

# How to Subvert Backdoored Encryption: Security Against Adversaries that Decrypt All Ciphertexts

Thibaut Horel\*  
Harvard University

Sunoo Park†  
MIT

Silas Richelson  
UC Riverside

Vinod Vaikuntanathan‡  
MIT

## Abstract

In this work, we examine the feasibility of secure and undetectable point-to-point communication in a world where governments can read all the encrypted communications of their citizens. We consider a world where the only permitted method of communication is via a government-mandated encryption scheme, instantiated with government-mandated keys. Parties cannot simply encrypt ciphertexts of some other encryption scheme, because citizens caught trying to communicate outside the government’s knowledge (*e.g.*, by encrypting strings which do not appear to be natural language plaintexts) will be arrested. The one guarantee we suppose is that the government mandates an encryption scheme which *is* semantically secure against outsiders: a perhaps reasonable supposition when a government might consider it advantageous to secure its people’s communication against foreign entities. But then, what good is semantic security against an adversary that holds all the keys and has the power to decrypt?

We show that even in the pessimistic scenario described, citizens *can* communicate securely and undetectably. In our terminology, this translates to a positive statement: all semantically secure encryption schemes support *subliminal communication*. Informally, this means that there is a two-party protocol between Alice and Bob where the parties exchange ciphertexts of what appears to be a normal conversation even to someone who knows the secret keys and thus can read the corresponding plaintexts. And yet, at the end of the protocol, Alice will have transmitted her secret message to Bob. Our security definition requires that the adversary not be able to tell whether Alice and Bob are just having a normal conversation using the mandated encryption scheme, or they are using the mandated encryption scheme for subliminal communication.

Our topics may be thought to fall broadly within the realm of *steganography*: the science of hiding secret communication within innocent-looking messages, or *cover objects*. However, we deal with the non-standard setting of an adversarially chosen distribution of cover objects (*i.e.*, a stronger-than-usual adversary), and we take advantage of the fact that our cover objects are ciphertexts of a semantically secure encryption scheme to bypass impossibility results which we show for broader classes of steganographic schemes. We give several constructions of subliminal communication schemes under the assumption that key exchange protocols with pseudorandom messages exist (such as Diffie-Hellman, which in fact has truly random messages). Each construction leverages the assumed semantic security of the adversarially chosen encryption scheme, in order to achieve subliminal communication.

---

\*Supported, in part, by the National Science Foundation under grants CAREER IIS-1149662, and CNS-1237235, by the Office of Naval Research under grants YIP N00014-14-1-0485 and N00014-17-1-2131, and by a Google Research Award.

†Supported by the Center for Science of Information STC (CSoI), an NSF Science and Technology Center (grant agreement CCF-0939370), MACS project NSF grant CNS-1413920, and a Simons Investigator Award Agreement dated 2012-06-05.

‡Supported in part by NSF Grants CNS-1350619, CNS-1414119 and CNS-1718161, Alfred P. Sloan Research Fellowship, Microsoft Faculty Fellowship and a Steven and Renee Finn Career Development Chair from MIT.

# 1 Introduction

Suppose that we lived in a world where the government wished to read all the communications of its citizens, and thus decreed that citizens must not communicate in any way other than by using a specific, government-mandated encryption scheme with government-mandated keys. Even face-to-face communication is not allowed: in this Orwellian world, anyone who is caught speaking to another person will be arrested for treason. Similarly, anyone whose communications appear to be hiding information will be arrested: *e.g.*, if the plaintexts encrypted using the government-mandated scheme are themselves ciphertexts of a different encryption scheme. However, the one assumption that we entertain in this paper, is that the government-mandated encryption scheme is, in fact, semantically secure: this is a tenable supposition with respect to a government that considers secure encryption to be in its interest, in order to prevent foreign powers from spying on its citizens' communications.

A natural question then arises: is there any way that the citizens would be able to communicate in a fashion undetectable to the government, based only on the semantic security of the government-mandated encryption scheme, and *despite the fact that the government knows the keys and has the ability to decrypt all ciphertexts*?<sup>1</sup> What can semantic security possibly guarantee in a setting where the adversary has the private keys?

This question may appear to fall broadly within the realm of *steganography*: the science of hiding secret communications within other innocent-looking communications (called “cover objects”), in an undetectable way. Indeed, it can be shown that if two parties have a shared secret, then based on slight variants of existing techniques for *secret-key steganography*, they can conduct communications hidden from the government.<sup>2</sup>

However, the question of whether two parties who have never met before can conduct hidden communications is more interesting. This is related to the questions of *public-key steganography* and *steganographic key exchange* which were both first formalized by von Ahn and Hopper [vAH04]. Public-key steganography is inadequate in our setting since exchanging or publishing public keys is potentially conspicuous and thus is not an option in our setting. All prior constructions of steganographic key exchange require the initial sampling of a public random string that serves as a public parameter of the steganographic scheme. Intuitively, in these constructions, the public random string can be thought to serve the purpose of selecting a specific steganographic scheme from a family of schemes *after* the adversary has chosen a strategy. That is, the schemes crucially assume that the adversary (the dystopian government, in our story above) cannot choose its covert distribution as a function of the public parameter.

It is conservative and realistic to expect a malicious adversary to choose the covert distribution *after* the honest parties have decided on their communication protocol (including the public parameters). After all, malice never sleeps [Mic16]. Alas, we show that if the covert distribution is allowed to depend on the communication protocol, steganographic communication is impossible. In other words, for every purported steganographic communication protocol, there is a covert distribution (even one with high min-entropy) relative to which the communication protocol fails to embed subliminal messages. The relatively simple counterexample we construct is inspired by the impossibility of deterministic extraction.

**Semantic Security to the Rescue?** However, this impossibility result does not directly apply to our setting, as our covert distribution is restricted to be a sequence of ciphertexts (that may encrypt arbitrary messages). Moreover, the ciphertexts are semantically secure against entities that are not privy to the private keys. We define the notion of a *subliminal communication*

---

<sup>1</sup>We note that one could, alternatively, consider an adversary with decryption capabilities arising from possession of some sort of “backdoor.” For the purposes of this paper, we opted for the simpler and still sufficiently expressive model where the adversary’s decryption power comes from knowledge of all the decryption keys.

<sup>2</sup>We refer the reader to Section 1.3 for more details.

*scheme* (Definition 3.1) as a steganographic communication scheme where security holds relative to covertext distributions that are guaranteed to be ciphertexts of some semantically secure encryption scheme. Is there a way to use semantic security to enable subliminal communication?

Our first answer to this question is negative. In particular, consider the following natural construction: first, design an extractor function  $f$ ; then, to subliminally transmit a message bit  $b$ , sample encryptions  $c$  of a (even adversarially prescribed) plaintext  $m$  using independent randomness every time, until  $f(c) = b$ . There are two reasons this idea does not work. First, if the plaintext bit  $b$  is not random, the adversary can detect this fact by simply applying the extractor function  $f$  to the transmitted covertext. Second, the government can pick an adversarial (semantically secure) encryption scheme where the extractor function  $f$  is constant on all ciphertexts; this is again similar to the impossibility of deterministic extraction.

Nevertheless, we show how to circumvent these difficulties and use the semantic security of the underlying (adversarial) encryption scheme and construct a subliminal communication scheme.

**Theorem 1.1** (Informal version of Theorem 5.1). *Under the decisional Diffie-Hellman (DDH) assumption—or any other assumption that gives rise to a key exchange protocol with messages indistinguishable from random—there is a subliminal communication scheme which allows the transmission of  $O(\log \kappa)$  many bits per ciphertext after a setup phase of  $\tilde{O}(\log \kappa)$  ciphertexts ( $\kappa$  is the security parameter).*

We then show how to improve our first construction to reduce the length of the setup phase under additional assumptions.

## 1.1 Overview of Our Construction

The first idea in our construction is implicit in essentially all the works in steganography starting from [Sim83]: namely, to achieve subliminal communication of arbitrary messages, it is sufficient to be able to undetectably communicate *uniformly randomly distributed strings of one's choice*. In other words, Alice samples a string  $r$  which is randomly distributed, produces some ciphertext(s) to be sent to Bob, such that Bob is able to learn  $r$  from them, and yet a PPT eavesdropper Eve who sees the entire communication transcript cannot distinguish between the following two cases:

1. Alice is indeed sending (hereafter, “embedding”) random strings to Bob, or
2. Alice is producing ciphertexts using the unmodified government-mandated encryption algorithm, without embedding such random strings.

To be more precise, the indistinguishability requirement holds for any given (adversarially specified) distribution  $\mathcal{M}$  of message sequences that Alice may choose to encrypt using the government-mandated encryption scheme. Notice that this does not preclude that Eve may be able to learn  $r$  and indeed, our constructions do allow an eavesdropper to learn the embedded strings. Given the ability to undetectably communicate randomly distributed strings, Alice and Bob can then embed to each other the messages of a key-exchange protocol with randomly distributed messages (such as Diffie-Hellman) to establish a shared secret, and then embed to each other ciphertexts of a secret-key encryption scheme with pseudorandom ciphertexts, using the established secret as the key.

All known constructions of such *undetectable random string embedding* rely on the sampling of a public random seed after the adversarial strategy is fixed. In this paper, however, we are interested in bootstrapping hidden communications from the very ground up, and we are not willing to assume that the parties start from a state where such a seed is already present.

We observe that the ability to embed randomly distributed strings *of one’s choice* — rather than, *e.g.*, to apply a deterministic function to ciphertexts of the government-mandated encryption scheme, and thereby obtain randomly distributed strings which the creator of the ciphertexts did not choose — is crucial to the above-outlined scheme. The notion of undetectably embedding *exogenous* random strings — *i.e.*, strings that are randomly distributed outside of Alice’s control, but both Alice and Bob can read them — is seemingly much weaker, and certainly cannot be used to embed key exchange messages or secret-key ciphertexts. However, we observe that this weaker primitive turns out to be achievable, for our specific setting, without the troublesome starting assumption of a public random seed. We identify a method for embedding *exogenous* random strings into ciphertexts of an adversarially chosen encryption scheme (interestingly, our method does not generalize to embedding into arbitrary min-entropy distributions). We then exploit this method to allow the communicating parties to establish a random seed — from which point they can proceed to embed random strings *of their choice*, as described above.

In building this weaker primitive, in order to bypass our earlier-described impossibility result, we extract from two ciphertexts at a time, instead of one. We begin with the following simple idea: for each consecutive pair of ciphertexts  $c$  and  $c'$ , a single hidden (random) bit  $b$  is defined by  $b = f(c, c')$  where  $f$  is some two-source extractor. It is initially unclear why this should work because (1)  $c$  and  $c'$  are encryptions of messages  $m$  and  $m'$  which are potentially dependent, and two-source extractors are not guaranteed to work without independence; and (2) even if this difficulty could be overcome, ciphertexts of semantically secure encryption scheme can have min-entropy as small as  $\omega(\log \kappa)$  (where  $\kappa$  is the security parameter) and no two-source extractor known to this day can extract from such a small min-entropy.

We overcome difficulty (1) by relying on the semantic security of the ciphertexts of the adversarially chosen encryption scheme. Paradoxically, even though the adversary knows the decryption key, we exploit the fact that semantic security still holds against the *extractor*, which does not have the decryption key. The inputs in our case are ciphertexts which are not necessarily independent, but semantic security implies that they are computationally indistinguishable from being independent. Thus, the output of  $f(c, c')$  is pseudorandom. Indeed, when  $f$  outputs a single bit (as in our construction), the output is also statistically close to random. The crucial point here is that the semantic security of the encryption scheme is used not against the government, but rather against the extraction function  $f$ .

Our next observation, to address difficulty (2), is that the ciphertexts are not only computationally independent, but they are also computationally indistinguishable from i.i.d. In particular, each pair of ciphertexts is indistinguishable from a pair of encryptions of 0, by semantic security. Based on this observation, we can use a very simple “extractor”, namely, the greater-than function GT. In fact, GT is an extractor with two input sources, whose output bit has negligible bias when the sources have  $\omega(\log \kappa)$  min-entropy and are *independently and identically distributed* (this appears to be a folklore observation; see, *e.g.*, [BIW04]). Because of the last condition, GT is not a true two-source extractor according to standard definitions, but is still suitable for our setting.

By repeatedly extracting random bits from pairs of consecutive ciphertexts using GT, Alice and Bob can construct a shared random string  $s$ . Note that in this process, Alice and Bob generate ciphertexts using the unmodified government-mandated encryption scheme, so the indistinguishability requirement clearly holds. We stress again that  $s$  is also known to a passive eavesdropper of the communication. This part of our construction, up to the construction of the string  $s$ , is presented in details in Section 5.1. From there, constructing a subliminal communication scheme is not hard: Alice and Bob use  $s$  as the seed of a strong seeded extractor to subliminally communicate random strings *of their choice* as explained in Section 5.2. The complete description of our protocol is given in Section 5.3.

## 1.2 Improved Constructions for Specific Cases

While our first construction has the advantage of simplicity, the initial phase to agree on shared random string (using the GT function) transmits only one hidden bit per ciphertext of the government-mandated encryption scheme. A natural question is whether this rate of transmission can be improved. We show that if the government-mandated encryption scheme is *succinct* in the sense that the ciphertext expansion factor is at most 2, then it is possible to improve the rate of transmission in this phase to  $O(\log \kappa)$  hidden bits per ciphertext using an alternative construction based on the extractor from [DEOR04]. In other words, our first result showed that if the government-mandated encryption scheme is semantically secure, we can use it to communicate subliminally; the second result shows that if the government-mandated encryption scheme is efficient, that is even better for us, in the sense that it can be used for more efficient subliminal communication.

**Theorem 1.2** (Informal version of Theorem 6.1). *If there is a secure key exchange protocol whose message distribution is pseudorandom, then there is a subliminal communication scheme in which a shared seed is established in two exchanges of ciphertexts of a succinct encryption scheme.*

Theorem 1.1 exploited the specific nature of the cover object distribution in our setting (specifically, that a sequence of encryptions of arbitrary messages is indistinguishable from an i.i.d. sequence of encryptions of zero). Theorem 1.2 exploits an additional consequence of the semantic security of the government-mandated encryption scheme: if it is succinct, then ciphertexts are computationally indistinguishable from sources of high min-entropy (*i.e.*, they have large HILL-entropy).

It may be possible to use more advanced two-source extractors to work with a larger class of government-mandated encryption schemes (with larger expansion factors); however, the best known such extractors have an inverse polynomial error rate [CZ16] (whereas our construction’s extractor has negligible error). Consequently, designing a subliminal communication protocol using these extractors seems to require additional ideas, and we leave this as an open problem.

Finally, we show yet another approach in cases where the distribution of “innocent” messages to be encrypted under the government-mandated encryption scheme has a certain amount of conditional min-entropy. For such cases, we construct an alternative scheme that leverages the semantic security of the encryption scheme in a rather different way: namely, the key fact for this alternative construction is that (in the absence of a decryption key) a ciphertext appears independent of the message it encrypts. In this case, running a two-source extractor on the message and the ciphertext works. The resulting improvement in the efficiency of the scheme is comparable to that of Theorem 1.2.

