  when the input is sampled according to  $U$ . In this case, the watchdog can in fact test the behavior of the algorithm on a *sequence* of such inputs from  $U$ , which is how we handle stateful subversion. Now, to defend against an input-trigger attack, we must ensure that the probability that any *particular*  $m$  causes an inconsistency is negligible. This sort of “worst-case vs. average-case” relationship arises in the theory of *self-correcting programs* [Rub91], where one wants to transform a program with a negligible fraction of bad inputs into a negligible error probability for every input. With these observations, we return to the large-message space challenge.

**Constructions.** First, we consider straightforward bit-wise encryption. When encrypting a message  $m = m_1 \dots m_\ell$  for  $\ell = |m|$ , the user generates the ciphertext by calling  $\text{Enc}_{\text{IMPL}}$   $\ell$  times, yielding  $C := (c_1, \dots, c_\ell)$ , where  $c_i = \text{Enc}_{\text{IMPL}}(K, m_i)$ . Since the bit encryption scheme we developed above preserves IND-CPA security; IND-CPA security follows via a simple hybrid argument. (Note, however, that in this case it is important that the encryption algorithm is stateless.)



To develop a more efficient scheme that can further handle long message (also stateful encryption) we augment the model by permitting the user (i.e., the challenger) to carry out *one trusted addition* operation for each encryption/decryption, see Fig. 9. (We continue to assume trusted amalgamation, as before.) We augment the specification of the encryption algorithm with a random message generator  $\text{MG}_{\text{SPEC}}$ . Specifically, the specification of the encryption algorithm  $\text{Enc}_{\text{SPEC}}$  has the form  $(\text{RG}_{\text{SPEC}}, \text{dEnc}_{\text{SPEC}}, \text{MG}_{\text{SPEC}})$ , where  $\text{MG}_{\text{SPEC}}$  has the specification  $(\text{MG}_{\text{SPEC}}^0, \text{MG}_{\text{SPEC}}^1, \Phi_{\text{SPEC}}^{\text{MG}})$ , and  $\text{RG}_{\text{SPEC}}$  is as before (as in Fig. 4). When encrypting a message  $m$ , the user first runs  $\text{MG}_{\text{IMPL}}$  (the implementation) to sample a random message  $u$ , and computes  $m' = m \oplus u$ . The user will call  $\text{Enc}_{\text{IMPL}}$  to encrypt the new message  $m'$  and obtains the ciphertext  $c'$ . This includes calls to  $\text{KG}_{\text{IMPL}}, \text{RG}_{\text{IMPL}}$  and passing the corresponding outputs  $K, r$  together with  $m'$  to  $\text{dEnc}_{\text{IMPL}}$ ; see Figure 9. Observe that  $m'$  is a *uniformly chosen message* (as the watchdog can ensure that  $u$  is safely generated). The new ciphertext  $C$  now includes  $u$  together with the ciphertext  $c'$ . For decryption, the user first runs  $\text{Dec}_{\text{IMPL}}$  on  $c'$  and obtains  $m'$ ; then the user computes  $m = m' \oplus u$ .

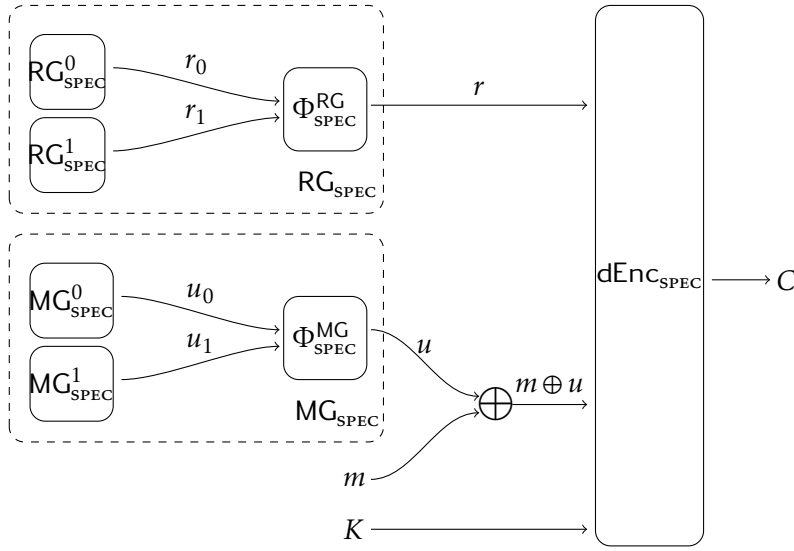


Figure 9: Stego-free encryption specification supporting large messages, where  $K \leftarrow \text{KG}_{\text{IMPL}}$ .

Security analysis. The intuition that this simple method works is as follows: First, we generalize Theorem 4.2 that symmetric-key encryption for *random messages* are also subversion resistant. To see this, we analyze the stego-freeness algorithm by algorithm. The key generation is the same as the bit encryption. The encryption algorithm now takes input a uniform message, together with the key and security parameter. It means that the encryption algorithm (the deterministic part  $\text{dEnc}_{\text{SPEC}}$ ) now takes inputs from public input distributions, i.e.,  $\text{KG}_{\text{IMPL}} \times \text{RG}_{\text{IMPL}} \times \mathcal{U}_{\mathcal{M}}$ , where  $\mathcal{U}_{\mathcal{M}}$  is the uniform distribution over message space  $\mathcal{M}$ . Following Theorem 3.3, such encryption algorithm will be stego-free as long as the specification is designed as Figure 5. Now for the decryption algorithm, since the encryption is for uniform messages, thus the decryption algorithm now also takes a public input distribution. Next, we show the encryption specification defined in Figure 9 indeed takes uniform messages as input. (1.) Following Theorem 3.1, the uniform message sampler  $\text{MG}_{\text{SPEC}}$  is stego-free. (2.) With the trusted addition operation, when a that  $\text{Enc}_{\text{IMPL}}$  takes as input will be  $m' = m \oplus u$ , where  $u$  looks uniform even to the adversary, thus  $m'$  would look uniform to  $\text{Enc}_{\text{IMPL}}$  (actually the deterministic part  $\text{dEnc}_{\text{IMPL}}$ ). Similar to the analysis of Theorem 4.2, we can



show a stronger result that handles the correctness and subversion resistance for symmetric-key encryption supporting long messages. The full stateful case was handled in detail in remark 3.2.

