

# Message Transmission with Reverse Firewalls— Secure Communication on Corrupted Machines

Yevgeniy Dodis<sup>1\*</sup>, Ilya Mironov<sup>2</sup>, and Noah Stephens-Davidowitz<sup>1\*\*</sup>

<sup>1</sup> Dept. of Computer Science, New York University.

<sup>2</sup> Google.

**Abstract.** A secure reverse firewall, as recently defined by Mironov and Stephens-Davidowitz, is a third party that “sits between a user and the outside world” and modifies the user’s sent and received messages so that *even if the user’s machine has been corrupted*, her security is still guaranteed. In other words, reverse firewalls allow us to provide meaningful (and, indeed, very strong) security guarantees against powerful adversaries that may have tampered with the user’s hardware or software (or adversaries that are aware of bugs in the user’s protocol implementation). A long list of recent events and disclosures shows that such threats are extremely common in practice, and they present a serious, arguably existential, threat to cryptography. Importantly, reverse firewalls defend against such threats without sharing any secrets with the user, and in general we expect the user to place no more trust in the firewall than she places in her communication channel.

While Mironov and Stephens-Davidowitz demonstrated that reverse firewalls can be constructed for very strong cryptographic primitives (which are of mostly theoretical interest), we study reverse firewalls for perhaps the most natural cryptographic task: secure message transmission. We find a rich structure of solutions that vary in efficiency, security, and setup assumptions, in close analogy with message transmission in the classical setting. Our strongest and most important result shows a protocol that achieves interactive, concurrent CCA-secure message transmission with a reverse firewall—i.e., CCA-secure message transmission on a possibly compromised machine! Surprisingly, this protocol is quite efficient and simple, requiring only a small constant number of public-key operations. It could easily be used in practice. Behind this result is a technical composition theorem that shows how key agreement with a sufficiently secure reverse firewall can be used to construct a message-transmission protocol with its own secure reverse firewall.

## 1 Introduction

In the past few years, the cryptographic community has learned of a disturbingly wide array of new security vulnerabilities. The revelations of Edward Snowden show that the United States National Security Agency successfully gained access to secret information by extraordinary means, including subverting cryptographic standards [PLS13,BBG13] and intercepting and tampering with hardware on its way to users [Gre14]. Meanwhile, many (apparently accidental) security flaws have been found in widely deployed pieces of cryptographic software, leaving users completely exposed [LHA<sup>+</sup>12, CVE14b, CVE14a, CVE14c]. Due to the complexity of modern cryptographic software, such vulnerabilities are extremely hard to detect in practice, and, ironically, cryptographic modules are often the easiest to attack, as attackers can often use cryptographic mechanisms to mask their activities or opportunistically hide their communications within encrypted traffic (as in the case of the Heartbleed vulnerability). This leads to cryptography’s current existential crisis, summarized by the following (terrifying) question: “How can we provide any meaningful security guarantees when the adversary may have arbitrarily tampered with a victim’s computer?”

---

\* Partially supported by gifts from VMware Labs and Google, and NSF grants 1319051, 1314568, 1065288, 1017471.

\*\* Partially supported by National Science Foundation under Grant No. CCF-1320188. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

Motivated by these concerns, Mironov and Stephens-Davidowitz recently proposed the novel concept of (*cryptographic*) *reverse firewalls*, designed to protect machines from *insider* attacks [MS15]. Reverse firewalls are autonomous intermediaries that intercept and “sanitize” transit traffic to back-stop security of network clients. A cryptographic protocol equipped with a reverse firewall has the desirable guarantee that its security is preserved even if it is run on a compromised machine. In other words, a correctly implemented reverse firewall ensures the security of a cryptographic protocol *even if the user’s own implementation of the protocol is buggy, or even maliciously compromised*. Such a tool is a natural and powerful solution to the problems discussed above. We stress that the firewall is untrusted in the sense that it shares no secrets with the user. (It is best to think of the reverse firewall as “trusted no more than the communication channel.”)

More concretely, Mironov and Stephens-Davidowitz start by considering an arbitrary cryptographic protocol that satisfies some notions of functionality (i.e., correctness) and security.<sup>1</sup> For example, perhaps the simplest non-trivial case is semantically secure message transmission from Alice to Bob, which has the functionality requirement that Bob should receive the correct plaintext message from Alice and the security requirement that a computationally bounded adversary “should not learn anything about Alice’s plaintext message” from the transcript of a run of the protocol. Formally, we can model this functionality by providing Alice with an input plaintext and requiring Bob’s output to match this on all inputs, and we can model semantic security by a standard indistinguishability security game.