**Theorem 1.3** (Informal version of Theorem 6.2). *If there is a secure key exchange protocol whose message distribution is pseudorandom, then there is a subliminal communication scheme:*

- *for any cover distribution consisting of ciphertexts of a semantically secure encryption scheme, if the innocent message distribution  $\mathcal{M}$  has conditional min-entropy rate  $1/2$ , or*
- *for any cover distribution consisting of ciphertexts of a semantically secure and succinct encryption scheme, if the innocent message distribution  $\mathcal{M}$  has conditional min-entropy  $\omega(\log \kappa)$ .*

*In both cases, the shared seed is established during the setup phase in only two exchanges of ciphertexts.*

We conclude this introductory section with some discussion of our results in a wider context.

**On Our Modeling Assumptions.** Our model considers a relatively powerful adversary that, for example, has the ability to choose the encryption scheme using which all parties must communicate, and to decrypt all such communications. We believe that this can be very realistic in certain scenarios, but it is also important to note the limitations that our model places on the adversary.

The most obvious limitation is that the encryption scheme chosen by the adversary must be semantically secure (against third parties that do not have the ability to decrypt). Another assumption is that citizens are able to run algorithms of their choice on their own computers without, for instance, having every computational step monitored by the government. Moreover, citizens may use encryption randomness of their choice when producing ciphertexts of the government-mandated encryption scheme: in fact, this is a key fact that our construction exploits. Interestingly, secrecy of the encryption randomness from the adversary is irrelevant: after all, the adversary can always choose an encryption scheme where the encryption randomness is recoverable given the decryption key. Despite this, the ability of the encryptor to choose the randomness to input to the encryption algorithm can be exploited—as by our construction—to allow for subliminal communication.

**The Meaning of Semantic Security when the Adversary Can Decrypt.** In an alternate light, our work may be viewed as asking the question: *what guarantee, if any, does semantic security provide against adversary in possession of the decryption key?* Our results find, perhaps surprisingly, that some meaningful guarantee is still provided by semantic security even against an adversary is able to decrypt: more specifically, that *any* communication channel allowing transmission of ciphertexts can be leveraged to allow for undetectable communications between two parties that have never met. From this perspective, our work may be viewed as the latest in a scattered series of recent works that consider what guarantees can be provided by cryptographic primitives that are somehow “compromised”—examples of recent works in this general flavor are cited in Section 1.3 below.

**Concrete Security Parameters.** From a more practical perspective, it may be relevant to consider that the government in our hypothetical Orwellian scenario would be incentivized to opt for an encryption scheme with the least possible security level so as to ensure security against foreign powers. In cases where the government considers itself to have more computational power than foreign adversaries (perhaps by a constant factor), this could create an interesting situation where the security parameter with which the government-mandated scheme must be instantiated is *below* what is necessary to ensure security against the government’s own computational power.

Such a situation could be risky for citizens’ hidden communications: intuitively, our constructions guarantee indistinguishability *against the citizens’ own government* between an “innocent” encrypted conversation and one which is carrying hidden subliminal messages. However, the distinguishing advantage in this indistinguishability game *depends on the security parameter* of the government-mandated encryption scheme. Thus, it could be that the two distributions are far enough apart for the citizens’ own government to distinguish (though not for foreign governments to distinguish). We observe that citizens cognizant of this situation can further reduce the distinguishing advantage beyond that provided by our basic construction, using the standard technique of amplifying the proximity of a distribution (which is far from random) to uniformly random, by taking the XOR of several samples from the far-from-random distribution.

Having outlined this potential concern and solution, in the rest of the paper we will disregard these issues in the interest of clarity of exposition, and present a purely asymptotic analysis.

**Open Problems.** Our work suggests a number of open problems. A natural one is the extent to which the modeling assumptions that this work makes — such as the ability of honest encryptors to use true randomness for encryption — can be relaxed or removed, while preserving

the ability to communicate subliminally. For example, one could imagine yet another alternate universe, in which the hypothetical Orwellian government not only mandates that citizens use the prescribed encryption scheme, but also that their encryption randomness must be derived from a specific government-mandated pseudorandom generator.

The other open problems raised by our work are of a more technical nature and better understood in the context of the specific details of our constructions; for this reason we defer their discussion to Section 7.

### 1.3 Other Related Work

The scientific study of steganography was initiated by Simmons more than thirty years ago [Sim83], and is the earliest mention of the term “subliminal channel” referring to the conveyance of information in a cryptosystem’s output in a way that is different from the intended output,<sup>3</sup> of which we are aware. Subsequent works such as [Cac98, Mit99, ZFK<sup>+</sup>98] initially explored information-theoretic treatments of steganography, and then Hopper, Langford, and von Ahn [HLv02] gave the first complexity-theoretic (secret-key) treatment almost two decades later. Public-key variants of steganographic notions—namely, public-key steganography and steganographic key exchange—were first defined by [vAH04]. There is very little subsequent literature on public-key steganographic primitives; one notable example is by Backes and Cachin [BC05], which considers public-key steganography against active attacks (their attack model, which is stronger than that of [vAH04], was also considered in [HLv02] but had never been applied to the public-key setting).

The alternative perspective of our work as addressing the question of whether any sort of secret communication can be achieved via transmission of ciphertexts of an adversarially designed cryptosystem alone fits into a scattered series of recent works that consider what guarantees can or cannot be provided by compromised cryptographic primitives. For example, Goldreich [Gol11], and later, Cohen and Klein [CK16], consider what unpredictability guarantee is achieved by the classic GGM construction [GGM86] when the traditionally secret seed is known; Austrin et al. [ACM<sup>+</sup>14] study whether certain cryptographic primitives can be secure even in the presence of an adversary that has limited ability to tamper with honest parties’ randomness; Dodis et al. [DGG<sup>+</sup>15] consider what cryptographic primitives can be built based on backdoored pseudorandom generators; and Bellare, Jaeger, and Kane [BJK15] present attacks that work against any symmetric-key encryption scheme, that completely compromise security by undetectably corrupting the algorithms of the encryption scheme (such attacks might, for example, be feasible if an adversary could generate a bad version of a widely used cryptographic library and install it on his target’s computer).

The last work mentioned above, [BJK15], is actually part of the broader field of kleptography, originally introduced by Young and Yung [YY97, YY96b, YY96a], which is also relevant context for the present work. Broadly speaking, a *kleptographic attack* “uses cryptography against cryptography” [YY97] — *i.e.*, changes the behavior of a cryptographic system in a fashion undetectable to an honest user with black-box access to the cryptosystem, such that the use of the modified system leaks some secret information (*e.g.*, plaintexts or key material) to the attacker who performed the modification. An example of such an attack might be to modify the key generation algorithm of an encryption scheme such that an adversary in possession of a “back door” can derive the private key from the public key, yet an honest user finds the generated key pairs to be indistinguishable from correctly produced ones. Kleptography has enjoyed renewed research activity since [BPR14] introduced a formal model of a specific type of kleptographic attack called *algorithm substitution attacks* (ASAs), motivated by recent revelations suggesting that intelligence agencies have successfully implemented attacks of this

<sup>3</sup>This phrasing is loosely borrowed from [YY97].

nature at scale. Recently, [BL17] formalized an equivalence between certain variants of ASA and steganography.

Our setting differs significantly from kleptography in that the encryption algorithms are public and not tampered with (*i.e.*, adhere to a purported specification), and in fact may be *known* to be designed by an adversarial party.

## 2 Preliminaries

**Notation.**  $\kappa$  is the security parameter throughout. PPT means “probabilistic polynomial time.”  $[n]$  denotes the set  $\{1, \dots, n\}$ .  $U_n$  is a uniform variable over  $\{0, 1\}^n$ , independent of every other variable in this paper. We write  $X \sim Y$  to express that  $X$  and  $Y$  are identically distributed. Given two variables  $X$  and  $Y$  over  $\{0, 1\}^k$ , we denote by  $\|X - Y\|_s$  the statistical distance defined by:

$$\|X - Y\|_s = \frac{1}{2} \sum_{x \in \{0, 1\}^k} |\Pr[X = x] - \Pr[Y = x]| = \max_{S \subseteq \{0, 1\}^k} |\Pr[X \in S] - \Pr[Y \in S]|.$$

For a random variable  $X$ , we define the min-entropy of  $X$  by  $H_\infty(X) = -\log \max_x \Pr[X = x]$ . The collision probability is  $\text{CP}(X) = \sum_x \Pr[X = x]^2$ .

### 2.1 Encryption and Key Exchange

We assume familiarity with the standard notions of semantically secure public-key and private-key encryption, and key exchange. This subsection defines notation and additional terminology.

**Public-Key Encryption.** We use the notation  $\mathsf{E} = (\mathsf{E.Gen}, \mathsf{E.Enc}, \mathsf{E.Dec})$  for the public-key encryption scheme mandated by the adversary.

**Secret-key Encryption.** We write  $\mathsf{SKE} = (\mathsf{SKE.Gen}, \mathsf{SKE.Enc}, \mathsf{SKE.Dec})$  to denote a secret-key encryption scheme. We define a *pseudorandom secret-key encryption scheme* to be a secret-key encryption scheme whose ciphertexts are indistinguishable from random. It is a standard result that pseudorandom secret-key encryption schemes can be built from one-way functions.

**Key Exchange.** We define a *pseudorandom key-exchange protocol* to be a key-exchange protocol whose transcripts are distributed indistinguishably from random messages. Recall that the standard security guarantee for key-exchange protocols requires that  $(T, K) \stackrel{c}{\approx} (T, K_\S)$ , where  $T$  is a key-exchange protocol transcript,  $K$  is the shared key established in  $T$ , and  $K_\S$  is a random unrelated key. A pseudorandom key-exchange protocol instead requires that  $(T, K) \stackrel{c}{\approx} (U, K_\S)$  where  $U$  is the uniform distribution over strings of the appropriate length.

Most known key agreement protocols are pseudorandom; in fact, most have truly random messages. This is the case, for example, for the classical protocol of Diffie and Hellman [DH76].

### 2.2 Extractors

We will need the following definitions of two-source and seeded extractors.

**Definition 2.1.** The family  $2\text{Ext} : \{0, 1\}^n \times \{0, 1\}^m \rightarrow \{0, 1\}^\ell$  is a  $(k_1, k_2, \varepsilon)$  *two-source extractor* if for all  $\kappa \in \mathbb{N}$  and for all pairs  $(X, Y)$  of independent random variables over  $\{0, 1\}^{n(\kappa)} \times \{0, 1\}^{m(\kappa)}$  such that  $H_\infty(X) \geq k_1(\kappa)$  and  $H_\infty(Y) \geq k_2(\kappa)$ , it holds that:

$$\|2\text{Ext}_\kappa(X, Y) - U_{\ell(\kappa)}\|_s \leq \varepsilon(\kappa). \quad (1)$$

We say that  $2\text{Ext}$  is *strong* w.r.t. the first input if it satisfies the following stronger property:

$$\|(X, 2\text{Ext}_\kappa(X, Y)) - (X, U_{\ell(\kappa)})\|_s \leq \varepsilon(\kappa).$$

A strong two-source extractor w.r.t. the second input is defined analogously. Finally, we say that  $2\text{Ext}$  is a  $(k, \varepsilon)$  *same-source* extractor if  $n = m$  and (1) is only required to hold when  $(X, Y)$  is a pair of i.i.d. random variables with  $H_\infty(X) = H_\infty(Y) \geq k(\kappa)$ .

**Definition 2.2.** The family  $\text{Ext} : \{0, 1\}^n \times \{0, 1\}^m \rightarrow \{0, 1\}^\ell$  is a  $(k, \varepsilon)$  *seeded extractor* if for all  $\kappa \in \mathbb{N}$  and any random variable  $X$  over  $\{0, 1\}^{m(\kappa)}$  such that  $H_\infty(X) \geq k(\kappa)$ , it holds that:

$$\|\text{Ext}_\kappa(U_{n(\kappa)}, X) - U_{\ell(\kappa)}\|_{\mathfrak{s}} \leq \varepsilon(\kappa).$$

We say moreover that  $\text{Ext}$  is *strong* if it satisfies the following stronger property:

$$\|(U_{n(\kappa)}, \text{Ext}_\kappa(U_{n(\kappa)}, X)) - (U_{n(\kappa)}, U_{\ell(\kappa)})\|_{\mathfrak{s}} \leq \varepsilon(\kappa).$$

### 3 Subliminal Communication

**Conversation Model.** The protocols we will construct take place over a communication between two parties  $P_0$  and  $P_1$  alternately sending each other ciphertexts of a public-key encryption scheme. *W.l.o.g.*, we assume that  $P_0$  initiates the communication, and that communication occurs over a sequence of *exchange-rounds* each of which comprises two sequential messages: in each exchange-round, one party  $P_b$  sends a message to  $P_{1-b}$  and then  $P_{1-b}$  sends a message to  $P_b$ . Let  $m_{b,i}$  denote the plaintext message sent by  $P_b$  to  $P_{1-b}$  in exchange-round  $i$ , and let  $\mathbf{m}_i = (m_{0,i}, m_{1,i})$  denote the pair of messages exchanged. For  $i \geq 1$ , let us denote by

$$\boldsymbol{\tau}_{0,i} = (\mathbf{m}_1, \dots, \mathbf{m}_{i-1}) \text{ and } \boldsymbol{\tau}_{1,i} = (\mathbf{m}_1, \dots, \mathbf{m}_{i-1}, m_{0,i})$$

the plaintext transcripts available to  $P_0$  and  $P_1$  respectively during exchange-round  $i$ , in the case when  $P_0$  sends the first message in exchange-round  $i$ .<sup>4</sup> We define  $\boldsymbol{\tau}_{0,0}$  and  $\boldsymbol{\tau}_{1,0}$  to be empty lists (*i.e.*, empty starting transcripts). (Note that when a notation contains both types of subscripts, we write the subscripts denoting the party and round in **blue** and **red** respectively, to improve readability.)

Recall that our adversary has the power to decrypt all ciphertexts under its chosen public-key encryption scheme  $\mathsf{E}$ . Intuitively, it is therefore important that the plaintext conversation between  $P_0$  and  $P_1$  appears innocuous (and does not, for example, consist of ciphertexts of another encryption scheme). To model this, we assume the existence of a next-message distribution  $\mathcal{M}$ , which outputs a next innocuous message given the transcript of the plaintext conversation so far. This is denoted by  $m_{b,i} \leftarrow \mathcal{M}(\boldsymbol{\tau}_{b,i})$ .

**Remark 1.** We emphasize that our main results make no assumptions at all on the distribution  $\mathcal{M}$ , and require only that the parties have oracle access to their own next-message distributions. Our main results hold in the presence of *arbitrary* message distributions: for example, they hold even in the seemingly inauspicious case when  $\mathcal{M}$  is constant, meaning the parties are restricted to repeatedly exchanging a fixed message.