**Theorem 4.3.** *For any IND-CPA secure symmetric-key encryption, the specification described as above is subversion resistant and correct according to Def. A.3, assuming a trusted  $\oplus$  operation and amalgamation.*

## 4.2 Public key encryption preserving IND-CPA security under subversion

Now we turn to public-key encryption, which follows fairly directly from the previous construction. The major difference arises with key generation, as asymmetric key generation has to be treated with more care than simple randomness generation; see Figure 10, which indicates the construction. Specifically, the basic techniques used for symmetric key encryption above can be adapted for public key encryption. Key generation must be considered as a randomized algorithm producing structured output (with only security parameter as input). With these tools at hand, we resolve the major open problem to construct a IND-CPA secure PKE when facing subversions that we asked at the beginning of the paper. We remark that the assumption of “trusted  $\oplus$ ” in the theorem below can be removed if the message space is small.

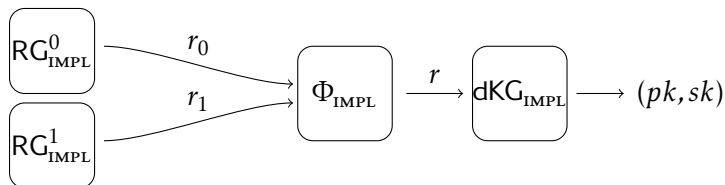


Figure 10: Stego-free asymmetric key generation.

**Corollary 4.4.** *For any IND-CPA secure public key bit encryption, there exists a specification design such that it is subversion resistant and correct with trust amalgamation. Furthermore, for any IND-CPA public key encryption (supporting large input space), there exists a specification design such that it is subversion resistant and correct according to Def. A.4, if the user can do a trusted  $\oplus$  and amalgamation.*

## 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.

- [Bou05] Jean Bourgain. More on the sum-product phenomenon in prime fields and its applications. *International Journal of Number Theory*, 1:1–32, 2005.
- [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.
- [CG88] Benny Chor and Oded Goldreich. Unbiased bits from sources of weak randomness and probabilistic communication complexity. *SIAM J. Comput.*, 17(2):230–261, 1988.
- [CZ15] Eshan Chattopadhyay and David Zuckerman. Explicit two-source extractors and resilient functions. Technical report, TR15-119, Electronic Colloquium on Computational Complexity, 2015.
- [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>.
- [Doc13] Docker.Inc. Docker, 2013. <https://www.docker.com/>.
- [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.
- [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.
- [NZ96] Noam Nisan and David Zuckerman. Randomness is linear in space. *Journal of Computer and System Sciences*, 52(1):43–52, 1996.
- [Rog15] Phillip Rogaway. The moral character of cryptographic work. *Cryptology ePrint Archive*, Report 2015/1162, 2015. <http://eprint.iacr.org/2015/1162>.
- [RTYZ15] Alexander Russell, Qiang Tang, Moti Yung, and Hong-Sheng Zhou. Cliptography: Clipping the power of kleptographic attacks. *Cryptology ePrint Archive*, Report 2015/695, 2015. <http://eprint.iacr.org/2015/695>.
- [Rub91] Ronitt A. Rubinfeld. *A Mathematical Theory of Self-checking, Self-testing and Self-correcting Programs*. PhD thesis, University of California at Berkeley, Berkeley, CA, USA, 1991. UMI Order No. GAX91-26752.

- [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.
- [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 Omitted Definitions

### A.1 Stego-free randomness generation.

One of the most fundamental components of a randomized algorithm is the randomness generation. In particular, the devastating attacks using subliminal channels are exploring the subverted randomness generation. Moreover, if we have clean randomness, many tasks (including the CPA-secure public key encryption) which are previously considered impossible may become reachable. Formally,

**Definition A.1.** For any randomness generation algorithm  $\text{RG}$ , consider a specification  $\text{RG}_{\text{SPEC}} := (\text{RG}_{\text{SPEC}}^0, \text{RG}_{\text{SPEC}}^1, \Phi_{\text{SPEC}})$ . We say such specification  $\text{RG}_{\text{SPEC}}$  is stego-free with a trusted amalgamation in the offline watchdog model, if there exists an offline PPT watchdog  $\mathcal{W}$ , for any PPT adversary  $\mathcal{A}$  playing the following game (see Figure 11) with challenger  $\mathcal{C}$ , such that, at least one of the following conditions hold:

$\text{Det}_{\mathcal{W}, \mathcal{A}}$  is non-negligible, or,  $\text{Adv}_{\mathcal{A}}$  is negligible,

where  $\text{Adv}_{\mathcal{A}}(1^\lambda) = \left| \Pr[b_{\mathcal{C}} = 1] - \frac{1}{2} \right|$ , and  $\text{Det}_{\mathcal{W}, \mathcal{A}}(1^\lambda) = \left| \Pr[\mathcal{W}^{\text{RG}_{\text{IMPL}}}(1^\lambda) = 1] - \Pr[\mathcal{W}^{\text{RG}_{\text{SPEC}}}(1^\lambda) = 1] \right|$ ,

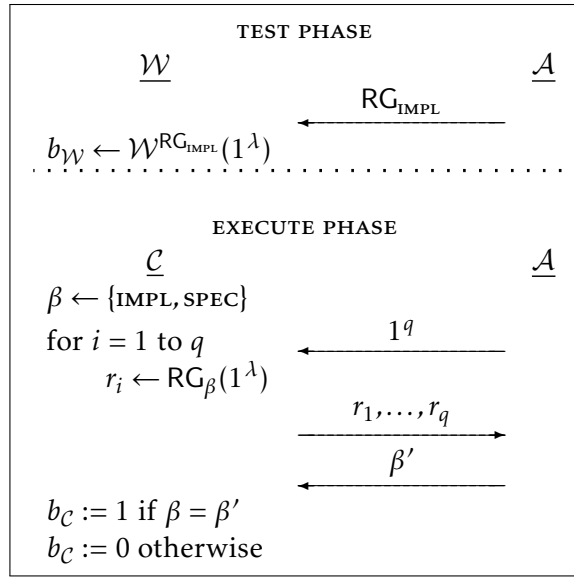


Figure 11: Game for stego-free randomness generation with trusted amalgamation.