A reverse firewall for Alice in such a protocol is a party that “sits between” Alice and the outside world and modifies the messages that Alice sends and receives. Such a firewall *maintains functionality* if the resulting protocol achieves the same functionality as the original. E.g., in the case of message transmission, Bob’s output should still match Alice’s input plaintext. More interestingly, the firewall *preserves security* if the protocol remains secure when we replace Alice’s computer with some arbitrarily corrupted party. For example, a reverse firewall for Alice preserves semantic security of message transmission if a computationally bounded adversary “learns nothing about Alice’s plaintext message” from the transcript of messages sent between the firewall and Bob, *regardless of how Alice behaves*. (E.g., the firewall may rerandomize the messages that Alice sends in a way that makes them indistinguishable from random from the adversary’s perspective, regardless of Alice’s original message. We analyze such protocols in a stronger setting in Section 3.1.) Finally, we say that the reverse firewall *resists exfiltration* if a compromised implementation “cannot leak any information to the outside world” beyond the functionality of the protocol.<sup>2</sup> Naturally, we consider reverse firewalls for Bob as well. (See Section 2.1 for the formal definitions and [MS14] for a thorough discussion.)

It is convenient to identify a class of *functionality-maintaining corruptions*: compromised implementations that are “technically legal” in the sense that they may deviate arbitrarily from the protocol, as long as they do not break its functionality. Some of our reverse firewalls are only secure against this type of corruption. (This model is introduced by [MS15], and the authors call security against unrestricted compromise *strong* security.) We emphasize that, while this restricted class of compromised implementations is not ideal, it is still quite large. In particular, all of the real-world com-

<sup>1</sup> The notion of functionality in [MS15] is quite simple, and it should not be confused with the much more complicated concept of functionality used in the universal composability framework. Formally, Mironov and Stephens-Davidowitz define a functionality requirement as any condition on the output of the parties that may depend on the input, and in practice, these requirements are straightforward.

<sup>2</sup> Often, exfiltration resistance is implied by security preservation, but not always. For example, a reverse firewall for Alice for a protocol in which Alice sends a message to Bob preserves semantic security if and only if it is exfiltration resistant. But, semantic security does not guarantee that *Bob*, the receiving party, will not leak any information. It obviously can be desirable for Bob to have a reverse firewall in such a protocol that resists exfiltration, even though the security of the underlying protocol does not require this.

promises mentioned above fall into this category [PLS13, BBG13, Gre14, Sup15, LHA<sup>+</sup>12, CVE14b, CVE14a, CVE14c], as do essentially all other forms of compromise considered in prior work, such as faulty or backdoored PRNGs [DGG<sup>+</sup>15], Algorithm Substitution Attacks [BPR14a], subliminal channels [Sim84], etc. (We discuss functionality-maintaining corruption in our setting in Section 2.4. See [MS14] for a detailed comparison of the general reverse firewall framework with prior work.)

Mironov and Stephens-Davidowitz demonstrate feasibility of this framework by constructing reverse firewalls for parties participating in Oblivious Transfer and Secure Function Evaluation protocols—very strong cryptographic primitives. The fact that such strong primitives can be made secure in this model is quite surprising and bodes well for the reverse firewalls framework. However, these protocols are rather inefficient and therefore mostly of theoretical interest. To fulfill the promise of reverse firewalls, we need to consider protocols of more practical importance, which is the focus of the present work. We demonstrate efficient protocols with reverse firewalls for perhaps the most natural cryptographic task: secure message transmission. Our strongest protocols provide security guarantees in the reverse firewall setting comparable to the strongest guarantees obtained in the classical setting (e.g., security against adaptive chosen ciphertext attacks with forward secrecy). Surprisingly, our protocols are also quite efficient, requiring only a few rounds and a small constant number of public-key operations on bit strings whose length does not depend on the length of the plaintext. Indeed, our protocols are comparably efficient to those used in practice, though they satisfy a much stronger notion of security. To achieve this, we also develop efficient key-agreement protocols with efficient, secure reverse firewalls, which might be of independent interest. We group our solutions according to a variety of setup assumptions, in which (1) the parties share a secret key; (2) there is no public-key infrastructure; (3) both parties have a public-key/private-key pair; or (4) only the recipient has a public-key/private-key pair. Non-trivial levels of security are achievable in almost all cases, and the various schemes mirror closely those that appear in the classical setting.