In Section 6, we discuss other more efficient constructions that can be used in settings where a stronger assumption — namely, that  $\mathcal{M}$  has a certain amount of min-entropy — is acceptable. This stronger assumption, while not without loss of generality, might be rather benign in certain contexts (for example, if the messages exchanged are images).

In all the protocols we consider, the symbol  $\mathfrak{s}$  is used to denote internal state kept locally by  $P_0$  and  $P_1$ . It is implicitly assumed that each party's state contains an up-to-date transcript of all messages received during the protocol. Parties may additionally keep other information

<sup>4</sup>If instead  $P_1$  spoke first in round  $i$ , then  $\boldsymbol{\tau}_{0,i}$  would contain  $m_{1,i}$ , and  $\boldsymbol{\tau}_{1,i}$  would not contain  $m_{0,i}$ .

in their internal state, as a function of the local computations they perform. For  $i \geq 1$ ,  $\mathfrak{s}_{b,i}$  denotes the state of  $P_b$  at the conclusion of exchange-round  $i$ . Initial states  $\mathfrak{s}_{b,0} = \emptyset$  are empty.

We begin with a simpler definition that only syntactically allows for the transmission of a single message (Definition 3.1). This both serves as a warm-up to the multi-message definition presented next (Definition 3.3), and will be used in its own right to prove impossibility results. See Remark 2 for further discussion of the relationship between these two definitions.

**Definition 3.1.** A *subliminal communication scheme* is a two-party protocol:

$$\Pi^{\mathbb{E}} = (\Pi_{0,1}^{\mathbb{E}}, \Pi_{1,1}^{\mathbb{E}}, \Pi_{0,2}^{\mathbb{E}}, \Pi_{1,2}^{\mathbb{E}}, \dots, \Pi_{0,r}^{\mathbb{E}}, \Pi_{1,r}^{\mathbb{E}}; \Pi_{1,\text{out}}^{\mathbb{E}})$$

where  $r \in \text{poly}$  is the number of exchange-rounds and each  $\Pi_{b,i}^{\mathbb{E}}$  is a PPT algorithm with oracle access to the algorithms of a public-key encryption scheme  $\mathbb{E}$ . Party  $P_0$  is assumed to receive as input a message  $\text{msg}$  (of at least one bit) that is to be conveyed to  $P_1$  in an undetectable fashion. The algorithms  $\Pi_{b,i}^{\mathbb{E}}$  are used by  $P_b$  in round  $i$ , respectively, and  $\Pi_{1,\text{out}}^{\mathbb{E}}$  denotes the algorithm run by  $P_1$  to produce an output  $\text{msg}'$  at the end of the protocol.

A subliminal communication scheme must satisfy the following syntax, correctness and security guarantees.

- **Syntax.** In each exchange-round  $i = 1, \dots, r$ :

$P_0$  performs the following steps:

1. Sample “innocuous message”  $m_{0,i} \leftarrow \mathcal{M}(\tau_{0,i-1})$ .
2. Generate ciphertext and state  $(c_{0,i}, \mathfrak{s}_{0,i}) \leftarrow \Pi_{0,i}^{\mathbb{E}}(\text{msg}, m_{0,i}, \text{pk}_1, \mathfrak{s}_{0,i-1})$ .
3. Locally store  $\mathfrak{s}_{0,i}$  and send  $c_{0,i}$  to  $P_1$ .

Then,  $P_1$  performs the following steps:<sup>5</sup>

1. Sample “innocuous message”  $m_{1,i} \leftarrow \mathcal{M}(\tau_{1,i-1})$ .
2. Generate ciphertext and state  $(c_{1,i}, \mathfrak{s}_{1,i}) \leftarrow \Pi_{1,i}^{\mathbb{E}}(m_{1,i}, \text{pk}_0, \mathfrak{s}_{1,i-1})$ .
3. Locally store  $\mathfrak{s}_{1,i}$  and send  $c_{1,i}$  to  $P_0$ .

After  $r$  rounds,  $P_1$  computes  $\text{msg}' = \Pi_{1,\text{out}}^{\mathbb{E}}(\text{sk}_1, \mathfrak{s}_{1,r})$  and halts.

- **Correctness.** For any  $\text{msg} \in \{0,1\}^\kappa$ , if  $P_0$  and  $P_1$  play  $\Pi^{\mathbb{E}}$  honestly, then  $\text{msg}' = \text{msg}$  with probability  $1 - \text{negl}(\kappa)$ . The probability is taken over the key generation  $(\text{pk}_1, \text{sk}_1), (\text{pk}_2, \text{sk}_2) \leftarrow \text{E.Gen}$  and the randomness of the protocol algorithms, as well as the message distribution  $\mathcal{M}$ .
- **Subliminal Indistinguishability.** For any semantically secure public-key encryption scheme  $\mathbb{E}$ , any  $\text{msg} \in \{0,1\}^\kappa$  and any next-message distribution  $\mathcal{M}$ , for  $(\text{pk}_1, \text{sk}_1), (\text{pk}_2, \text{sk}_2) \leftarrow \text{E.Gen}$ , the following distributions are computationally indistinguishable:

Ideal( $\text{pk}_1, \text{sk}_1, \text{pk}_2, \text{sk}_2, \mathcal{M}$ ):	Subliminal $_{\Pi}$ ( $\text{msg}, \text{pk}_1, \text{sk}_1, \text{pk}_2, \text{sk}_2, \mathcal{M}$ ):
for $i = 1, \dots, r$ :	for $i = 1, \dots, r$ :
$m_{0,i} \leftarrow \mathcal{M}(\tau_{0,i})$	$m_{0,i} \leftarrow \mathcal{M}(\tau_{0,i})$
$m_{1,i} \leftarrow \mathcal{M}(\tau_{1,i})$	$m_{1,i} \leftarrow \mathcal{M}(\tau_{1,i})$
$c_{0,i} \leftarrow \text{E.Enc}(\text{pk}_1, m_{0,i})$	$(c_{0,i}, \mathfrak{s}_{0,i}) \leftarrow \Pi_{0,i}^{\mathbb{E}}(\text{msg}, m_{0,i}, \text{pk}_1, \mathfrak{s}_{0,i-1})$
$c_{1,i} \leftarrow \text{E.Enc}(\text{pk}_0, m_{1,i})$	$(c_{1,i}, \mathfrak{s}_{1,i}) \leftarrow \Pi_{1,i}^{\mathbb{E}}(m_{1,i}, \text{pk}_0, \mathfrak{s}_{1,i-1})$
output $(\text{pk}_1, \text{sk}_1, \text{pk}_2, \text{sk}_2; (c_{b,i})_{b \in \{0,1\}, i \in [r]})$	output $(\text{pk}_1, \text{sk}_1, \text{pk}_2, \text{sk}_2; (c_{b,i})_{b \in \{0,1\}, i \in [r]})$

<sup>5</sup>Note that the steps executed by  $P_0$  and  $P_1$  are entirely symmetric except in the following two aspects: first,  $P_0$ 's input  $\text{msg}$  is present in step 2 but not in step 2; and secondly, the state  $\mathfrak{s}_{1,i-1}$  used in step 2 contains the round- $i$  message  $c_{0,i}$ , whereas the state  $\mathfrak{s}_{0,i-1}$  used in step 2 depends only on the transcript until round  $i - 1$ .

If the *subliminal indistinguishability* requirement is satisfied only for next-message distributions  $\mathcal{M}$  in a restricted set  $\mathbb{M}$ , rather than for any  $\mathcal{M}$ , then  $\Pi$  is said to be a *subliminal communication scheme* for  $\mathbb{M}$ .

**Definition 3.2.** The *rate* of a subliminal communication protocol  $\Pi$  is defined as  $\frac{2r}{\kappa}$ , where  $r$  is defined as in Definition 3.1.<sup>6</sup> This is the average number of bits which are subliminally communicated per ciphertext of  $\mathbb{E}$ .

For simplicity, Definition 3.1 presents a communication scheme in which only a single hidden message  $\text{msg}$  is transmitted. More generally, it is desirable to transmit multiple messages, and bidirectionally, and perhaps in an adaptive manner.<sup>7</sup> In multi-message schemes, it may be beneficial for efficiency that the protocol have a two-phase structure where some initial preprocessing is done in the first phase, and then the second phase can thereafter be invoked many times to transmit different hidden messages.<sup>8</sup> This will be a useful notion later in the paper, for our constructions, so we give the definition of a multi-message scheme here.

**Definition 3.3.** A *multi-message subliminal communication scheme* is a two-party protocol defined by a pair  $(\Phi, \Xi)$  where  $\Phi$  (“Setup Phase”) and  $\Xi$  (“Communication Phase”) each define a two-party protocol. Each party outputs a state at the end of  $\Phi$ , which it uses as an input in each subsequent invocation of  $\Xi$ . An execution of a multi-message subliminal communication scheme consists of an execution of  $\Phi$  followed by one or more executions of  $\Xi$ . More formally:

$$\begin{aligned}\Phi^{\mathbb{E}} &= (\Phi_{0,1}^{\mathbb{E}}, \Phi_{1,1}^{\mathbb{E}}, \Phi_{0,2}^{\mathbb{E}}, \Phi_{1,2}^{\mathbb{E}}, \dots, \Phi_{0,r}^{\mathbb{E}}, \Phi_{1,r}^{\mathbb{E}}) \\ \Xi^{\mathbb{E}} &= (\Xi_{0,1}^{\mathbb{E}}, \Xi_{1,1}^{\mathbb{E}}, \Xi_{0,2}^{\mathbb{E}}, \Xi_{1,2}^{\mathbb{E}}, \dots, \Xi_{0,r'}^{\mathbb{E}}, \Xi_{1,r'}^{\mathbb{E}}; \Xi_{1,\text{out}}^{\mathbb{E}})\end{aligned}$$

where  $r, r' \in \text{poly}$  are the number of exchange-rounds in  $\Phi$  and  $\Xi$  respectively. and where each  $\Phi_{b,i}^{\mathbb{E}}, \Xi_{b,i}^{\mathbb{E}}$  is a PPT algorithm with oracle access to the algorithms of a public-key encryption scheme  $\mathbb{E}$ . The protocol must satisfy the following syntax, correctness and security guarantees.

- **Syntax.** In each exchange-round  $i = 1, \dots, r$  of  $\Phi$ :  $P_0$  executes the following steps for  $b = 0$ , and then  $P_1$  executes the same steps for  $b = 1$ .
  1. Sample “innocuous message”  $m_{b,i} \leftarrow \mathcal{M}(\tau_{b,i-1})$ .
  2. Generate ciphertext and state  $(c_{b,i}, \mathfrak{s}_{b,i}) \leftarrow \Phi_{b,i}^{\mathbb{E}}(m_{b,i}, \text{pk}_{1-b}, \mathfrak{s}_{b,i-1})$ .
  3. Locally store  $\mathfrak{s}_{b,i}$  and send  $c_{b,i}$  to  $P_{1-b}$ .

After the completion of  $\Phi$ , either party may initiate  $\Xi$  by sending a first message of the  $\Xi$  protocol (with respect to a message  $\text{msg}$  to be steganographically hidden, known to the initiating party). Let  $P_S$  and  $P_R$  denote the initiating and non-initiating parties in an execution of  $\Xi$ , respectively.<sup>9</sup> Let  $\text{msg} \in \{0, 1\}^\kappa$  be the hidden message that  $P_S$  is to transmit to  $P_R$  in an undetectable fashion during an execution of  $\Xi$ .

The execution of  $\Xi$  proceeds as follows over exchange-rounds  $i' = 1, \dots, r'$ :

<sup>6</sup>The factor of two comes from the fact that each exchange-round contains two messages.

<sup>7</sup>That is, the messages to be transmitted may become known as the protocol progresses, rather than all being known at the outset. This is the case, for example, if future messages depend on responses to previous ones.

<sup>8</sup>As a concrete example: consider a simple protocol for transmitting a single encrypted message, consisting of key exchange followed by the transmission of message encrypted under the established key. When adapting this protocol to support multiple messages, it is beneficial to split the protocol into a one-time “phase 1” consisting of key exchange, and a “phase 2” encompassing the ciphertext transmission which can be invoked many times on different messages using the same phase-1 key. Such a protocol has much better amortized efficiency than simply repeating the single-message protocol many times, *i.e.*, establishing a new key for each ciphertext.

<sup>9</sup>Subscripts  $S, R \in \{0, 1\}$  stand for “sender” and “receiver,” respectively.

- $P_S$  acts as follows:
  1. Sample  $m_{S,r+i'} \leftarrow \mathcal{M}(\tau_{S,r+i'-1})$ .
  2. Generate  $(c_{S,r+i'}, \mathfrak{s}_{S,r+i'}) \leftarrow \Xi_{0,i'}^E(\text{msg}, m_{S,r+i'}, \text{pk}_R, \mathfrak{s}_{S,r+i'-1})$ .
  3. Locally store  $\mathfrak{s}_{S,r+i'}$  and send  $c_{S,r+i'}$  to  $P_R$ .
- $P_R$  acts as follows:
  1. Sample  $m_{R,r+i'} \leftarrow \mathcal{M}(\tau'_{R,r+i'-1})$ .
  2. Generate  $(c_{R,r+i'}, \mathfrak{s}_{R,r+i'}) \leftarrow \Xi_{1,i'}^E(m_{R,r+i'}, \text{pk}_S, \mathfrak{s}_{R,r+i'-1})$ .
  3. Locally store  $\mathfrak{s}_{R,r+i'}$  and send  $c_{R,r+i'}$  to  $P_S$ .

At the end of an execution of  $\Xi$ ,  $P_R$  computes  $\text{msg}' = \Xi_{1,\text{out}}^E(\text{sk}_1, \mathfrak{s}_{1,r+r'})$ .