## A.2 Stego-free randomized algorithms with public input distributions

**Definition A.2.** For a randomize algorithm  $G$ , consider a specification  $G_{\text{SPEC}} := (\text{RG}_{\text{SPEC}}, \text{dG}_{\text{SPEC}})$ , where  $\text{RG}_{\text{SPEC}} := (\text{RG}_{\text{SPEC}}^0, \text{RG}_{\text{SPEC}}^1, \Phi_{\text{SPEC}})$ , and  $\text{dG}_{\text{SPEC}}$  is deterministic. We say such specification  $G_{\text{SPEC}}$  is stego-free with a trusted amalgamation in the offline watchdog model, if there exists an offline PPT watchdog  $\mathcal{W}$ , for any PPT adversary  $\mathcal{A}$  playing the following game (see Figure 12) with challenger  $\mathcal{C}$ , such that, either,

$\text{Adv}_{\mathcal{A}}$  is negligible, or,  $\text{Det}_{\mathcal{W}, \mathcal{A}}$  is non-negligible.

where  $\text{Adv}_{\mathcal{A}}(1^\lambda) = \left| \Pr[b_{\mathcal{C}} = 1] - \frac{1}{2} \right|$ , and,  $\text{Det}_{\mathcal{W}, \mathcal{A}}(1^\lambda) = \left| \Pr[\mathcal{W}^{\text{G}_{\text{IMPL}}, \text{IG}}(1^\lambda) = 1] - \Pr[\mathcal{W}^{\text{G}_{\text{SPEC}}}(1^\lambda) = 1] \right|$ ,

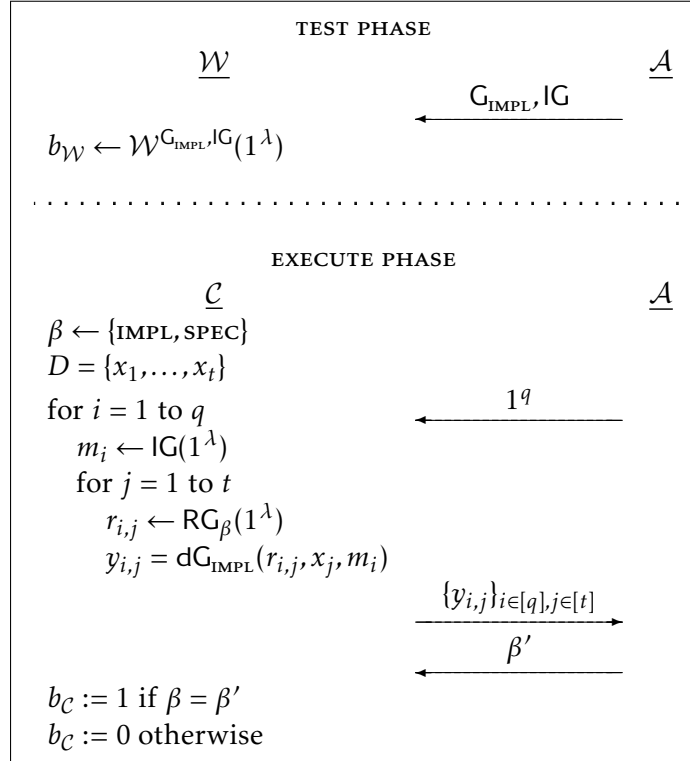


Figure 12: Game for stego-free randomized algorithm with trusted amalgamation. IG is the instance generator, and  $D$  is a poly-size input domain.

### A.3 Subversion resistant symmetric-key encryption with long messages

**Definition A.3.** For any randomized symmetric-key encryption scheme, consider a specification  $\mathbf{E}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}, \text{Enc}_{\text{SPEC}}, \text{Dec}_{\text{SPEC}})$ . We say specification  $\mathbf{E}_{\text{SPEC}}$  is subversion resistant with a trusted addition and amalgamation in the offline watchdog model, if there exists an offline PPT watchdog  $\mathcal{W}$ , for any PPT adversary  $\mathcal{A}$ , playing the game defined in Figure 13, either,

$\text{Adv}_{\mathcal{A}}$  is negligible, or,  $\text{Det}_{\mathcal{W}, \mathcal{A}}$  is non-negligible.

where  $\text{Adv}_{\mathcal{A}}(1^\lambda) = |\Pr[b_C = 1] - \frac{1}{2}|$ , and  $\text{Det}_{\mathcal{W}, \mathcal{A}}(1^\lambda) = |\Pr[\mathcal{W}^{\text{E}_{\text{IMPL}}}(1^\lambda) = 1] - \Pr[\mathcal{W}^{\text{E}_{\text{SPEC}}}(1^\lambda) = 1]|$ ,