## 1.1 Related work

Message transmission in the classical setting (i.e., without reverse firewalls) is of course extremely well-studied, and a summary of such work is beyond the scope of this paper. (We note, however, that our notions of security of interactive message transmission protocols follow closely Dodis and Fiore [DF14].) Similarly, there have been many different approaches to cryptography in the presence of compromise. We therefore provide only a limited discussion of related work. (Both [BPR14a] and [MS15] contain a more thorough discussion of alternative approaches to security in the presence of a compromise, though they naturally do not mention the many relevant papers that appeared simultaneously with or after their publication, such as [AMV15, BJK15, RTYZ15, DFP15, BH15, DGG<sup>+</sup>15, SFKR15]. In particular, [MS15] contains a detailed comparison of the reverse firewalls framework with many prior models, showing that this framework generalizes much of the prior work on insider attacks and various related notions.)

Bellare, Paterson, and Rogaway consider the problem of securely encrypting a message when the encrypting party might be compromised [BPR14a]. However, they do not allow for reverse firewalls (as their work predates the RF framework). As a result, they necessarily restrict themselves to the case in which the parties share a secret key and the compromised party maintains functionality. (I.e., a tampered encryption algorithm must produce ciphertexts that decrypt to the correct plaintext.) Their main result is a lower bound, showing roughly that even a relatively weak adversary can break any scheme that “non-trivially uses randomness” in this setting. (Of course, this lower bound no longer applies in the reverse firewalls framework, as our results demonstrate.) They also provide a positive result that is more relevant to our work, in which they show how to build a CCA-secure deterministic stateful scheme. The main idea behind their positive result is the observation that, if

an encryption scheme is deterministic and has a unique ciphertext that decrypts to each plaintext, then a tampered implementation of the encryption algorithm can only maintain functionality if its behavior is identical to the honest encryption algorithm. In short, they show that “randomness is dangerous” in this setting.

Our work can be viewed as a generalization of [BPR14a] in a number of directions. We consider multi-round protocols in which the parties might not share secret keys, and we consider arbitrarily compromised adversaries. In order to do this, we use the reverse firewalls framework of [MS15]. While Bellare et al.’s work stresses the danger of randomness in secure message transmission, our work highlights the benefits of randomness. In particular, our schemes rely heavily on “rerandomization” by the reverse firewall. However, we do use the deterministic encryption scheme of Bellare et al. as part of two of our protocols.

Recently, Ateniese, Magri, and Venturi studied reverse firewalls for signature schemes [AMV15]. Their work can be considered as complementary to ours, as we are concerned with privacy, while they consider authentication. We also note that our more advanced key-agreement scheme uses unique signatures, and we implicitly rely on the fact that unique signatures have a reverse firewall, as [AMV15] prove. Indeed, the more general primitive of rerandomizable signatures that they consider would also suffice for our purposes and might be more efficient in practice.

Finally, we note that some of our study of key agreement is similar to work on key-agreement protocols secure against active insiders, and the study of key control (e.g., [PW03,KS05,DPSW06]). These works consider key-agreement protocols involving at least three parties, in which one or more of the participants wishes to maliciously fix the key or otherwise subvert the security of the protocol. Some of the technical challenges that we encounter are similar to those encountered in the key control literature, and indeed, the simple toy protocol that we present in Section 3.4 can be viewed as a simple instantiation of some of the known solutions to the key-control problem (e.g., [DPSW06]). However, since prior work approached this problem from a different perspective—with three or more parties and without reverse firewalls—our more technical solutions presented in Section 4.2 are quite different. In particular, the key-control literature (and most literature on key agreement in general) typically focuses on creating protocols that produce “non-malleable” keys, whereas our protocols need some type of malleability specifically to allow the firewall to maul the keys! (We accomplish this without sacrificing security.)

## 1.2 Our results

Secure transmission of a message between two parties is arguably the most basic and essential cryptographic task. In this section, we walk through several common setup assumptions for this problem and outline our solutions for secure message-transmission protocols with reverse firewalls against end-point compromises. In what follows, Alice is always the sender and Bob is always the receiver of the message. All of our security notions apply to the concurrent setting, in which the adversary may instantiate many runs of the protocol simultaneously.