- **Correctness.** For any  $\text{msg} \in \{0, 1\}^\kappa$ , if  $P_0$  and  $P_1$  execute  $(\Phi, \Xi)$  honestly, then for every execution of  $\Xi$ , the transmitted and received messages  $\text{msg}$  and  $\text{msg}'$  are equal with overwhelming probability. The probability is taken over the key generation  $(\text{pk}_1, \text{sk}_1), (\text{pk}_2, \text{sk}_2) \leftarrow \text{E.Gen}$  and the randomness of the protocol algorithms, as well as the message distribution  $\mathcal{M}$ .
- **Subliminal Indistinguishability.** For any semantically secure public-key encryption scheme  $\text{E}$ , any polynomial  $p = p(\kappa)$ , any sequence of hidden messages  $\vec{\text{msg}} = (\text{msg}_i)_{i \in [p]} \in (\{0, 1\}^\kappa)^p$ , any sequence of bits  $\vec{b} = (b_1, \dots, b_p) \in \{0, 1\}^p$  and any next-message distribution  $\mathcal{M}$ , for  $(\text{pk}_b, \text{sk}_b) \leftarrow \text{E.Gen}$ ,  $b \in \{0, 1\}$  the following distributions are computationally indistinguishable:

Ideal( $\text{pk}_1, \text{sk}_1, \text{pk}_2, \text{sk}_2, \mathcal{M}$ ):	Subliminal $_{\Phi, \Xi}(\vec{\text{msg}}, \vec{b}, \text{pk}_1, \text{sk}_1, \text{pk}_2, \text{sk}_2, \mathcal{M})$ :
for $i = 1, \dots, r + pr'$ : $m_{0,i} \leftarrow \mathcal{M}(\tau_{0,i})$ $m_{1,i} \leftarrow \mathcal{M}(\tau_{1,i})$ $c_{0,i} \leftarrow \text{E.Enc}(\text{pk}_1, m_{0,i})$ $c_{1,i} \leftarrow \text{E.Enc}(\text{pk}_0, m_{1,i})$ output: $(\text{pk}_1, \text{sk}_1, \text{pk}_2, \text{sk}_2; (c_{b,i})_{b \in \{0,1\}, i \in [r+pr']})$	for $i = 1, \dots, r$ : $m_{0,i} \leftarrow \mathcal{M}(\tau_{0,i})$ $m_{1,i} \leftarrow \mathcal{M}(\tau_{1,i})$ $(c_{0,i}, \mathfrak{s}_{0,i}) \leftarrow \Phi_{0,i}^E(\text{msg}, m_{0,i}, \text{pk}_1, \mathfrak{s}_{0,i-1})$ $(c_{1,i}, \mathfrak{s}_{1,i}) \leftarrow \Phi_{1,i}^E(m_{1,i}, \text{pk}_0, \mathfrak{s}_{1,i-1})$ for $j = 1, \dots, p$ : let $\beta = b_j$ and $\bar{\beta} = 1 - b_j$ for $i' = 1, \dots, r'$ : let $\iota = r + (j - 1)r' + i'$ $m_{\beta,\iota} \leftarrow \mathcal{M}(\tau_{\beta,\iota})$ $m_{\bar{\beta},\iota} \leftarrow \mathcal{M}(\tau_{\bar{\beta},\iota})$ $(c_{\beta,\iota}, \mathfrak{s}_{\beta,\iota}) \leftarrow \Xi_{\beta,i'}^E(\text{msg}, m_{\beta,\iota}, \text{pk}_{\bar{\beta}}, \mathfrak{s}_{\beta,\iota-1})$ $(c_{\bar{\beta},\iota}, \mathfrak{s}_{\bar{\beta},\iota}) \leftarrow \Xi_{\bar{\beta},i'}^E(m_{\bar{\beta},\iota}, \text{pk}_{\beta}, \mathfrak{s}_{\bar{\beta},\iota-1})$ output: $(\text{pk}_1, \text{sk}_1, \text{pk}_2, \text{sk}_2; (c_{b,i})_{b \in \{0,1\}, i \in [r+pr']})$

If the *subliminal indistinguishability* requirement is satisfied only for  $M$  in a restricted set  $\mathbb{M}$ , rather than for any  $\mathcal{M}$ , then  $(\Phi, \Xi)$  is said to be a *multi-message subliminal communication scheme* for  $\mathbb{M}$ .

**Definition 3.4.** The *asymptotic rate* of a multi-message subliminal communication protocol  $(\Phi, \Xi)$  is defined as  $\frac{\kappa}{2r'}$ , where  $r'$  is defined as in Definition 3.3. The asymptotic rate is the average number of bits which are subliminally communicated per ciphertext exchanged between  $P_0$  and  $P_1$  after the one-time setup phase is completed.

**Definition 3.5.** The *setup cost* of a multi-message subliminal communication protocol  $(\Phi, \Xi)$  is defined as  $r$ , *i.e.*, the number of rounds in  $\Phi$ . The setup cost is the number of ciphertexts which must be sent back and forth between  $P_0$  and  $P_1$  in order to complete the setup phase.

**Remark 2.** Definition 3.3 is *equivalent* to Definition 3.1 in the sense that the existence of any single-message scheme trivially implies a multi-message scheme and vice versa. We present Definition 3.3 as it will be useful for presenting and analyzing asymptotic efficiency of our constructions, but note that this equivalence means that the simpler Definition 3.1 suffices in the context of impossibility (or possibility) results, such as that given in Section 4.

## 4 Impossibility Results

### 4.1 Locally Decodable Subliminal Communication Schemes

A first attempt at achieving subliminal communication might consider schemes with the following natural property: the receiving party  $P_1$  extracts hidden bits *one ciphertext at a time*, by the application of a single (possibly randomized) decoding function. We refer to such schemes as *locally decodable* and our next impossibility theorem shows that non-trivial locally decodable schemes do not exist if the encryption scheme  $E$  is chosen adversarially.

**Theorem 4.1.** *For any locally decodable protocol  $\Pi$  satisfying the syntax of a single-message<sup>10</sup> subliminal communication scheme, there exists a semantically secure public-key encryption scheme  $E$  dependent on the public randomness of  $\Pi$ , such that  $E$  violates the correctness condition of Definition 3.1. Therefore, no locally decodable protocol  $\Pi$  is a subliminal communication scheme.*

*Proof.* Let us consider a locally decodable scheme such as in the statement of the theorem, and let us denote by  $\Pi_{2,\text{out}} : \{0,1\}^n \times \{0,1\}^m \rightarrow \{0,1\}$  the decoding function of the scheme where the second input consists of random bits (the public randomness) and the first input is a ciphertext  $c$ . Since we allow the encryption scheme to depend on the public randomness of the subliminal scheme, define the partial function  $f_r(c) = \Pi_{1,\text{out}}(c, r)$ .  $f_r$  is now a deterministic function of the ciphertext and we conclude the proof by constructing an encryption scheme which biases the output of  $f_r$  arbitrarily close to a constant bit. This is a contradiction, since by correctness and subliminal indistinguishability,  $f_r(c)$  should have negligible bias when subliminally communicating a uniformly random message  $\text{msg} \leftarrow \{m_1, m_2\}$ .

Let  $E$ , be a semantically secure encryption scheme with ciphertext space  $\mathcal{C} = \{0,1\}^n$  and message space  $\mathcal{M}$ . Without loss of generality we assume that for at least half the messages  $m \in \mathcal{M}$ , we have  $\Pr[f_r(E.\text{Enc}(pk, m)) = 1] \geq \frac{1}{2}$  (otherwise we can just replace 1 by 0 in the construction below). We now define the encryption scheme  $E'$  which is identical to  $E$  except for  $E'.\text{Enc}$  which on input  $(pk, m)$  runs as follows for some constant  $t$ .

1. Repeat at most  $t$  times:
  - (a) Sample encryption  $c \leftarrow E.\text{Enc}(pk, m)$ .
  - (b) If  $f_r(c) = 1$ , exit the loop; otherwise, continue.
2. Output  $c$ .

It is clear that  $E'$  is also semantically secure: oracle access to  $E'.\text{Enc}$  can be simulated with oracle access to  $E.\text{Enc}$ , so a distinguisher which breaks the semantic security of  $E'$  can also be used to break the semantic security of  $E$ . Finally, for a message  $m$  such that  $\Pr[f_r(C_m) = 1] \geq \frac{1}{2}$ , by definition of  $E'.\text{Enc}$ , it holds that

$$\Pr[f_r(E'.\text{Enc}(PK, m)) = 1] \geq 1 - \frac{1}{2^t}.$$

This shows that the output of  $f_r$  can be arbitrarily biased and concludes the proof.  $\square$

<sup>10</sup>Remark 2 discusses the sufficiency of proving impossibility for single-message schemes.

**Remark 3.** The essence of the above theorem is the impossibility of deterministic extraction: no single deterministic function can deterministically extract from ciphertexts of arbitrary encryption schemes. The way to bypass this impossibility is to have the extractor depend on the encryption scheme. Note that multiple-source extraction, which is used in our constructions in the subsequent sections, implicitly do depend on the underlying encryption scheme, since the additional sources of input depend on the encryption scheme and thus can be thought of as “auxiliary input” that is specific to the encryption scheme at hand.

## 4.2 Steganography for Adversarial Cover Distributions

Our second impossibility result concerns a much more general class of communication schemes, which we call *steganographic communication schemes*. Subliminal communication schemes, as well as the existing notions of public-key steganography and steganographic key exchange from the steganography literature, are instantiations of the more general definition of a (multi-message) steganographic communication scheme. To our knowledge, the general notion of a steganographic communication scheme has not been formalized in this way in prior work. In the context of this work, the general definition is helpful for proving broad impossibilities across multiple types of steganographic schemes.

As mentioned in the introduction, a limitation of all existing results in the steganographic literature, to our knowledge, is that they assume that the *cover distribution* — *i.e.*, the distribution of innocuous objects in which steganographic communication is to be embedded — is fixed *a priori*. In particular, the cover distribution is assumed not to depend on the description of the steganographic communication scheme. The impossibility result given in Section 4.1 is an example illustrative of the power of adversarially choosing the cover distribution: Theorem 4.1 says that by choosing the encryption scheme  $E$  to depend on a given subliminal communication scheme, an adversary can rule out the possibility of any hidden communication at all.

Our next impossibility result (Theorem 4.2) shows that if the cover distribution is chosen adversarially, then non-trivial steganographic communication is impossible.

**Theorem 4.2.** *Let  $\Pi$  be a steganographic communication scheme. Then for any  $k \in \mathbb{N}$ , there exists a cover distribution  $\mathcal{C}$  of conditional min-entropy  $k$  such that the steganographic indistinguishability of  $\Pi$  does not hold for more than one message.*

In Appendix A, we give the formal definition of a *steganographic communication scheme*, along with the proof of Theorem 4.2. We have elected to present these in the appendix as the definition introduces a set of new notation only used for the corresponding impossibility result, and both the definition and the impossibility result are somewhat tangential to the main results of this work, whose focus is on subliminal communication schemes.

## 5 Construction of the Subliminal Scheme

The goal of this section is to establish the following theorem, which states that our construction  $(\Phi^*, \Xi^*)$  is a subliminal communication scheme when instantiated with a pseudorandom key-exchange protocol (such as Diffie-Hellman).

**Theorem 5.1.** *The protocol  $(\Phi^*, \Xi^*)$  given in Definition 5.13, when instantiated with a pseudorandom key-exchange protocol  $\Lambda$ , is a multi-message subliminal communication scheme.*

The detailed description and proofs of security and correctness of our scheme can be found in the following subsections. Our construction makes no assumption on the message distribution  $\mathcal{M}$  and in particular holds when the exchanged plaintexts (of the adversarially mandated encryption scheme  $E$ ) are a fixed, adversarially chosen sequence of messages. An informal outline of the construction is given next.

**Definition 5.2.** Outline of the construction.

1. **Setup Phase  $\Phi^*$**

- (a) A  $\tilde{O}(\log \kappa)$ -bit string  $S$  is established between  $P_0$  and  $P_1$  by extracting randomness from pairs of consecutive ciphertexts. (*Protocol overview in Section 5.1.*)
- (b) Let  $\text{Ext}$  be a strong seeded extractor, and let  $S$  serve as its seed. By rejection-sampling ciphertexts  $c$  until  $\text{Ext}_S(c) = \text{str}$ , either party can embed a random string  $\text{str}$  of their choice in the conversation. (*Protocol overview in Section 5.2.*) By embedding in this manner the messages of a pseudorandom key-exchange protocol, both parties establish a shared secret  $\text{sk}^*$ .<sup>11</sup>

2. **Communication Phase  $\Xi^*$**

Both parties can now communicate arbitrary messages of their choice by (1) encrypting them using a pseudorandom secret-key encryption scheme  $\text{SKE}$  using  $\text{sk}^*$  as the secret key, and (2) embedding the ciphertexts of  $\text{SKE}$  using the rejection-sampling technique described in Step 1b.<sup>12</sup> (*Detailed protocol in Section 5.3.*)

The full protocol is given, and proven to be a subliminal communication scheme, in Section 5.3.

## 5.1 Establishing a Shared Seed

In this section, we give a protocol which allows  $P_0$  and  $P_1$  to establish a random public parameter which will be used in subsequent phases of our subliminal scheme. As such, this can be thought of as drawing a subliminal scheme at random from a family of subliminal schemes. The parameter is public in the sense that anyone eavesdropping on the channel between  $P_0$  and  $P_1$  gains knowledge of it. A crucial point is that the random draw occurs *after* the adversarial encryption scheme  $\text{E}$  is fixed, thus bypassing the impossibility results of Section 4.

Our strategy is simple: extract randomness from pairs of ciphertexts. Since the extractor does not receive the key, semantic security holds with respect to the extractor: a pair of ciphertexts for two arbitrary messages is indistinguishable from two encryptions of a fixed message; thus, a same-source extractor suffices for our purposes (see Lemma 5.6). Even though semantic security guarantees only  $\omega(\log \kappa)$  min-entropy of ciphertexts (see Lemma 5.4), we will be able to make use of the “greater-than” extractor (Definition 5.3) applied to pairs of ciphertexts, and obtain Theorem 5.8.

**Definition 5.3.** The *greater-than extractor*  $\text{GT}$  is defined by  $\text{GT}(x, y) = \mathbf{1}[x \geq y]$ .

**Lemma 5.4** (Ciphertexts have super-logarithmic min-entropy). *Let  $\text{PKE}$  be a semantically secure encryption scheme. Then there exists a negligible function  $\varepsilon$  such that for all  $\kappa \in \mathbb{N}$ ,  $m \in \mathcal{M}_\kappa$ , writing  $C_m^{pk} \sim \text{PKE.Enc}(pk, m)$ :*

$$\Pr \left[ (pk, sk) \leftarrow \text{PKE.Gen}(1^\kappa) : H_\infty(C_m^{pk}) \geq \log \frac{1}{\varepsilon(\kappa)} \right] \geq 1 - \varepsilon(\kappa).$$

*Proof.* In Appendix B. □

---

<sup>11</sup>Note that the random string  $\text{str}$  is known to an eavesdropper who has knowledge of the seed  $S$ . Nonetheless, (1) the established secret  $\text{sk}^*$  is unknown to the eavesdropper by the security of the key-exchange protocol and (2) the transcript is indistinguishable to the eavesdropper from one in which no key exchange occurred at all, due to the pseudorandomness of the key-exchange messages.

<sup>12</sup>Again, an eavesdropper could know the  $\text{SKE}$  ciphertexts exchanged, if he knew the seed  $S$ , but could not distinguish the  $\text{SKE}$  ciphertexts from truly random strings, and thus could not tell whether any subliminal communication was occurring at all. *Cf.* footnote 11.

Given that ciphertexts of semantically secure encryption schemes have min-entropy  $\omega(\log \kappa)$ , we will consider extractors which have negligible bias on such sources. This motivates the following definition.