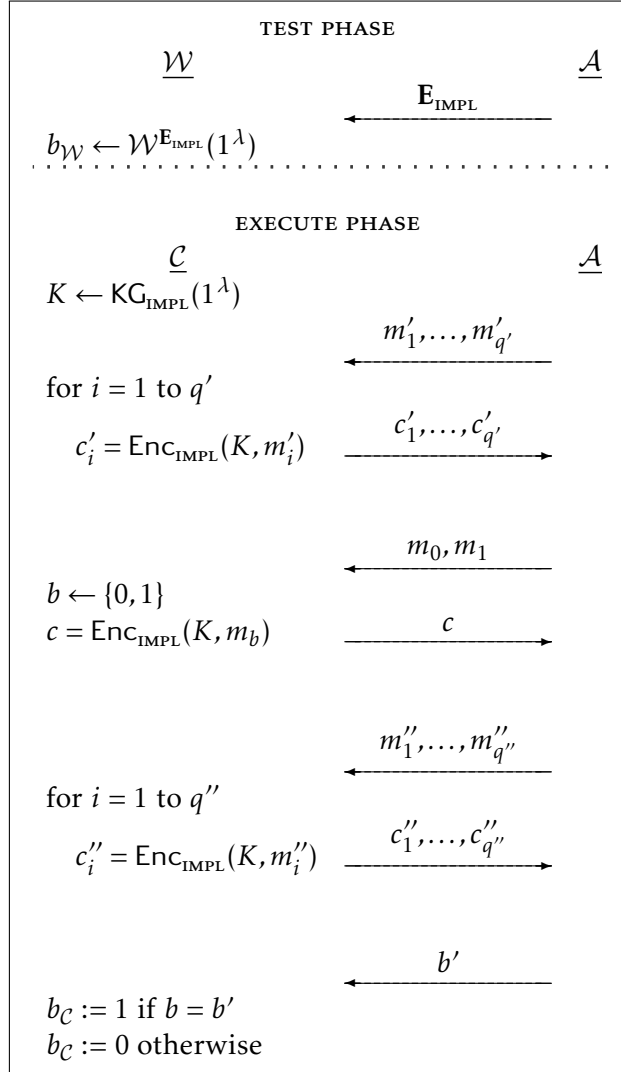


Figure 13: Subversion-resistant symmetric-key encryption with trusted addition.

#### A.4 Subversion resistant public key encryption with long messages

**Definition A.4.** For any public key encryption scheme, consider a specification  $\mathbf{E}_{\text{SPEC}} := (\text{KG}_{\text{SPEC}}, \text{Enc}_{\text{SPEC}}, \text{Dec}_{\text{SPEC}})$ . We say specification  $\mathbf{E}_{\text{SPEC}}$  is subversion resistant with a trusted addition and amalgamation in the offline watchdog model, if there exists an offline PPT watchdog  $\mathcal{W}$ , for any PPT adversary  $\mathcal{A}$ , such that, at least one of the following conditions hold in the game defined in Figure 14:

**Det** $_{\mathcal{W}, \mathcal{A}}$  is non-negligible, or, **Adv** $_{\mathcal{A}}$  is negligible,

where  $\mathbf{Adv}_{\mathcal{A}}(1^\lambda) = |\Pr[b_C = 1] - \frac{1}{2}|$ , and  $\mathbf{Det}_{\mathcal{W}, \mathcal{A}}(1^\lambda) = |\Pr[\mathcal{W}^{\mathbf{E}_{\text{IMPL}}}(1^\lambda) = 1] - \Pr[\mathcal{W}^{\mathbf{E}_{\text{SPEC}}}(1^\lambda) = 1]|$ .

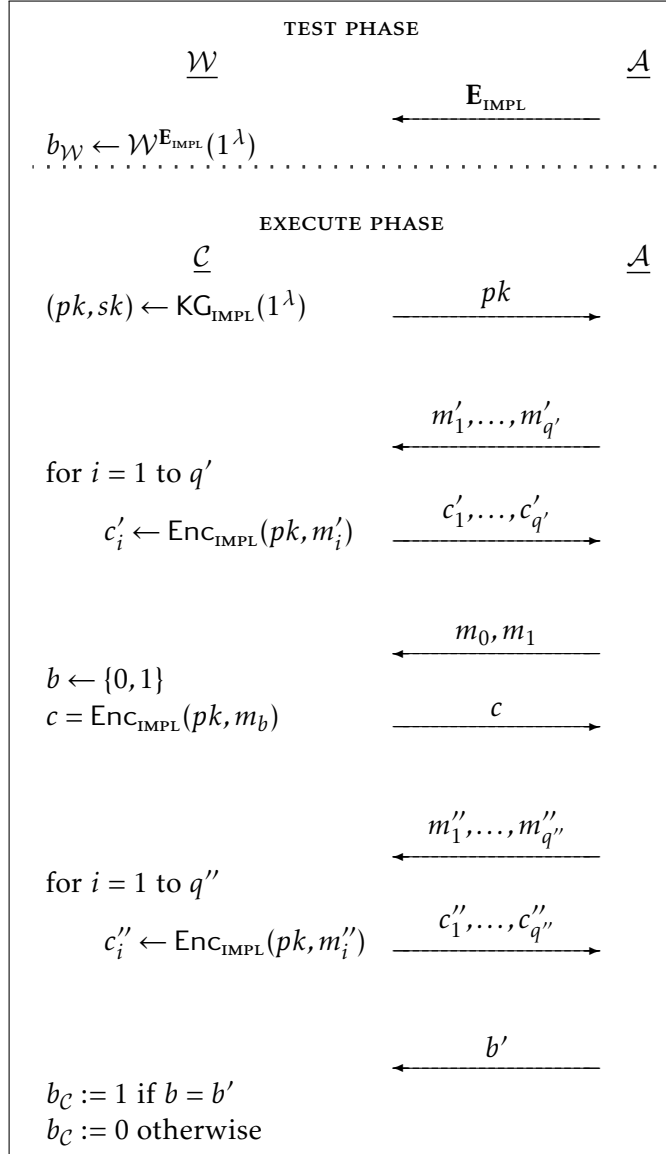


Figure 14: Subversion-resistant public key encryption supporting large messages.



## B An attack on single-source randomness primitives

**Subverted randomness generation attack:** In the following attack on public key encryption, the adversary honestly implements the key generation and decryption, and only subverts the encryption algorithm. Suppose the specification of the (public-key) encryption algorithm is defined as  $\text{Enc}_{\text{SPEC}} := (\text{RG}_{\text{SPEC}}, \text{dEnc}_{\text{SPEC}}, \Phi_{\text{SPEC}})$ . The meaning of each component is self-evident:  $\text{RG}_{\text{SPEC}}$  generates uniformly random bits  $r_0$ , the function  $\Phi_{\text{SPEC}}$  “cleans”  $r_0$  to produce the final random bits  $r$ , and, finally,  $\text{dEnc}_{\text{SPEC}}$  takes the random bits  $r$ , the encryption key  $pk$ , and the message bit as inputs and produces a ciphertext  $C$ .