**Symmetric-key setting.** In the first and simplest scenario, Alice and Bob have a shared secret key. (See Appendix A.) Quite naturally, Alice may want to use a symmetric-key encryption scheme to communicate with Bob. Using a standard scheme (e.g., AES-CBC) would, however, expose her to a number of “algorithm-substitution attacks” (ASA) described by Bellare, Paterson, and Rogaway [BPR14a], such as IV-replacement or a biased-ciphertext attacks. To defend against ASA Bellare et al. propose using a *stateful, deterministic* encryption scheme based on either a counter or a nonce, as we described above. We briefly consider this case, observing that their solution corresponds to a one-round protocol in our model (in which the firewall simply lets messages

pass unaltered). We also show that strong security (when Alice’s implementation is not necessarily functionality-maintaining) is not achievable within this framework without using (less efficient) public-key primitives.

**Unkeyed setting.** In the next scenario, Alice and Bob have no shared secret keys and no pre-existing public-key infrastructure. (See Section 3.) This is the first setting in which reverse firewalls meaningfully improve security. Since neither the sender nor the receiver can be authenticated without public keys, the strongest guarantee achievable in this scenario is security against passive adversaries, i.e. security against chosen-plaintext attacks (CPA). While this level of security is typically considered to be insufficient in practice, the ideas that we develop here will be useful for solving the harder problems that we discuss below. As such, we present two solutions in this setting.

First, we consider a simple two-round protocol in which Bob sends a (freshly generated) public key and Alice responds with an encryption of her plaintext under Bob’s key. Of course, if Alice’s computer is compromised, then the ciphertext that she sends to Bob can be a vulnerability. It may be used as a channel to leak a secret, or it may be encrypted using “bad randomness,” under a completely different (adversarial) key, or otherwise improperly. There is of course a simple fix for this: Alice and Bob can use a rerandomizable encryption scheme, and Alice’s reverse firewall can rerandomize her ciphertext. If, on the other hand, Bob’s computer is compromised, then the key that he sends can also be used as a channel to leak his secrets—either to Alice or an eavesdropper—or he can simply send Alice an insecure key, so that an eavesdropper will be able to read Alice’s plaintext. If Alice and Bob use an encryption scheme that is “key malleable,” then Bob’s reverse firewall can fix this by rerandomizing the key that Bob sends to Alice and then converting Alice’s response from a ciphertext under the rerandomized key into a ciphertext under Bob’s original key. (We note that both the keys and ciphertexts used in ElGamal encryption can easily be rerandomized in the way that we require here.) Indeed, both reverse firewalls guarantee *strong* security for Alice and Bob—even if the players’ implementations produce malformed messages or just refuse to cooperate, their respective firewalls can still guarantee security!

The above solution is simple and elegant, but we suspect that such protocols will be of little practical value because they require the computation of public-key operations on the entire plaintext. Since plaintexts are often quite long and public-key operations tend to be much slower than symmetric-key operations, this can be quite inefficient. In practice, it is much faster to use public-key operations to transmit a (relatively short) key for a symmetric-key encryption scheme and then to send the plaintext encrypted under this symmetric key. There are two general methods for transmitting this key in the classical setting: hybrid encryption and key agreement.

Before we describe our solution based on key agreement, it will be instructive to consider why hybrid encryption seems to fail in this setting. Recall that in a hybrid encryption scheme, Bob sends a public key to Alice (as above), Alice selects a uniformly random key  $rk$  for a symmetric-key scheme and sends  $rk$  encrypted under Bob’s public key together with the encryption of her message under  $rk$ . Such a scheme fails in our setting because a corrupted implementation of Alice can choose a “bad key”  $rk^*$  with which to encrypt the message. The “bad key”  $rk^*$  might be known to an adversary; might be chosen so that the ciphertext takes a specific form that leaks some information; or might otherwise compromise Alice’s security. It seems that a firewall cannot hope to fix this problem without resorting to slow public-key primitives.