**Definition 5.5.** Let  $2\text{Ext} : \{0, 1\}^n \times \{0, 1\}^m \rightarrow \{0, 1\}^\ell$  be a two-source extractor, we say that  $2\text{Ext}$  is an *extractor for super-logarithmic min-entropy* if  $2\text{Ext}$  is a  $(d \log \kappa, d \log \kappa, \frac{1}{\kappa^d})$  extractor for any  $d \in \mathbb{N}$ . In particular, for any negligible function  $\varepsilon$ , there exists a negligible function  $\varepsilon'$  such that  $2\text{Ext}$  is a  $(\log \frac{1}{\varepsilon}, \log \frac{1}{\varepsilon}, \varepsilon')$  extractor.

The following lemma shows that the output of a same-source extractor for super-logarithmic min-entropy on two ciphertexts is statistically indistinguishable from uniform, even in the presence of the key.

**Lemma 5.6.** Let  $\text{PKE}$  be a semantically secure encryption scheme with ciphertext length  $n$ , and let  $2\text{Ext} : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}^\ell$  be a same-source extractor for super-logarithmic min-entropy with  $\ell(\kappa) = O(\log \kappa)$ , then there exists a negligible function  $\varepsilon$  such that, for any  $\kappa \in \mathbb{N}$ ,  $(m_0, m_1) \in \mathcal{M}_\kappa^2$ , writing  $(PK, SK) \sim \text{PKE.Gen}(1^\kappa)$ ,  $C_i^{pk} \sim \text{PKE.Enc}(pk, m_i)$ ,  $i \in \{0, 1\}$ :

$$\left\| (PK, SK, 2\text{Ext}(C_0^{PK}, C_1^{PK})) - (PK, SK, U_{\ell(\kappa)}) \right\|_s \leq \varepsilon(\kappa).$$

*Proof.* We will prove that for any polynomial  $p$  and for large enough  $\kappa$ :

$$\Pr \left[ (pk, sk) \leftarrow \text{PKE.Gen}(1^\kappa) \quad : \quad \left\| 2\text{Ext}(C_0^{pk}, C_1^{pk}) - U_{\ell(\kappa)} \right\|_s \leq \frac{1}{p(\kappa)} \right] \geq 1 - \frac{1}{p(\kappa)}.$$

Assume by contradiction that there exists  $p \in \text{poly}(\kappa)$  and an infinite set  $I \subseteq \mathbb{N}$  such that for  $\kappa \in I$ :

$$\Pr \left[ (pk, sk) \leftarrow \text{PKE.Gen}(1^\kappa) \quad : \quad \left\| 2\text{Ext}(C_0^{pk}, C_1^{pk}) - U_{\ell(\kappa)} \right\|_s \geq \frac{1}{p(\kappa)} \right] \geq \frac{1}{p(\kappa)}. \quad (2)$$

We now construct an adversary  $D$  distinguishing between  $(PK, C_0^{PK})$  and  $(PK, C_1^{PK})$  with non-negligible advantage. On input  $(pk, c)$ ,  $D$  runs as follows:

1. Sample two encryptions of  $m_0$ :  $c_0, c'_0 \stackrel{iid}{\leftarrow} \text{PKE.Enc}(pk, m_0)$ , and  $c_1 \leftarrow \text{PKE.Enc}(pk, m_1)$ .
2. If  $2\text{Ext}(c_0, c_1) = 2\text{Ext}(c'_0, c)$  output 1, otherwise output 0.

First, note that on input  $(PK, C_1^{PK})$ ,  $D$  outputs 1 iff a collision occurs at step 2. By (2), with probability at least  $\frac{1}{p}$  over the draw of the  $(pk, sk)$ , a collision occurs with probability at least  $\frac{1}{2^\ell} + \frac{1}{2^\ell p^2}$ . Otherwise, a collision occurs with probability at least  $\frac{1}{2^\ell}$ . Overall, for  $\kappa \in I$ :

$$\Pr[D(PK, C_1^{PK}) = 1] \geq \frac{1}{2^\ell} + \frac{1}{2^\ell p^3}. \quad (3)$$

By Lemma 5.4, after conditioning on the event that  $H_\infty(C_0^{pk}) \geq \log q(\kappa)$ , the guarantee of  $2\text{Ext}$  applies to a pair of independent encryptions of  $m_0$  under  $pk$  and we obtain, for large enough  $\kappa$ :

$$\Pr \left[ (pk, sk) \leftarrow \text{PKE.Gen}(1^\kappa) \quad : \quad \left\| 2\text{Ext}(C_0^{pk}, C'_0^{pk}) - U_{\ell(\kappa)} \right\|_s \leq \frac{1}{q(\kappa)} \right] \geq 1 - \frac{1}{q(\kappa)},$$

This implies that for large enough  $\kappa$ :

$$\Pr[D(PK, C_0^{PK}) = 1] \leq \frac{1}{2^\ell} + \frac{2}{q}. \quad (4)$$

Together, (3) and (4) imply, after choosing  $q = 2^{\ell+2}p^3 \in \text{poly}$ , that for large enough  $\kappa \in I$ :

$$\Pr[\mathsf{D}(PK, C_0^{PK}) = 1] - \Pr[\mathsf{D}(PK, C_1^{PK}) = 1] \geq \frac{2}{q}.$$

This contradicts the security of PKE and concludes the proof.  $\square$

Finally, we observe that the “greater-than” extractor is a same-source extractor for super-logarithmic min-entropy (Lemma 5.7). To the best of our knowledge, this is a folklore fact which is for example mentioned in [BIW04].

**Lemma 5.7.** *For any  $k \leq n$ , GT is a  $(k, \frac{1}{2k})$  same-source extractor.*

*Proof.* In Appendix C.  $\square$

We now conclude this section with a full description of our method for establishing the public parameter  $S$  introduced in Step 1a.

**Theorem 5.8.** *Let  $\mathsf{E}$  be a semantically secure public-key encryption scheme and let  $\rho \in \text{poly}$ . Define random variables as follows.*

- For  $b \in \{0, 1\}$ , let  $K_b = (PK_b, SK_b) = \mathsf{E.Gen}(1^\kappa)$ .
- For  $b \in \{0, 1\}$  and  $i \in [2\rho]$ , let  $C_{b,i} = \mathsf{E.Enc}(PK_{1-b}, m_{b,i})$  representing the ciphertexts exchanged between  $P_0$  and  $P_1$  during  $2\rho$  exchange-rounds.
- Let  $S = (\text{GT}(C_{0,1}, C_{0,2}), \text{GT}(C_{1,1}, C_{1,2}), \dots, \text{GT}(C_{0,2\rho-1}, C_{0,2\rho}), \text{GT}(C_{1,2\rho-1}, C_{1,2\rho}))$ .

There exists a negligible function  $\varepsilon$  such that:

$$\|(K_0, K_1, S) - (K_0, K_1, U_{2\rho})\|_s \leq \varepsilon.$$

*Proof.* Writing  $S = (S_1, S'_1, \dots, S_\rho, S'_\rho)$ , we have:

$$\begin{aligned} \|(K_0, K_1, S) - (K_0, K_1, U_{2\rho})\|_s &\leq \sum_{i=1}^{\rho} (\|(K_0, K_1, i) - (K_0, K_1, U_i)\|_s \\ &\quad + \|(K_0, K_1, S'_i) - (K_0, K_1, S_i)\|_s) \leq 2\rho\varepsilon, \end{aligned}$$

where the first inequality follows by independence of the ciphertexts conditioned on the keys, and the second inequality follows by Lemma 5.6.  $\square$

**Remark 4.** In the construction of Theorem 5.8, the ciphertexts exchanged between  $P_0$  and  $P_1$  are sent without any modification, so subliminal indistinguishability clearly holds at this point.

## 5.2 Embedding Random Strings

In this section, we assume that both parties have access to a public parameter  $S$  and construct a protocol which allows for embedding of uniformly random strings into ciphertexts of an adversarially chosen encryption scheme  $\mathsf{E}$ , as required by Steps 1b and 2 of the construction outline (Definition 5.2). The security guarantee is that for a uniformly random parameter  $S$  and uniformly random strings to be embedded, the ciphertexts of  $\mathsf{E}$  with embedded random strings are indistinguishable from ciphertexts of  $\mathsf{E}$  produced by direct application of  $\mathsf{E.Enc}$ , even to an adversary who knows the decryption keys of  $\mathsf{E}$ . This can be thought of as a relaxation of subliminal indistinguishability (Definition 3.1) where the two main differences are that (1) the parties have shared knowledge of a random seed, and (2) indistinguishability only holds when

embedding a *random* string, rather than for arbitrary strings. We first present a construction to embed logarithmically many random bits (Theorem 5.9) and then show how to sequentially compose it to embed arbitrarily polynomially many random bits (Theorem 5.10). These constructions rely on a strong seeded extractor that can extract logarithmically many bits from sources of super-logarithmic min-entropy. Almost universal hashing is a simple such extractor, as stated in Proposition 5.11.

**Theorem 5.9.** *Let  $\text{Ext} : \{0, 1\}^d \times \{0, 1\}^n \rightarrow \{0, 1\}^v$  be a strong seeded extractor for super-logarithmic min-entropy with  $v = O(\log \kappa)$ , and let  $\mathbf{E}$  be a semantically secure encryption scheme with ciphertext space  $\mathcal{C} = \{0, 1\}^n$ . Let  $\Sigma^{\mathbf{E}, S}$  be defined as in Algorithm 1, then the following guarantees hold:*

---

**Algorithm 1** Rejection sampler  $\Sigma^{\mathbf{E}, S}$

---

PUBLIC PARAMETER:  $S$  (a  $d$ -bit seed).

INPUT:  $(\text{str}, m, \text{pk})$  where  $\text{str}$  is the string to be embedded.

1. Generate encryption  $c \leftarrow \mathbf{E}.\text{Enc}(\text{pk}, m)$ .
  2. If  $\text{Ext}(r, c) = \text{str}$ , then output  $c$ . Else, go back to step 1.
- 

1. Correctness: for any  $S \in \{0, 1\}^d$  and  $\text{str} \in \{0, 1\}^v$ , if  $c = \Sigma^{\mathbf{E}, S}(\text{str}, m, \text{pk})$ , and  $\text{str}' = \text{Ext}(S, c)$ , then  $\text{str}' = \text{str}$ .
2. Security: there exists a negligible function  $\varepsilon$  such that writing  $(PK, SK) = \mathbf{E}.\text{Gen}(1^\kappa)$ ,  $C = \mathbf{E}.\text{Enc}(PK, m)$  and  $C' = \Sigma^{\mathbf{E}, U_d}(U_v, m, PK)$ , the following holds:

$$\|(PK, SK, U_d, C) - (PK, SK, U_d, C')\|_{\mathfrak{s}} \leq \varepsilon(\kappa).$$

*Proof.* Define  $C'' = \mathbf{E}.\text{Enc}(PK_2, m)$ , an encryption of  $m$  independent of  $C$ . By definition of rejection sampling,  $C \sim \Sigma^{\mathbf{E}, S}(\text{Ext}(S, C''), m, PK_2)$ . Since  $\text{Ext}$  is a strong extractor for super-logarithmic min-entropy, and since  $C$  has super-logarithmic min-entropy, there exists a negligible function  $\varepsilon$  such that:

$$\|(PK_2, SK_2, U_d, \text{Ext}(U_d, C'')) - (PK_2, SK_2, U_d, U_v)\|_{\mathfrak{s}} \leq \varepsilon(\kappa).$$

The statistical distance can only decrease by applying  $\Sigma^{\mathbf{E}}$  on both sides, hence:

$$\|(PK_2, SK_2, U_d, C) - (PK_2, SK_2, U_d, C')\|_{\mathfrak{s}} \leq \varepsilon(\kappa).$$

which proves the security guarantee. Correctness is immediate.  $\square$

**Remark 5.** Rejection sampling is a simple and natural approach that has been used by prior work in the steganographic literature, such as [BC05]. Despite the shared use of this common technique, our construction is more different from prior art than it might seem at first glance. The novelty of our construction arises from the challenges of working in a model with a stronger adversary who can choose the distribution of ciphertexts (*i.e.*, the adversary gets to choose the public-key encryption scheme  $\mathbf{E}$ ). We manage to bypass the impossibilities outlined in Section 4 notwithstanding this stronger adversarial model, and in contrast to prior work, construct a protocol to established a shared seed from scratch, rather than simply assuming that one has been established in advance.

We now sequentially compose Theorem 5.9 to embed longer strings.

**Theorem 5.10.** *Let  $\Sigma$  be the rejection sampler defined in Algorithm 1. Let  $\ell \in \text{poly}$  and  $U_\ell$  be a uniformly random message of  $\ell$  bits. For  $v \leq \ell$ , we write  $U_\ell = U_{\ell,1} \parallel \dots \parallel U_{\ell,\nu}$  where  $U_{\ell,i}$  is a block of  $v$  bits from  $U_\ell$  and  $\nu = \frac{\ell}{v}$ . Given cover messages  $m_1, \dots, m_\ell$ , define  $(PK, SK) = \text{E.Gen}(1^\kappa)$ ,  $C'_i = \Sigma^{\text{E}, U_d}(U_{\ell,i}, m_i, \text{pk})$ ,  $C_i = \text{E.Enc}(PK, m_i)$ , then there exists a negligible function  $\varepsilon$  such that:*

$$\left\| (PK, SK, U_d, (C_i)_{i \in [\nu]}) - (PK, SK, U_d, (C'_i)_{i \in [\nu]}) \right\|_{\mathfrak{s}} \leq \varepsilon(\kappa).$$

*Proof.* Define  $K = (PK, SK)$ , then:

$$\left\| (K, U_d, (C_i)_{i \in [\nu]}) - (K, U_d, (C'_i)_{i \in [\nu]}) \right\|_{\mathfrak{s}} \leq \sum_{i=1}^{\nu} \left\| (K, C_i) - (K, C'_i) \right\|_{\mathfrak{s}} \leq \nu \varepsilon,$$

where the first inequality is by independence of the sequences  $(C_i)_{i \in [\nu]}$  and  $(C'_i)_{i \in [\nu]}$  conditioned on the keys, and the second inequality is by Theorem 5.9.  $\square$

Finally, we observe that almost universal hashing is a strong seeded extractor for super-logarithmic min-entropy which has negligible error when the output length is  $O(\log(\kappa))$  (Proposition 5.11). This exactly satisfies the requirement of Theorem 5.9. Moreover, the seed length of this extractor is only super-logarithmic, meaning that the seed can be established, in Step 1a of the Setup Phase (Definition 5.2), in  $\tilde{O}(\log \kappa)$  many exchange-rounds of communication.