*The attack:* The adversary  $\mathcal{A}$  first randomly chooses a backdoor  $z$ , and prepares the subverted implementation  $\text{Enc}_{\text{IMPL}}^z := (\text{RG}_{\text{IMPL}}, \text{dEnc}_{\text{IMPL}}, \Phi_{\text{IMPL}})$  with the backdoor embedded. In particular,  $\text{RG}_{\text{IMPL}} := \text{RG}_{\text{IMPL}}^z$  carries out rejection sampling to ensure that the ciphertext encrypting 0 and the ciphertext encrypting 1 can be distinguished by applying a PRF (using  $z$  as the key); the algorithms  $\text{dEnc}_{\text{IMPL}}$  and  $\Phi_{\text{IMPL}}$  are honestly implemented (that is, identical to the specifications). Later the adversary  $\mathcal{A}$  can easily learn secret information (indeed, the plaintext) from the ciphertext generated by the subverted algorithms by applying the PRF (using her backdoor  $z$  as the key). See Figure 15 for detailed description.

```

RGIMPLz(pk):
-----
Repeat:
  r0 ← RGSPEC
  c0 = dEncSPEC(pk, 0; ΦSPEC(r0))
  c1 = dEncSPEC(pk, 1; ΦSPEC(r0))
Until: PRF(z, c0) = 0 ∧ PRF(z, c1) = 1
Return r0

A(z, C):
-----
b = PRF(z, C)
Return b

```

Figure 15: Subverted randomness generation and the message recovery algorithms

**Security analysis.** Due to the rejection sampling condition, it is easy to see that the adversary defined in Figure 15 can determine the plaintext bit perfectly from the ciphertext. As for the *detection probability*, the randomness output by  $\text{RG}_{\text{SPEC}}$  is a uniform  $\lambda$ -bit string; in contrast, the randomness output by  $\text{RG}_{\text{IMPL}}$  is a string selected uniformly from a (random) subset  $S$  of  $\{0, 1\}^\lambda$  (determined by the PRF). The subset  $S$  consists of all strings that carry 0 and 1 to ciphertexts satisfying a criterion given by the PRF. Let us think of the PRF as a random function, that means the rejection sampling condition will be satisfied with probability  $1/4$  for each  $r_0$  uniformly sampled. Essentially, we can consider  $S$  as a random subset of  $\{0, 1\}^\lambda$  with (expected) size  $2^{\ell-2}$ . If there is no collision when  $\mathcal{W}$  asks  $q$  queries, then the  $q$  different bit strings observed by  $\mathcal{W}$  can come from either of the whole space  $\{0, 1\}^\lambda$  or a subset  $S$  (with size larger than  $q$ ). This means conditioned on no collision, no watchdog can tell apart  $\text{RG}_{\text{IMPL}}$  from  $\text{RG}_{\text{SPEC}}$ . Using the total probability theorem, we can bound the distinguishing advantage by the probability that a collision appears in  $q$  queries when sample from  $S$ . We leave the full proof to Sec. C in the appendix.

**Proposition B.1.** *For any CPA-secure (public-key) encryption scheme, for any public function  $\Phi$ , the adversary  $\mathcal{A}$  shown in Figure 15 can learn the plaintext with probability 1 given the ciphertext generated by  $\text{Enc}_{\text{IMPL}}^z$  even if the randomness generation is separated and immunized by a random oracle. Furthermore, suppose  $\text{RG}_{\text{IMPL}}$  outputs  $\ell$  bits of randomness; then the detection advantage is  $q^2/2^{\ell-4} + \text{negl}(\lambda)$  for all PPT watchdogs making  $q$  queries during the interrogation, assuming PRF is a secure pseudorandom function.*

*Proof.* The effectiveness of  $\mathcal{A}$  is straightforward, since the randomness generated by  $\text{RG}_{\text{IMPL}}$  makes the ciphertext to be distinguishable using the PRF, thus the adversary  $\mathcal{A}$  who knows the backdoor recovers the message bit perfectly.

Next, we will argue that no offline watchdog can notice the subversion, particularly  $\text{RG}_{\text{IMPL}}$ . (All other components are honestly implemented). We define the game sequence as follows:

Game-0 is the game that  $\mathcal{W}$  is interacting with  $\text{RG}_{\text{IMPL}}$ .

Game-1 is the same as Game-0 except that the PRF used in  $\text{RG}_{\text{IMPL}}$  is replaced with a random function  $R$ .

Game-2 is the same as Game-1 except that  $\text{RG}_{\text{IMPL}}$  resamples, if  $\mathcal{W}$  notices a collision in the responses to the  $q$  queries.

Game-3 is the same as Game-2 except that  $\text{RG}_{\text{IMPL}}$  is replaced with  $\text{RG}_{\text{SPEC}}$ .

Game-4 is the same as Game-3 except removing the resampling condition.

Analyzing the gaps, we see: Game-0 and Game-1 are indistinguishable because of the PRF security; Game-1 and Game-2 are identical, if there is no collision among the  $q$  queries in Game-2. In Game 2, for any  $r_0$ , the probability that  $R(c_0) = 0 \wedge R(c_1) = 1$  is  $1/4$ , where  $c_0, c_1$  are ciphertexts encrypting 0, 1 respectively using  $\Phi(r_0)$  as the coin. Suppose  $S$  is the set that contains all the values that satisfies  $R(c_0) = 0 \wedge R(c_1) = 1$ , then the expected size of  $S$ ,  $E[|S|] = 2^\ell/4 = 2^{\ell-2}$ . It follows that with negligible probability,  $|S| \leq \frac{E[|S|]}{2} = 2^{\ell-3}$ . Then the probability that there exists a collision among  $q$  samples from  $S$  would be bounded by  $q^2/|S| \leq q^2/2^{\ell-3}$ . Game-2 and Game-3 is identical, since the responses can either appear in a random subset or the whole space. Game-3 and Game-4 are identical except there is collision when sampling  $q$  uniform points. The probability of such collision exist is  $q^2/2^\ell$ . Combining them above, we have the proposition.  $\square$