**Proposition 5.11.** *Let  $\delta$  be a negligible function and let  $\mathcal{H}$  be a family of  $\delta$ -almost pairwise independent hash functions mapping  $\{0, 1\}^n$  to  $\{0, 1\}^{c \cdot \log n}$ , then the extractor  $\text{Ext} : \mathcal{H} \times \{0, 1\}^n \rightarrow \{0, 1\}^{c \cdot \log n}$  defined by  $\text{Ext}(h, x) = h(x)$  is a strong seeded extractor for super-logarithmic min-entropy. Furthermore, there exists an explicit family  $\mathcal{H}$  of such hash functions, such that sampling uniformly from  $\mathcal{H}$  requires  $O(c \cdot \log n + \log \frac{1}{\delta})$  bits.*

*Proof.* See [SZ94].  $\square$

### 5.3 Full Protocol $(\Phi^*, \Xi^*)$

First, we establish some notation for the syntax of a key-exchange protocol.

**Definition 5.12** (Key-exchange protocol syntax). A key-exchange protocol is a two-party protocol defined by

$$\Lambda = (\Lambda_{0,1}, \Lambda_{1,1}, \Lambda_{0,2}, \Lambda_{1,2}, \dots, \Lambda_{0,k}, \Lambda_{1,k}, \Lambda_{0,\text{out}}, \Lambda_{1,\text{out}}).$$

We assume  $k$  simultaneous communication rounds, where  $\Lambda_{b,i}$  represents the computation performed by  $P_b$  in the  $i$ th round. The parties are stateful and their state is implicitly updated at each round to contain the transcript so far and any local randomness generated so far. Each  $\Lambda_{b,i}$  takes as input the transcript up to round  $i - 1$  and the state of  $P_b$ , and outputs a message  $\lambda_{b,i}$  to be sent in the  $i$ th round. For notational simplicity, we write explicitly only the first input to  $\Lambda_{b,i}$ , and leave the second input (*i.e.*, the state) implicit.  $\Lambda_{0,\text{out}}, \Lambda_{1,\text{out}}$  are run by  $P_0, P_1$  respectively to compute the shared secret at the conclusion of the protocol.

Next, we give the full construction of  $(\Phi^*, \Xi^*)$  following the outline in Definition 5.2.

**Definition 5.13.**  $(\Phi^*, \Xi^*)$  is parametrized by the following.

- $\text{Ext} : \{0, 1\}^d \times \{0, 1\}^n \rightarrow \{0, 1\}^v$ , a strong seeded extractor.
- $\Lambda$ , a pseudorandom key-exchange protocol with  $\ell$ -bit messages.<sup>13</sup>

<sup>13</sup>In presenting our construction  $(\Phi^*, \Xi^*)$ , we do not denote the state of parties w.r.t. the key-exchange protocol  $\Lambda$  by a separate variable, but assume that it is part of the state  $\mathfrak{s}_{b,i}$  of the overall protocol.

- SKE, a pseudorandom secret key encryption scheme with  $\xi$ -bit ciphertexts.

We define each phase of our construction in turn.

## 1. Setup Phase $\Phi^*$

### (a) Establishing a $d$ -bit shared seed

- For  $b \in \{0, 1\}$  and  $i \in \{1, \dots, d\}$ ,  $\Phi_{b,i}^*(m_{b,i}, \text{pk}_{1-b}, \mathfrak{s}_{b,i-1})$  outputs a ciphertext  $c_{b,i} = \text{E.Enc}(\text{pk}_{1-b}, m_{b,i})$  and sets the updated state  $\mathfrak{s}_{b,i}$  to be the transcript of all protocol messages sent and received so far.
- At the conclusion of the  $d$  exchange-rounds, each party updates his state to contain the seed  $S$  which is defined by

$$S = (\text{GT}(c_{0,1}, c_{0,2}), \text{GT}(c_{1,1}, c_{1,2}), \dots, \text{GT}(c_{0,d-1}, c_{0,d}), \text{GT}(c_{1,d-1}, c_{1,d})) .$$

This seed  $S$  is assumed to be accessible in all future states throughout both phases during the remainder of the protocol.

### (b) Subliminal key exchange

Let  $\nu = \frac{\xi}{v}$ . Subliminal key exchange occurs over  $k \cdot \nu$  exchange-rounds.

- For  $j \in \{1, \dots, k\}$  and  $b \in \{0, 1\}$ :
  - $P_b$  retrieves from his state the key-exchange transcript so far  $(\lambda_{b,j'})_{b \in \{0,1\}, j' < j}$ .
  - $P_b$  computes the next key-exchange message

$$\lambda_{b,j} \leftarrow \Lambda_{b,j}((\lambda_{b,j'})_{b \in \{0,1\}, j' < j}) .$$

- $P_b$  breaks  $\lambda_{b,j}$  into  $v$ -bit blocks  $\lambda_{b,j} = \lambda_{b,j}^1 || \dots || \lambda_{b,j}^\nu$ .
- The  $\nu$  blocks are transmitted sequentially as follows. For  $\iota \in \{1, \dots, \nu\}$ :

Let  $i = d + (j - 1)\nu + \iota$ .

$\Phi_{b,i}^*(m_{b,i}, \text{pk}_{1-b}, \mathfrak{s}_{b,i-1})$  outputs  $c_{b,i} \leftarrow \Sigma^{\text{E}, S}(\lambda_{b,j}^\iota, m_{b,i}, \text{pk}_{1-b})$  and sets the updated state  $\mathfrak{s}_{b,i}$  to contain the transcript of all protocol messages sent and received so far.

- At the conclusion of the  $\iota$  exchange-rounds, each party  $b \in \{0, 1\}$  updates his state to contain the  $j$ th key-exchange message  $\lambda_{1-b,j}$  computed as follows:

$$\lambda_{1-b,j} = \text{Ext}(S, c_{b,d+(j-1)\nu+1}) || \dots || \text{Ext}(S, c_{b,d+j\nu})$$

- At the conclusion of the  $k \cdot \nu$  exchange rounds, each party updates his state to contain the secret key  $\text{sk}^*$  computed as follows:

$$\text{sk}^* = \text{SKE.Gen}(1^\kappa; \Lambda_{\text{out}}((\lambda_{b,j})_{b \in \{0,1\}, j \in [k]})) .$$

## 2. Communication Phase $\Xi^*$

Each communication phase occurs over  $r' = \xi/v$  exchange-rounds.

Let  $\beta \in \{0, 1\}$  be the initiating party and let  $\bar{\beta} = 1 - \beta$ .

$P_\beta$  performs the following steps.

- Generate  $c^* \leftarrow \text{SKE.Enc}(\text{sk}^*, \text{msg})$ .
- Break  $c^*$  into  $v$ -bit blocks  $c^* = c_1^* || \dots || c_{r'}^*$ .

For  $i' \in \{1, \dots, r'\}$ :

- Let  $i'' = r + i'$ .

- 30 •  $\Xi_{0,i'}^*(\text{msg}, m_{0,i'}, \text{pk}_{\beta}, \mathfrak{s}_{\beta,i'-1})$  outputs  $c_{\beta,i'} \leftarrow \Sigma^{E,S}(c_{i'}, m_{\beta,i'}, \text{pk}_{\beta})$ .
- 31 •  $\Xi_{1,i'}^*(m_{\beta,i'}, \text{pk}_{\beta}, \mathfrak{s}_{\beta,i'-1})$  outputs  $c_{\beta,i'} \leftarrow \text{E.Enc}(\text{pk}_{\beta}, m_{\beta,i'})$ .
- 32 • Both parties update their state to contain the transcript of all protocol messages
- 33 exchanged so far.

34 After the  $r'$  exchange-rounds,  $P_{\beta}$  computes  $c^{**}$  as follows:

$$c^{**} = \text{Ext}(S, c_{\beta,r+1}) || \dots || \text{Ext}(S, c_{\beta,r+r'}) .$$

35 Then,  $P_{\beta}$  outputs  $\text{msg}' \leftarrow \text{SKE.Dec}(\text{sk}^*, c^{**})$ . (That is,  $\Xi_{1,\text{out}}^*(\mathfrak{s}_{\beta,r'}) = \text{msg}'$ .)

#### 5.4 Proof that $(\Phi^*, \Xi^*)$ Is a Subliminal Communication Scheme

Finally, we give the proof of our main theorem. We recall the statement here.

**Theorem 5.1.** Assume there exists a pseudorandom key-exchange protocol. Then there is a multi-message subliminal communication scheme  $(\Phi^*, \Xi^*)$ , given in Definition 5.13.

*Proof.* We define three hybrids.

- HYBRID 0 (“REAL WORLD”): Parties execute  $(\Phi^*, \Xi^*)$ .
- HYBRID 1: Exactly like Hybrid 0, except that the seed  $S$  in Phase 1a is replaced by a truly random  $d$ -bit string (the same string for both parties).
- HYBRID 2: Exactly like Hybrid 1, except that the key exchange messages  $\lambda_{b,j}$  in  $\Phi^*$  are replaced by random strings. That is, Line 12 is replaced by:

$$P_b \text{ samples } \lambda_{b,j} \leftarrow \{0, 1\}^{\ell} \text{ at random.}$$

- HYBRID 3: Exactly like Hybrid 2, except that the ciphertexts of SKE in  $\Xi^*$  are replaced by random strings. That is, Line 26 is replaced by:

$$\text{Sample } c^* \leftarrow \{0, 1\}^{\xi} \text{ at random.}$$

Hybrids 0 and 1 are indistinguishable by direct application of Theorem 5.8.

Hybrids 1 and 2 are indistinguishable by the pseudorandomness of the key-exchange protocol (defined in Section 2.1). Note that it is essential that the  $\Lambda$ -transcript’s indistinguishability from random holds even in the presence of the established key  $\text{sk}^*$ , since in our protocol  $(\Phi^*, \Xi^*)$ , the later protocol messages are produced as a function of  $\text{sk}^*$ .

Hybrids 2 and 3 are indistinguishable because ciphertexts of SKE are pseudorandom in the absence of the corresponding secret key  $\text{sk}^*$ . Because we already replaced the messages of  $\Lambda$  with random messages independent of  $\text{sk}^*$ , the distribution of all protocol messages in Hybrid 1 can be generated based on just the SKE-ciphertexts  $c^*$  (of line 26).

Finally, Hybrid 3 is indistinguishable from the ideal distribution

$$\text{Ideal}(\text{pk}_1, \text{sk}_1, \text{pk}_2, \text{sk}_2, \mathcal{M})$$

from Definition 3.3 by Theorem 5.10. Note that the rejection sampler  $\Sigma^{E,U_d}$  of Algorithm 1 has a truly random  $d$ -bit seed, so to invoke Theorem 5.10 we rely on the fact that  $S$  is truly random in Hybrid 3.  $\square$

## 5.5 On the Setup Cost and Asymptotic Rate of $(\Phi^*, \Xi^*)$

**Setup Cost.** The setup cost of our scheme can be broken down into the costs of Step 1a and Step 1b as follows.

- *Step 1a:* If our scheme is instantiated with the extractor Ext from Proposition 5.11, then we need to establish a seed of length  $\tilde{O}(\log \kappa)$ , which implies that  $\tilde{O}(\log \kappa)$  exchange-rounds are required in Step 1a. This is arguably the least efficient step in our scheme; this inefficiency stems from the use of the GT extractor which only outputs one bit: to the best of our knowledge this the only extractor which applies to our setting. In Section 6, we discuss ways in which this cost can be reduced under additional assumptions on the next-message distribution or the encryption scheme E by replacing GT by extractors with longer outputs.
- *Step 1b:* The cost of this step is  $k \cdot \frac{\ell}{v}$ , where  $k$  is the number of rounds of the key-exchange protocol  $\Lambda$  that we use,  $\ell$  is the length of messages in  $\Lambda$  and  $v$  is the output length of Ext. (Concretely, the Diffie-Hellman key exchange protocol achieves  $k = 1$  and  $\ell = O(\kappa)$ .) If we use the extractor from Proposition 5.11 as Ext, then we can achieve  $v = c \log \kappa$  for any  $c > 0$ . This implies that the cost of Step 1b is upper-bounded by  $c' \frac{\kappa}{\log \kappa}$  for any  $c' > 0$ . Note that because the min-entropy of ciphertexts from E can be as small as  $\omega(\log \kappa)$  and is *a priori* unknown to the designer of the subliminal scheme, the output of Ext must be  $O(\log \kappa)$ .<sup>14</sup> The extractor Ext from Proposition 5.11 is optimal in this respect.

**Asymptotic Rate.** The asymptotic rate of our scheme depends on the output length of Ext. Using the extractor from Proposition 5.11, we can subliminally embed  $c \log \kappa$  random bits per ciphertext of E and hence achieve an asymptotic rate of  $c \log \kappa$  for any  $c > 0$ . As in the discussion regarding the cost of Step 1b, this is optimal given that the min-entropy of ciphertexts can be as small as  $\omega(\log \kappa)$ .

**Trade-off Between Running Time and Rate.** Note that the parameter  $c$  from the previous paragraph controls the trade-off between the running time of our scheme and its asymptotic rate. Indeed, the expected running time of the rejection sampler defined in Algorithm 1 is  $O(\kappa^c)$ , when embedding  $c \log \kappa$  random bits. This trade-off is inherent to rejection sampling, and it is an interesting open question to determine whether it can be improved by an alternative technique.

## 6 Improving Setup Cost

In this section, we present alternative constructions which improve the setup cost of subliminal communication under additional assumptions, either on the next-message distribution  $\mathcal{M}$ , or on the public-key encryption scheme E. These additional assumptions allow us to replace the “greater-than” extractor in Step 1a of Definition 5.2 with extractors with longer output, thus reducing the number of exchange-rounds required to establish the shared seed. Section 6.1 gives a construction when E is “succinct” (*i.e.*, has a constant expansion factor). Section 6.2 presents a construction when  $\mathcal{M}$  has a known amount of min-entropy. Both these constructions yield a seed establishment in two exchange-rounds.

<sup>14</sup>Suppose not, *i.e.*, suppose that the output of Ext were  $v \in \omega(\log \kappa)$  bits long. Then the adversary could choose an encryption scheme whose ciphertexts have min-entropy  $z \in \omega(\log \kappa) \cap o(v)$  (*e.g.*,  $z = \sqrt{v(\kappa) \log \kappa}$ ). Since the extractor output cannot have more min-entropy than its input, the extractor’s output when evaluated on ciphertexts would be distinguishable from random  $v$ -bit strings.

## 6.1 Succinct Encryption Schemes

Let us suppose that the adversarially chose scheme  $E$  is *succinct*. Here we define *succinct* as having an expansion factor less than 2. Recall that the expansion factor is defined to be the ratio of ciphertext length to plaintext length. Under this assumption, we can improve the number of rounds in the construction of Section 5.1 by replacing the extractor  $GT$  by the extractor  $BLE$  from [DEOR04]. Note that this extractor requires the min-entropy rate of the sources to be slightly above  $\frac{1}{2}$ , yet ciphertexts of  $E$  only have min-entropy  $\omega(\log \kappa)$ . However, the succinctness assumption combined with semantic security implies that ciphertexts have sufficiently large HILL entropy: they are computationally indistinguishable from sources of min-entropy rate slightly above  $\frac{1}{2}$ . Since the extractor is a polynomial time algorithm, its output when computed on ciphertexts from  $E$  will be computationally indistinguishable from the uniform distribution. Formally, we prove the following theorem.