**Theorem 3.3.** *For any randomized algorithm  $G$ , consider the specification  $G_{\text{SPEC}} := (\text{RG}_{\text{SPEC}}, \text{dG}_{\text{SPEC}})$ , where  $\text{RG}_{\text{SPEC}}$  generates  $\lambda = \lambda(n)$  bits of uniform randomness and  $\text{dG}_{\text{SPEC}}$  is deterministic. Let  $(\text{RG}_{\text{SPEC}}^0, \text{RG}_{\text{SPEC}}^1, \Phi_{\text{SPEC}})$  be the double-split specification of  $\text{RG}_{\text{SPEC}}$  as above. If (1)  $\text{RG}_{\text{SPEC}}^0(1^\lambda)$  and  $\text{RG}_{\text{SPEC}}^1(1^\lambda)$  output  $\lambda$  uniform bits; (2)  $\Phi_{\text{SPEC}}$  takes  $r_0 \circ r_1$  as input, and outputs  $r$  (so it maps strings of length  $2k$  to strings of length  $k$ ) (see Fig. 5), then  $G_{\text{SPEC}}$  is stego-free with a trusted amalgamation (according to Def. A.2 in appendix. A.2). Here  $\Phi_{\text{SPEC}}$  is modeled as a random oracle, and  $\Phi_{\text{IMPL}}^0, \Phi_{\text{IMPL}}^1$  are executed independently.*

*Proof.* The watchdogs will be a combination of the ones in Theorem 3.1 and Lemma 2.3 to guarantee unpredictability of  $\text{RG}_{\text{IMPL}}^b$  and the overwhelming consistency for deterministic algorithms with a public input distribution. Here  $\text{dG}_{\text{SPEC}}$  is a deterministic algorithm and the output of  $\text{RG}_{\text{SPEC}} \times \text{IG}$  would be the input distribution to  $\text{dG}_{\text{SPEC}}$ .

We here only describe briefly about the game changes.

Game-0 is the same as Figure 12, the adversary  $\mathcal{A}$  prepares every piece of the implementation, and the challenger simply calls them and passes the inputs to the next implementation as defined (doing the basic amalgamation).

Game-1 is the same as Game-0 except that the randomness  $r$  is uniformly sampled by  $\mathcal{C}$ .

Game-2 is the same as Game-1 except that  $dG_{\text{IMPL}}$  is replaced with  $dG_{\text{SPEC}}$ .

Note that in Game-0, it corresponds to  $b = 0$ , while in Game-2, every implementation of the algorithm (except input generation) is now the specification, it corresponds to  $b = 1$ .

From Theorem 3.1, with a trusted amalgamation, the output from the implementation  $RG_{\text{IMPL}} := (RG_{\text{IMPL}}^0, RG_{\text{IMPL}}^1, \Phi_{\text{IMPL}})$  is pseudorandom to the adversary  $\mathcal{A}$  who made  $RG_{\text{IMPL}}$ . Thus in the security game defined in Figure 12, we can let the challenger simply generate  $r$  uniformly to reach Game-1.

From Lemma 2.3,  $dG_{\text{SPEC}}$  will be a deterministic algorithm with a public input distribution, thus  $dG_{\text{IMPL}}$  would be consistent with  $dG_{\text{SPEC}}$  with an overwhelming probability when inputs are sampled according to  $RG_{\text{SPEC}} \times IG$ , thus Game-2 can be reached with only a negligible gap.

Once in Game-2, all components are specification.  $\square$

**Remark B.1.** *We remark that for this particular attack we assume that the implementation has access to a public key—this yields an intuitively natural attack against an encryption scheme which permits full recovery of the message. However, the basic structure of the attack can be adapted to randomized algorithms in generality.*

## C Omitted Proofs

**Theorem 3.3.** *For any randomized algorithm  $G$ , consider the specification  $G_{\text{SPEC}} := (\text{RG}_{\text{SPEC}}, \text{dG}_{\text{SPEC}})$ , where  $\text{RG}_{\text{SPEC}}$  generates  $\lambda = \lambda(n)$  bits of uniform randomness and  $\text{dG}_{\text{SPEC}}$  is deterministic. Let  $(\text{RG}_{\text{SPEC}}^0, \text{RG}_{\text{SPEC}}^1, \Phi_{\text{SPEC}})$  be the double-split specification of  $\text{RG}_{\text{SPEC}}$  as above. If (1)  $\text{RG}_{\text{SPEC}}^0(1^\lambda)$  and  $\text{RG}_{\text{SPEC}}^1(1^\lambda)$  output  $\lambda$  uniform bits; (2)  $\Phi_{\text{SPEC}}$  takes  $r_0 \circ r_1$  as input, and outputs  $r$  (so it maps strings of length  $2k$  to strings of length  $k$ ) (see Fig. 5), then  $G_{\text{SPEC}}$  is stego-free with a trusted amalgamation (according to Def. A.2 in appendix. A.2). Here  $\Phi_{\text{SPEC}}$  is modeled as a random oracle, and  $\Phi_{\text{IMPL}}^0, \Phi_{\text{IMPL}}^1$  are executed independently.*

*Proof.* The watchdogs will be a combination of the ones in Theorem 3.1 and Lemma 2.3 to guarantee unpredictability of  $\text{RG}_{\text{IMPL}}^b$  and the overwhelming consistency for deterministic algorithms with a public input distribution. Here  $\text{dG}_{\text{SPEC}}$  is a deterministic algorithm and the output of  $\text{RG}_{\text{SPEC}} \times \text{IG}$  would be the input distribution to  $\text{dG}_{\text{SPEC}}$ .

We here only describe briefly about the game changes.

Game-0 is the same as Figure 12, the adversary  $\mathcal{A}$  prepares every piece of the implementation, and the challenger simply calls them and passes the inputs to the next implementation as defined (doing the basic amalgamation).

Game-1 is the same as Game-0 except that the randomness  $r$  is uniformly sampled by  $\mathcal{C}$ .

Game-2 is the same as Game-1 except that  $\text{dG}_{\text{IMPL}}$  is replaced with  $\text{dG}_{\text{SPEC}}$ .

Note that in Game-0, it corresponds to  $b = 0$ , while in Game-2, every implementation of the algorithm (except input generation) is now the specification, it corresponds to  $b = 1$ .

From Theorem 3.1, with a trusted amalgamation, the output from the implementation  $\text{RG}_{\text{IMPL}} := (\text{RG}_{\text{IMPL}}^0, \text{RG}_{\text{IMPL}}^1, \Phi_{\text{IMPL}})$  is pseudorandom to the adversary  $\mathcal{A}$  who made  $\text{RG}_{\text{IMPL}}$ . Thus in the security game defined in Figure 12, we can let the challenger simply generate  $r$  uniformly to reach Game-1.

From Lemma 2.3,  $\text{dG}_{\text{SPEC}}$  will be a deterministic algorithm with a public input distribution, thus  $\text{dG}_{\text{IMPL}}$  would be consistent with  $\text{dG}_{\text{SPEC}}$  with an overwhelming probability when inputs are sampled according to  $\text{RG}_{\text{SPEC}} \times \text{IG}$ , thus Game-2 can be reached with only a negligible gap.

Once in Game-2, all components are specification.  $\square$

## D Purifying randomness in the standard model

The split-and-amalgamate paradigm proposed in Section 3.2 assumes the specification of the immunization function  $\Phi_{\text{SPEC}}$  to be modeled as a random oracle. Intuitively, it seems hard to remove this assumption when the randomness generation is split into only two segments, since a PPT offline watchdog can guarantee at most a  $\omega(\log \lambda)$  entropy in the output of each  $\text{RG}_{\text{IMPL}}^b$  for  $b = 0, 1$ , (and the concrete lower bound of entropy is not even clear). Nevertheless, it is interesting to consider immunization in the standard model if we are willing to have more than constant number of components.

**A simple approach via multi-splitting.** Observe that suppose the randomness generation implemented by the adversary only outputs one bit, the watchdog (who checks whether there is a bit appears significantly more frequently than the other after drawing certain amount of samples) can guarantee that the bit only has a negligible bias. If not, suppose the bias for outputting bit value 1 is  $1/2 + \delta$  for a noticeable  $\delta$ , then according to the Chernoff bound, after the watchdog draws enough samples, he will notice that there are significantly more 1s appear which is abnormal if the specification is used which outputs a uniform bit.

With the above observation, we can simply extract one bit from a  $\text{RG}_{\text{IMPL}}^1$  from each draw. If we now split the randomness generation into  $n$  copies  $\text{RG}_{\text{IMPL}}^1, \dots, \text{RG}_{\text{IMPL}}^n$ , we can collect  $n$ -bits if we independently draw one bit from each of them.

Security analysis. We now analyze how these  $n$ -bits are distributed. Suppose the output distribution of  $\text{RG}_{\text{IMPL}}^i$  is denoted by  $D_i$ , the user first samples  $d_i \leftarrow D_i$ , and outputs only the first bit of  $d_i$ ; each such bit has bias at most  $1/2 + \epsilon$ , for a negligible function  $\epsilon$ . We denote each bit using the random variable  $X_1, \dots, X_n$ . It follows that for any particular value  $b_1, \dots, b_n \in \{0, 1\}^n$ , as  $1 + \frac{x}{n} \leq e^x$  we have

$$\Pr[(X_1, \dots, X_n) = (b_1, \dots, b_n)] \leq (1/2 + \epsilon)^n \leq \frac{1}{2^n} e^{2n\epsilon},$$

Let  $U_n$  be the uniform distribution over  $\{0, 1\}^n$ , and  $Y_n$  be the distribution of  $X_1 \dots X_n$ . Recall that if  $Y_n(\bar{x}) \leq U_n(\bar{x})(1 + t)$  for all  $\bar{x} \in \{0, 1\}^n$ , the statistical distance satisfies:

$$\|Y_n - U_n\| = \frac{1}{2} \sum_{\bar{x}} |Y_n(\bar{x}) - U_n(\bar{x})| \leq \frac{1}{2} \sum_{\bar{x}} U_n(\bar{x}) \cdot t = \frac{t}{2}.$$

Putting back the value  $t = e^{2n\epsilon} - 1$ , we conclude that the statistical distance between the distribution of  $X_1 \dots X_n$  and the uniform distribution over  $n$ -bit strings, is no more than  $2n\epsilon \leq \text{negl}(\lambda)$ , the inequality holds because  $e^x - 1 \leq 2x$  for  $x < 1/2$ . Summarizing above, we have:

**Proposition D.1.** *Suppose the randomness generation specification  $\text{RG}_{\text{SPEC}}$  is defined as  $(\text{RG}_{\text{SPEC}}^1, \dots, \text{RG}_{\text{SPEC}}^n)$ , where each  $\text{RG}_{\text{SPEC}}^i$  is supposed to output uniform  $n$  bits, and  $\Phi_{\text{SPEC}} : \{0, 1\}^{n^2} \rightarrow \{0, 1\}^n$  simply outputs the first bit of each  $n$ -bit block. Such  $\text{RG}_{\text{SPEC}}$  is stego-free with a trusted amalgamation.*

**Remark D.1.** *We keep the immunization function for consistency with the general model. In this simple multi-splitting, we may not even need  $\Phi_{\text{SPEC}}$ ; we can define each  $\text{RG}_{\text{SPEC}}^i$  to output only 1 random bit.*

**A more efficient construction.** The above construction is extremely simple, but with a price that the randomness generation has to be split into  $n$  pieces if the output is  $n$ -bit long. Note that there are at least  $c \log n$  bits entropy (for any constant  $c \geq 2$ ) in the output of each  $\text{RG}_{\text{IMPL}}^i$ . This comes from the guarantee of a watchdog that makes  $O(n^{2c})$  queries during the interrogation. Moreover, the bitstring can be stretched (if it is close to uniform) to a polynomially longer output using a PRG.<sup>8</sup>

We may use more powerful machineries of a strong randomness extractor together with a pseudorandom generator. The intuition is to first get a short seed and then use the power of the extractor to bootstrap the randomness generation by a factor of  $\log n$  since we can get at least this many uniform bits from the output of each  $\text{RG}_{\text{IMPL}}^i$ . After collecting  $O(n^\epsilon)$  bits (for some small  $\epsilon$ ), we can apply the PRG to stretch it to  $O(n)$  bits. Suppose  $\text{RG}_{\text{SPEC}}$  now is defined as  $(\text{RG}_{\text{SPEC}}^1, \dots, \text{RG}_{\text{SPEC}}^\ell, \Phi_{\text{SPEC}})$ , where  $\ell = n^\epsilon / \log n + \log n$  for some small  $\epsilon$ . We describe the immunization function  $\Phi_{\text{SPEC}} : \{0, 1\}^{n \cdot \ell} \rightarrow \{0, 1\}^n$  below, when inputting  $r_1, \dots, r_\ell$  (generated by  $\text{RG}_{\text{IMPL}}^i$  respectively).