**Theorem 6.1.** *Let us denote by  $n$  (resp.  $p$ ) the bit-lengths of ciphertexts (resp. plaintexts) from  $E$ . Let  $BLE$  be the function constructed in [DEOR04] and let us denote by  $B : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}^v$  its truncation to  $v = O(\log \kappa)$  bits. With the same notations as Theorem 5.8, define:*

$$S = (B(C_{0,1}, C_{0,2}), B(C_{1,1}, C_{1,2}), \dots, B(C_{0,2\rho-1}, C_{0,2\rho}), B(C_{1,2\rho-1}, C_{1,2\rho})).$$

*There is a negligible function  $\varepsilon$  such that  $\|(K_0, K_1, S) - (K_0, K_1, U_{2\rho v})\|_s \leq \varepsilon$ .*

*Proof.* For any  $\gamma > 0$ , the function  $BLE$  from [DEOR04] can extract  $2p - n - \gamma p$  random bits with bias  $\varepsilon = 2^{-1/2\gamma p}$  from two independent  $n$ -bit sources of min-entropy  $p$ . Note that for  $n \leq (1 - \gamma)2p$ ,  $2p - n - \gamma p \geq \gamma p$ , and by semantic security of  $E$ , we have that  $n$  (and hence  $p$ ) are super-logarithmic. Hence the function  $B$  defined in Theorem 6.1 is well-defined and has bias  $\varepsilon \in \text{negl}$  for independent sources of min-entropy  $p$ .

For public key  $pk$ , define  $Y^{pk} = E.\text{Enc}(pk, U_p)$  the random variable obtained by encrypting a message chosen uniformly at random. By correctness of  $E$ ,  $Y^{pk}$  is a flat source supported on  $\{0, 1\}^p$ , hence  $H_\infty(Y^{pk}) = p$ . Hence for  $Z^{pk}$  an independent copy of  $Y^{pk}$ , we have  $\|B(Y^{pk}, Z^{pk}) - U_v\|_s \leq \varepsilon(\kappa)$ . By semantic indistinguishability of  $E$ , we have that for any ciphertext  $C$  of  $E$ ,  $C \stackrel{c}{\approx} Y^{pk}$ ; *i.e.*,  $C$  has HILL-entropy  $p$ . Adapting the proof of Lemma 5.6 and using the same notation, since  $v = O(\log \kappa)$ , we get  $\|(PK, SK, B(C_0^{PK}, C_1^{PK})) - (PK, SK, U_v)\|_s \leq \varepsilon$ . From there, the proof of Theorem 6.1 follows verbatim the proof of Theorem 5.8.  $\square$

**Remark 6.** Theorem 6.1 implies in particular that a  $\tilde{O}(\log \kappa)$  random seed can be established in step 1a of our subliminal scheme in only two exchange-rounds of communication, whereas the construction of Theorem 5.8 with the  $GT$  extractor requires  $\tilde{O}(\log \kappa)$  exchange-rounds.

## 6.2 Next-Message Distributions with Min-Entropy

The previous subsection reduced the number of rounds to establish the shared seed by using a two-source extractor with more than one bit of output (unlike the  $GT$  extractor). However, this required  $E$  to be succinct to guarantee that ciphertexts of  $E$  have sufficient HILL-entropy. In this section, we observe that if the message distribution itself has enough min-entropy, then we can use the *plaintext* and corresponding *ciphertext* as a pair of sources to extract from, and obtain a similar improvement without requiring that  $E$  be succinct. Note that this construction again exploits the semantic security of  $E$  to guarantee that the plaintext and the ciphertext are indistinguishable from independent sources to the two-source extractor.

All we require of  $E$  in this subsection is that the ciphertext distribution has min-entropy at least  $\omega(\log \kappa)$ , which follows from semantic security (Lemma 5.4); however, we additionally require that the next message distribution  $\mathcal{M}$  have min-entropy rate above  $\frac{1}{2}$ . While this

is quite a lot of min-entropy to demand from  $\mathcal{M}$ , we note that the precise requirement of  $\frac{1}{2}$  arises from the current state of the art in two-source extractors, and improved two-source extractor constructions would directly imply improvements in the min-entropy requirement of our construction. Recent research in two-source extraction has been quite productive; with luck, future advances will provide us with an alternative which demands much less entropy from  $\mathcal{M}$ .

**Theorem 6.2.** *Suppose that  $k_1$  is such that for all  $k_2 = \omega(\log \kappa)$  there exists a negligible function  $\varepsilon$  and a  $(k_1, k_2, \varepsilon)$ -two source extractor  $2\text{Ext}$  with output length  $\ell = O(\log \kappa)$ . and let  $\rho \in \text{poly}$ . Define random variables as follows.*

- For  $b \in \{0, 1\}$ , let  $K_b = (PK_b, SK_b) = \text{E.Gen}(1^\kappa)$ .
- For  $b \in \{0, 1\}$  and  $i \in [\rho]$ , let  $M_{b,i}$  and  $C_{b,i} = \text{E.Enc}(PK_{1-b}, m_{b,i})$  be the messages and ciphertexts exchanged between  $P_0$  and  $P_1$  in  $\rho$  exchange-rounds.
- $S = (2\text{Ext}(M_{0,1}, C_{0,1}), 2\text{Ext}(M_{1,1}, C_{1,1}), \dots, 2\text{Ext}(M_{0,\rho}, C_{0,\rho}), 2\text{Ext}(M_{1,\rho}, C_{1,\rho}))$ .

Then if the next message distribution  $\mathcal{M}(\text{conv})$  has min-entropy at least  $k_1$ :

$$\|(K_0, K_1, S) - (K_0, K_1, U_\rho)\|_s \leq \varepsilon.$$

*Proof.* Let us consider  $M$  and  $M'$  two independent samples from the next-message distribution and let us denote by  $C_M^{PK}$  and  $C_{M'}^{PK}$  their encryption under key  $PK$ . Consider a two-source extractor  $2\text{Ext}$  as in the statement of the theorem. By semantic security, it follows that for some negligible function  $\varepsilon'$ :

$$2\text{Ext}(M, C_M^{PK}) \stackrel{c}{\approx}_{\varepsilon'} 2\text{Ext}(M', C_{M'}^{PK}).$$

By property of  $2\text{Ext}$  it follows that  $\|2\text{Ext}(M', C_{M'}^{PK}) - U_\ell\|_s \leq \varepsilon$ . Since  $\ell = O(\log \kappa)$ , we can prove similarly to Lemma 5.6:

$$\|(PK, SK, 2\text{Ext}(M, C_M^{PK})) - (PK, SK, U_\ell)\|_s \leq \varepsilon + \varepsilon'.$$

We can now conclude similarly to the proof of Theorem 5.8.  $\square$

The following result from [Raz05] implies that for any  $\delta > 0$ , when  $k_1 = \frac{1}{2} + \delta$ , two-source extractors exist which may be used in the above theorem. The output length of such extractors is logarithmic, implying that a  $\tilde{O}(\log \kappa)$ -bit random seed  $S$  can be established in two exchange-rounds of communication.

**Lemma 6.3** (Theorem 1 in [Raz05]). *For any  $\delta > 0$ ,  $k_1 \geq \frac{1}{2} + \delta$  and  $k_2 = \omega(\log \kappa)$ , there exists a negligible function  $\varepsilon$  and a  $(k_1, k_2, \varepsilon)$ -two source extractor with output length  $\ell = O(\log \kappa)$ .*

## 7 Open problems

**Deterministic Extraction.** Our impossibility result in Theorem 4.1 holds because the adversary can choose the encryption scheme  $\text{E}$  as a function of a given candidate subliminal scheme. However, note that under the additional assumption that  $\text{E}$  is restricted to a predefined class  $\mathcal{C}$  of encryption schemes, we could bypass this impossibility as long as a deterministic extractor that can extract randomness from ciphertexts of any encryption scheme in  $\mathcal{C}$  exists. We are only aware of two deterministic extractors leading to a positive result for restricted classes of encryption schemes:

- if an upper bound on the circuit size of  $\text{E}$  is known, then we can use the deterministic extractor from [TV00]. This extractor relies on strong complexity-theoretic assumptions and requires the sources to have min-entropy  $(1 - \gamma)n$  for some unspecified constant  $\gamma$ .

- if  $E$  is computed by a circuit of constant depth ( $\mathbf{AC}^0$ ), then the deterministic extractor of [Vio11] can be used and requires  $\sqrt{n}$  min-entropy.

Note that both these extractors have a min-entropy requirement which is too strict to be directly applicable to ciphertexts of arbitrary encryption schemes, but it would be interesting to give improved constructions for the specific case of encryption schemes. This would also have direct implications for the efficiency of the subliminal scheme we construct in Section 5.2: indeed, one could then skip Step 1a and use a deterministic extractor directly in Steps 1b and 2, thus saving  $\tilde{O}(\log \kappa)$  exchange-rounds in the setup phase.

**Multi-Source Extraction.** Another interesting question is whether multi-source extractors for the specific case when the sources are independent and identically distributed can achieve better parameters than extractors for general independent sources. We already saw that a very simple extractor (namely, the “greater-than” function) works for i.i.d. sources and extracts one bit with negligible bias, even when the sources only have  $\omega(\log \kappa)$  min-entropy. The non-constructive result of [CG88] guarantees the existence of a two-source extractor of negligible bias and output length  $\omega(\log \kappa)$  for sources of min-entropy  $\omega(\log \kappa)$ . However, known *explicit* constructions are far from achieving the same parameters, and improving them in the specific case of identically distributed sources is an interesting open problem which was also mentioned in [BIW04].

## Acknowledgements

We are grateful to Omer Paneth and Adam Sealfon for insightful remarks on an earlier draft of this paper.

## References

- [ACM<sup>+</sup>14] Per Austrin, Kai-Min Chung, Mohammad Mahmoody, Rafael Pass, and Karn Seth. On the impossibility of cryptography with tamperable randomness. In Garay and Gennaro [GG14], pages 462–479.
- [Auc98] David Aucsmith, editor. *Information Hiding, Second International Workshop, Portland, Oregon, USA, April 14-17, 1998, Proceedings*, volume 1525 of *Lecture Notes in Computer Science*. Springer, 1998.
- [BC05] Michael Backes and Christian Cachin. Public-key steganography with active attacks. In Joe Kilian, editor, *Theory of Cryptography, Second Theory of Cryptography Conference, TCC 2005, Cambridge, MA, USA, February 10-12, 2005, Proceedings*, volume 3378 of *Lecture Notes in Computer Science*, pages 210–226. Springer, 2005.
- [BIW04] Boaz Barak, Russell Impagliazzo, and Avi Wigderson. Extracting randomness using few independent sources. In *45th Symposium on Foundations of Computer Science (FOCS 2004), 17-19 October 2004, Rome, Italy, Proceedings*, pages 384–393, 2004.
- [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, *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, Denver, CO, USA, October 12-6, 2015*, pages 1431–1440. ACM, 2015.

- [BL17] Sebastian Berndt and Maciej Liškiewicz. Algorithm substitution attacks from a steganographic perspective. In Bhavani M. Thuraisingham, David Evans, Tal Malkin, and Dongyan Xu, editors, *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, CCS 2017, Dallas, TX, USA, October 30 - November 03, 2017*, pages 1649–1660. ACM, 2017.
- [BPR14] Mihir Bellare, Kenneth G. Paterson, and Phillip Rogaway. Security of symmetric encryption against mass surveillance. In Garay and Gennaro [GG14], pages 1–19.
- [Cac98] Christian Cachin. An information-theoretic model for steganography. In Aucsmith [Auc98], pages 306–318.
- [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.
- [CK16] Aloni Cohen and Saleet Klein. The GGM function family is a weakly one-way family of functions. In Martin Hirt and Adam D. Smith, editors, *Theory of Cryptography - 14th International Conference, TCC 2016-B, Beijing, China, October 31 - November 3, 2016, Proceedings, Part I*, volume 9985 of *Lecture Notes in Computer Science*, pages 84–107, 2016.
- [CZ16] Eshan Chattopadhyay and David Zuckerman. Explicit two-source extractors and resilient functions. In *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2016, Cambridge, MA, USA, June 18-21, 2016*, pages 670–683, 2016.
- [DEOR04] Yevgeniy Dodis, Ariel Elbaz, Roberto Oliveira, and Ran Raz. Improved randomness extraction from two independent sources. In *7th International Workshop on Approximation Algorithms for Combinatorial Optimization Problems, APPROX 2004, and 8th International Workshop on Randomization and Computation, RANDOM 2004, Cambridge, MA, USA, August 22-24, 2004, Proceedings*, pages 334–344, 2004.
- [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, *Advances in Cryptology - EUROCRYPT 2015 - 34th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Sofia, Bulgaria, April 26-30, 2015, Proceedings, Part I*, volume 9056 of *Lecture Notes in Computer Science*, pages 101–126. Springer, 2015.
- [DH76] Whitfield Diffie and Martin E. Hellman. New directions in cryptography. *IEEE Trans. Information Theory*, 22(6):644–654, 1976.
- [GG14] Juan A. Garay and Rosario Gennaro, editors. *Advances in Cryptology - CRYPTO 2014 - 34th Annual Cryptology Conference, Santa Barbara, CA, USA, August 17-21, 2014, Proceedings, Part I*, volume 8616 of *Lecture Notes in Computer Science*. Springer, 2014.
- [GGM86] Oded Goldreich, Shafi Goldwasser, and Silvio Micali. How to construct random functions. *J. ACM*, 33(4):792–807, 1986.
- [Gol11] Oded Goldreich. *Studies in complexity and cryptography*. chapter The GGM Construction Does NOT Yield Correlation Intractable Function Ensembles, pages 98–108. Springer-Verlag, Berlin, Heidelberg, 2011.

- [HLv02] Nicholas J. Hopper, John Langford, and Luis von Ahn. Provably secure steganography. In Moti Yung, editor, *Advances in Cryptology - CRYPTO 2002, 22nd Annual International Cryptology Conference, Santa Barbara, California, USA, August 18-22, 2002, Proceedings*, volume 2442 of *Lecture Notes in Computer Science*, pages 77–92. Springer, 2002.
- [Mic16] Silvio Micali. ALGORAND: the efficient and democratic ledger. *CoRR*, abs/1607.01341, 2016.
- [Mit99] Thomas Mittelholzer. An information-theoretic approach to steganography and watermarking. In Andreas Pfitzmann, editor, *Information Hiding, Third International Workshop, IH'99, Dresden, Germany, September 29 - October 1, 1999, Proceedings*, volume 1768 of *Lecture Notes in Computer Science*, pages 1–16. Springer, 1999.
- [Raz05] Ran Raz. Extractors with weak random seeds. In Harold N. Gabow and Ronald Fagin, editors, *Proceedings of the 37th Annual ACM Symposium on Theory of Computing, Baltimore, MD, USA, May 22-24, 2005*, pages 11–20. ACM, 2005.
- [Sim83] Gustavus J. Simmons. The prisoners’ problem and the subliminal channel. In David Chaum, editor, *Advances in Cryptology, Proceedings of CRYPTO '83, Santa Barbara, California, USA, August 21-24, 1983.*, pages 51–67. Plenum Press, New York, 1983.
- [SZ94] Aravind Srinivasan and David Zuckerman. Computing with very weak random sources. In *35th Annual Symposium on Foundations of Computer Science, Santa Fe, New Mexico, USA, 20-22 November 1994*, pages 264–275. IEEE Computer Society, 1994.
- [TV00] Luca Trevisan and Salil P. Vadhan. Extracting randomness from samplable distributions. In *41st Annual Symposium on Foundations of Computer Science, FOCS 2000, 12-14 November 2000, Redondo Beach, California, USA*, pages 32–42, 2000.
- [Vad12] Salil P. Vadhan. Pseudorandomness. *Foundations and Trends in Theoretical Computer Science*, 7(13):1–336, 2012.
- [vAH04] Luis von Ahn and Nicholas J. Hopper. Public-key steganography. In Christian Cachin and Jan Camenisch, editors, *Advances in Cryptology - EUROCRYPT 2004, International Conference on the Theory and Applications of Cryptographic Techniques, Interlaken, Switzerland, May 2-6, 2004, Proceedings*, volume 3027 of *Lecture Notes in Computer Science*, pages 323–341. Springer, 2004.
- [Vio11] Emanuele Viola. Extractors for circuit sources. In Rafail Ostrovsky, editor, *IEEE 52nd Annual Symposium on Foundations of Computer Science, FOCS 2011, Palm Springs, CA, USA, October 22-25, 2011*, pages 220–229. IEEE Computer Society, 2011.
- [YY96a] Adam L. Young and Moti Yung. Cryptovirology: Extortion-based security threats and countermeasures. In *1996 IEEE Symposium on Security and Privacy, May 6-8, 1996, Oakland, CA, USA*, pages 129–140. IEEE Computer Society, 1996.
- [YY96b] Adam L. Young and Moti Yung. The dark side of ”black-box” cryptography, or: Should we trust capstone? In Neal Koblitz, editor, *Advances in Cryptology - CRYPTO '96, 16th Annual International Cryptology Conference, Santa Barbara, California, USA, August 18-22, 1996, Proceedings*, volume 1109 of *Lecture Notes in Computer Science*, pages 89–103. Springer, 1996.

- [YY97] Adam L. Young and Moti Yung. Kleptography: Using cryptography against cryptography. In Walter Fumy, editor, *Advances in Cryptology - EUROCRYPT '97, International Conference on the Theory and Application of Cryptographic Techniques, Konstanz, Germany, May 11-15, 1997, Proceeding*, volume 1233 of *Lecture Notes in Computer Science*, pages 62–74. Springer, 1997.
- [ZFK<sup>+</sup>98] Jan Zöllner, Hannes Federrath, Herbert Klimant, Andreas Pfitzmann, Rudi Piotraschke, Andreas Westfeld, Guntram Wicke, and Gritta Wolf. Modeling the security of steganographic systems. In Aucsmith [[Auc98](#)], pages 344–354.

## A Steganographic communication schemes

Here, we give a formal definition of a steganographic communication scheme, formalized as a protocol consisting of  $r$  exchange-rounds in which two parties  $P_0$  and  $P_1$  engage in communication distributed according to a cover distribution  $\mathcal{C}$ , and in which  $P_0$  steganographically transmits a hidden message  $\text{msg}$  to  $P_1$ . The cover distribution  $\mathcal{C}$  is a parameter of the communication scheme.

**Definition A.1** (Steganographic communication scheme). A *steganographic communication scheme* with respect to a cover distribution  $\mathcal{C}$  is a two-party protocol:

$$\Pi^{\mathcal{C}} = (\Pi_{0,1}^{\mathcal{C}}, \Pi_{1,1}^{\mathcal{C}}, \Pi_{0,2}^{\mathcal{C}}, \Pi_{1,2}^{\mathcal{C}}, \dots, \Pi_{0,r}^{\mathcal{C}}, \Pi_{1,r}^{\mathcal{C}}; \Pi_{1,\text{out}}^{\mathcal{C}})$$

where  $r \in \text{poly}$  and where each  $\Pi_{b,i}^{\mathcal{C}}$  is a PPT algorithm with oracle access to a  $\mathcal{C}$ -sampler. Party  $P_0$  is assumed to receive as input a message  $\text{msg}$  that is to be conveyed to  $P_1$  in an undetectable fashion. The protocol must satisfy the following syntax, correctness and security guarantees.

- **Syntax.** For each  $i = 1, \dots, r$ :

1.  $P_0$  draws  $(c_{0,i}, \mathfrak{s}_{0,i}) \leftarrow \Pi_{0,i}^{\mathcal{C}}(\text{msg}, \mathfrak{s}_{0,i-1})$ , locally stores  $\mathfrak{s}_{0,i}$ , and sends  $c_{0,i}$  to  $P_1$ .
2.  $P_1$  draws  $(c_{1,i}, \mathfrak{s}_{1,i}) \leftarrow \Pi_{1,i}^{\mathcal{C}}(\mathfrak{s}_{1,i-1})$ , locally stores  $\mathfrak{s}_{1,i}$ , and sends  $c_{1,i}$  to  $P_0$ .

After the  $r$ th exchange-round,  $P_1$  computes  $\text{msg}' = \Pi_{1,\text{out}}^{\mathcal{C}}(c_{0,1}, c_{1,1}, \dots, c_{0,r}, c_{1,r})$  and halts.

- **Correctness.** For any  $\text{msg} \in \{0, 1\}^{\kappa}$ ,  $\text{msg}' = \text{msg}$  with overwhelming probability. The probability is taken over the randomness of all of the  $\Pi_{b,i}^{\mathcal{C}}$  and their  $\mathcal{C}$ -samples.
- **Steganographic Indistinguishability w.r.t.  $\mathcal{M}$ .** The following distribution

$$\left\{ (c_{0,1}, c_{1,1}, \dots, c_{0,r}, c_{1,r}) : \text{msg} \leftarrow \mathcal{M}, (c_{0,i}, \mathfrak{s}_{0,i}) \leftarrow \Pi_{0,i}^{\mathcal{C}}(\text{msg}, \mathfrak{s}_{0,i-1}), (c_{1,i}, \mathfrak{s}_{1,i}) \leftarrow \Pi_{1,i}^{\mathcal{C}}(\mathfrak{s}_{1,i-1}) \right\}$$

is computationally indistinguishable from the cover distribution of length  $2r - 1$ . (Note that  $\mathcal{C}_{2r-1}$  is not necessarily a product distribution.)

Depending on the specific application at hand, a steganographic communication scheme might require that steganographic indistinguishability hold w.r.t. *all* message distributions  $\mathcal{M}$ , or alternatively, only require that it hold w.r.t. certain specific message distributions (*e.g.*, uniformly random messages).

Definition A.1 concerns the transmission of only a single message. There is a natural generalization of this definition to a more general *multi-message* version, analogous to how Definition 3.3 is a multi-message generalization of Definition 3.1 in the context of subliminal communication schemes. Much as observed in Remark 2 in the context of subliminal communication schemes, the natural multi-message generalization of Definition A.1 is *equivalent* to Definition A.1 in the sense that the existence of any single-message scheme easily implies a multi-message scheme and vice versa; the multi-message version of the definition is mainly useful when considering the asymptotic efficiency of schemes. We present only Definition A.1 here as the simpler single-message definition suffices in the context of an impossibility result.

The existing notions of public-key steganography and steganographic key exchange, as well as the *subliminal communication schemes* introduced in this work, are all instantiations of the more general definition of a (multi-message) steganographic communication scheme.<sup>15</sup>

<sup>15</sup>Public-key steganography protocols can be thought of as being parametrized by the public/private key pairs of

## A.1 Proof of Theorem 4.2

In this subsection we give the proof of impossibility of non-trivial steganographic communication in the presence of an adversarially chosen cover distribution. We recall the statement of Theorem 4.2 from Section 4.2.

**Theorem 4.2.** Let  $\Pi$  be a steganographic communication scheme. Then for any  $k \in \mathbb{N}$ , there exists a cover distribution  $\mathcal{C}$  of conditional min-entropy  $k$  such that the steganographic indistinguishability of  $\Pi$  does not hold for more than one message.

*Proof.* The proof exploits the impossibility of extracting more than one bit from Santha-Vazirani (SV) sources (cf. Proposition 6.6 in [Vad12]). Recall that a  $\delta$ -SV source  $X \in \{0, 1\}^n$  is defined by the following conditions:

$$\delta \leq \Pr[X_i = 1 \mid X_1 = x_1, \dots, X_{i-1} = x_{i-1}] \leq 1 - \delta$$

for all  $i \in [n]$  and all  $(x_1, \dots, x_{i-1}) \in \{0, 1\}^{i-1}$ .

Consider  $\delta > 0$  and  $\Pi$  a steganographic communication scheme with  $r$ -rounds. Let  $n \in \mathbb{N}$  be such that  $\delta$ -SV sources of length  $n$  have conditional min-entropy  $k$  when interpreted as  $2r - 1$  blocks of size  $\frac{n}{2r-1}$ . Assume for contradiction that  $\Pi$  can embed at least two messages  $m_1$  and  $m_2$ . Steganographic indistinguishability applied to the case where  $\text{msg} \leftarrow \{m_1, m_2\}$  implies the existence of a negligible  $\varepsilon$  such that  $\mathbf{1}[\Pi_{2,\text{out}}(\mathcal{C}_{2r-1}) = m_1]$  is a random bit with bias at most  $\varepsilon$ .

But Proposition 6.6 in [Vad12] implies the existence of a  $\delta$ -SV source  $\mathcal{C}_{2r-1}$  such that  $\Pr[\Pi_{2,\text{out}}(\mathcal{C}_{2r-1}) = m_1] \geq 1 - \delta$  or  $\Pr[\Pi_{2,\text{out}}(\mathcal{C}_{2r-1}) = m_1] \leq \delta$ . In other words,  $\mathbf{1}[\Pi_{2,\text{out}}(\mathcal{C}_{2r-1}) = m_1]$  has bias at least  $\frac{1}{2} - \delta$ . For  $\delta < \frac{1}{2} - \varepsilon$ , this is a contradiction.  $\square$

**Remark 7.** At first glance, it may seem that one could bypass this impossibility and communicate more than one bit of information by sequentially repeating the same steganographic scheme. However, one could then apply Theorem 4.2 to the sequential composition and obtain a new cover distribution which would break this scheme.

## B Min-Entropy of Ciphertexts

**Lemma 5.4.** Let PKE be a semantically secure encryption scheme. Then there exists a negligible function  $\varepsilon$  such that for all  $\kappa \in \mathbb{N}$ ,  $m \in \mathcal{M}_\kappa$ , writing  $C_m^{pk} \sim \text{PKE.Enc}(pk, m)$ :

$$\Pr \left[ (pk, sk) \leftarrow \text{PKE.Gen}(1^\kappa) : H_\infty(C_m^{pk}) \geq \log \frac{1}{\varepsilon(\kappa)} \right] \geq 1 - \varepsilon(\kappa).$$

*Proof.* Let us assume by contradiction that there exists  $p \in \text{poly}(\kappa)$ , a countably infinite set  $I$  and a family  $\{m_\kappa\}_{\kappa \in I}$  such that for all  $\kappa \in I$ :

$$\Pr \left[ (pk, sk) \leftarrow \text{PKE.Gen}(1^\kappa) : H_\infty(C_{m_\kappa}^{pk}) \leq \log p(\kappa) \right] \geq \frac{1}{p(\kappa)}. \quad (5)$$

Since PKE is non-trivial, for all  $\kappa \in I$  there exists  $m'_\kappa \neq m_\kappa$ .

Consider the distinguisher  $D$  which on input  $(pk, c)$  runs as follows:

---

the communicating parties, which are initially sampled by the parties according to some key generation algorithm. However, the definition of a steganographic communication protocol as stated is not parametrized in an analogous way. This syntactic discrepancy can be resolved by equivalently thinking of the public-key steganography as a *family* of steganographic communication protocols each of which has the key pairs hardwired, induced by the sampling of key pairs according to the key generation algorithm; and then requiring the steganographic indistinguishability guarantee to hold only with overwhelming probability over key generation.

1. Sample encryption  $e \leftarrow \text{PKE.Enc}(pk, m_\kappa)$ .
2. If  $e = c$  output 1, otherwise output 0.

Writing  $(PK, SK) \sim \text{PKE.Gen}(1^\kappa)$ , we claim that  $D$  distinguishes  $(PK, \text{Enc}(PK, m_\kappa))$  from  $(PK, \text{Enc}(PK, m'_\kappa))$  with non-negligible advantage. First note that by (5), with probability at least  $\frac{1}{p(\kappa)}$  over the draw of  $pk$ ,  $\text{Enc}(pk, m_\kappa)$  has collision probability at least  $\frac{1}{p(\kappa)^2}$ . By definition,  $D$  outputs 1 on input  $(pk, \text{Enc}(pk, m_\kappa))$  if and only if a collision occurs in step 1. Hence:

$$\Pr[D(PK, \text{Enc}(PK, m_\kappa)) = 1] \geq \frac{1}{p(\kappa)^3}.$$

Second, since PKE is correct, there exists a negligible function  $\varepsilon'$  such that  $\text{Enc}(pk, m_\kappa) = \text{Enc}(pk, m'_\kappa)$  with probability at most  $\varepsilon'$ . Hence:

$$\Pr[D(PK, \text{Enc}(PK, m'_\kappa)) = 1] \leq \varepsilon'(\kappa).$$

The previous two inequalities together imply that for large enough  $\kappa \in I$ :

$$\left| \Pr[D(PK, \text{Enc}(PK, m_\kappa)) = 1] - \Pr[D(PK, \text{Enc}(PK, m'_\kappa)) = 1] \right| \geq \frac{1}{2p(\kappa)^3},$$

which contradicts the semantic security of PKE.  $\square$

## C Greater-Than is a Same-Source Extractor

**Lemma 5.7.** For any  $k \leq n$ , GT is a  $(k, \frac{1}{2^k})$  same-source extractor.

*Proof.* Let  $X$  and  $Y$  be two i.i.d. variables with  $H_\infty(X) = H_\infty(Y) \geq k$ . Then

$$2 \Pr[X \geq Y] = \Pr[X \geq Y] + \Pr[X \leq Y] = 1 + \text{CP}(X) \leq 1 + \frac{1}{2^k},$$

where the last inequality uses the standard upper bound on the collision probability  $\text{CP}(X)$  with respect to the min-entropy of  $X$ .  $\square$