- First,  $\Phi_{\text{SPEC}}$  takes  $r_1, \dots, r_{\log n}$  and returns the first bit of each  $r_i$  and obtains a  $\log n$ -bit uniform string denoted by  $s_0$ .

<sup>8</sup>We remark that here the PRG is just one deterministic algorithm, and it is associated with a public input distribution. Such PRG can be based on e.g., some *fixed* hard function. This is in contrast with the (backdoored) PRG that requires some randomly generated public parameters  $PP$ , and has the danger of containing backdoors in  $PP$  [DGG<sup>+</sup>15, RTYZ15].

- Then,  $\Phi_{\text{SPEC}}$  applies a seeded strong extractor  $\text{Ext}$  to  $r_{\log n+1}, \dots, r_\ell$  respectively using  $s_0$  as a seed. It obtains  $\ell \log n = n^\epsilon$  uniform bits, ( $\log n$  bits from each of them), which are denoted by  $s_1$ .
- Last, we apply a pseudorandom generator PRG on  $s_1$  to stretch it to  $n$  bits, and output those as the final random coin.

Security analysis. Let us first assume  $\Phi_{\text{SPEC}}$  is not subverted. In the first step,  $s_0$  would be close to uniform due to Proposition D.1. In the second step, since  $\text{Ext}$  is a strong extractor, the output of  $\text{Ext}(s_0, r_i)$  is close to uniform and independent of  $s_0$ , via the simple union bound,  $s_1$  is also negligibly close to uniform. Last step follows easily because of the PRG property.

Next, observe that although  $\Phi_{\text{SPEC}}$  here is a complex function, it is still a deterministic algorithm, and it is with a public input distribution. Following Lemma 2.3, the watchdog can guarantee that the subverted implementation  $\Phi_{\text{IMPL}}$  will be consistent with the specification with an overwhelming probability when the input is sampled from  $\text{RG}_{\text{IMPL}}^1 \times \dots \times \text{RG}_{\text{IMPL}}^\ell$ .

**Proposition D.2.** *There exists a specification for the randomness generation that outputs  $n$  bits is stego-free with the trusted amalgamation and  $O(n^\epsilon/\log n)$  segments for any constant  $\epsilon$ .*

## E Further Applications

We show above how we can salvage randomized algorithms, in particular randomized encryption schemes, when they are under subversion. While the technique we develop is used to generically eliminate subliminal channels which has its own interest. In this section, we give two more examples to illustrate how the techniques can be applied to broader settings. We first show how to circumvent the impossibility result of [DGG<sup>+</sup>15] for publicly immunizing the *output* of a potentially subverted PRG (called backdoored PRG). Then we show how our decomposition-then-amalgamate paradigm can be used to simulate an intuitively more powerful semi-private model proposed in [DGG<sup>+</sup>15]. In such a model, there is a uniformly random string  $r$  generated by a trusted party and  $r$  is made public only after the adversary provides the implementations. This will motivate follow-up works to consider the power of the semi-private model, then realize it using our paradigm.

**Public immunization of output of a backdoored PRG.** In a backdoored PRG, the stretch algorithm takes a random seed and also a public key (generated by the adversary) as inputs. The goal of the adversary is to ensure that the output of the PRG looks uniform to everyone except herself by providing a malicious public key. In particular, it captured the notorious subversion attack on the Dual\_EC PRG. Unfortunately, as shown in [DGG<sup>+</sup>15], publicly immunizing the outputs of the backdoored PRG using a random oracle is impossible, i.e., the adversary can still distinguish the immunized string from a random string. The reason is similar to the attack shown in Section 3.1 since the subverted stretch algorithm can make random oracle queries during the execution. Russell et al. proposed a way to sidestep this impossibility by randomizing the public key [RTYZ15] with the caveat that the specification of the key generation should output uniform group elements. As pointed out by [DGG<sup>+</sup>15], immunizing the output of the PRG is preferable as one does not have to pay particular attention to the key/parameter generation.

It is not hard to see our split-and-amalgamate paradigm can be applied here to immunize the *output* of a backdoored PRG. We present a simplified version that the PRG is stateless, it is not hard generalize it to allow iterations to capture the Dual\_EC PRG, and we defer that to the full version.

Moreover, if we allow to split into more pieces, we can immunize the output in the standard model. Even in the semi-private model with a trusted randomness, the known result uses either random oracle or UCE, neither of which can be explicitly instantiated so far.

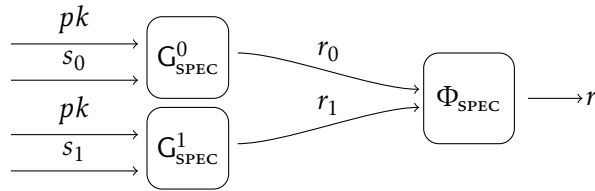


Figure 16: Immunizing backdoored PRG.

Security analysis. It is almost the same as the analysis of Theorem 3.1, if we can show (i) For each  $s_b$ , it is unpredictable to  $G_{\text{IMPL}}^{1-b}$ ; (ii)  $G_{\text{IMPL}}(s_0) \parallel G_{\text{IMPL}}(s_1)$  is unpredictable to the adversary who makes the implementations. This is easy to see since the watchdog can check the entropy contained in the output  $r_0, r_1$  respectively by feeding PRG with random seed and  $pk$ . We defer the proof to the full version.

**Simulating the semi-private model.** Semi-private model was introduced in [DGG<sup>+</sup>15], where a uniformly random string  $r$  is generated by a trusted party and  $r$  is made public only after the big brother provides the implementations. It was shown in [DGG<sup>+</sup>15] that in such a model, one can bypass the impossibility of public immunization of backdoored PRG output. We may expect more interesting results in this model due to the power of this extra trusted randomness.

Using our technique, we could ask the adversary to provide a stego-free specification of randomness generation, for example, we may ask the adversary to implement  $\text{RG}_{\text{IMPL}}^0, \text{RG}_{\text{IMPL}}^1, \Phi_{\text{IMPL}}$ , and the user can use them to generate bits that look uniform even to the adversary. In this way, we can simulate the semi-private model.