Luckily, key agreement can be combined with symmetric-key encryption to securely and efficiently transmit a message in the setting of reverse firewalls. A key-agreement protocol allows Alice and Bob to jointly select a secret key over an insecure channel. Such a protocol is often used in conjunction with symmetric-key encryption in the classical setting, where it is justified by

composition theorems relating the security of the message-transmission protocol to the underlying key-agreement protocol. Indeed, we give an analogous result (Theorem 2) that works in our setting, showing that a key-agreement protocol with sufficiently secure reverse firewalls can be combined with symmetric-key encryption to produce an efficient CPA-secure message-transmission protocol with secure reverse firewalls. This motivates the study of key agreement in the setting of reverse firewalls, and in Section 3.4, we show how to construct a relatively simple key-agreement protocol with reverse firewalls that suffices.

Intuitively, the reason that this solution does not suffer the same fate as hybrid encryption in this setting is because the key is chosen interactively, and the firewall therefore has the opportunity to rerandomize the key. To truly prevent any party from having any control over the final key, we use a three-round protocol in which Bob commits to his share of the randomness used to generate the key first. Alice then sends her share of the randomness, and finally Bob opens his commitment. We show that, when the commitment scheme and underlying key-agreement protocol are suitably rerandomizable, this idea yields a key-agreement protocol with a secure reverse firewall for both parties.

**Publicly keyed setting.** In the stronger setting in which we have a public-key infrastructure, with a public-key/private-key pair for each party (considered in Section 4), we can hope to achieve interactive CCA security (i.e., the adversary may interact with the parties arbitrarily and may view Bob’s output of decrypted messages). Perhaps surprisingly, we prove a generic composition theorem in this setting, which shows that it suffices to construct a key-agreement protocol with reverse firewalls that satisfy certain security properties. In analogy with the unkeyed setting, after agreeing to a key with Bob, Alice can use symmetric-key encryption to send the actual plaintext message. The resulting protocol is CCA-secure, and Alice’s reverse firewall preserves this security. (See Theorem 4.)

To instantiate this scheme, we then construct (in Section 4.2) a key-agreement protocol that satisfies the necessary security properties. In particular, the scheme is secure against active adversaries, and each party has a reverse firewall that preserves this security. This key-agreement protocol is significantly more technical than our scheme in the unkeyed setting, and common techniques from classical key agreement do not work. In particular, many key-agreement protocols achieve security against an active adversary by having both parties sign the transcript at the end of the protocol; intuitively, this allows the parties to know if the adversary has tampered with any messages, so that they will never agree to different keys or to a key chosen dishonestly. But in our setting, we actually *want* the firewall to be able to modify the parties’ messages! We therefore need to find some unique information that the parties can use to confirm that they have agreed to the same key without preventing the firewall from modifying the key. Furthermore, we need the firewall to be able to check signatures itself, so that it can block invalid messages. Therefore, our primary technical challenge in this context is to find some piece of information that (1) uniquely identifies the key; (2) respects the firewall’s changes to the parties’ messages; and (3) is efficiently computable from the transcript. And, of course, the key itself must still be indistinguishable from random in the presence of this information. Our solution uses hashed Diffie-Hellman (similar in spirit to [Kil07]) and bilinear maps. (We also use unique signatures to prevent the signatures from becoming a channel themselves.)

Surprisingly, in spite of the challenges, our protocol only requires four rounds with relatively short messages, and the parties themselves (including the firewall) only need to perform a small constant number of operations.<sup>3</sup> This compares quite favorably with protocols that are currently

---

<sup>3</sup> For example, Bob only needs to perform one group exponentiation, compute a bilinear map once, check the opening of a commitment, verify one (unique) signature, (uniquely) sign one string, and apply a hash function once.

implemented in practice (which of course are not secure in our setting), and we therefore believe that this protocol can and should be implemented and used in the real world.

**Singly keyed setting.** In the “intermediate” setting in which Bob has a public-key/private-key pair but Alice does not, we show how to achieve CCA security in a single round using rerandomizable RCCA-secure encryption [Gro04, PR07]. (See Section 5.) However, these schemes do not achieve forward secrecy, and they require relatively inefficient public-key primitives that must be applied to full plaintexts. We leave to future work the question of whether CCA-secure and forward-secret message transmission is achievable in this regime. We also do not know whether a composition theorem similar to Theorem 4 exists in this regime. (We discuss some additional open questions in Section 6.)

## 2 Definitions

### 2.1 Reverse firewalls

We use the definition of reverse firewalls from [MS14] (and we refer the reader to [MS14] for a longer discussion of the reverse firewall framework). In particular, a reverse firewall  $\mathcal{W}$  is a stateful algorithm that maps messages to messages. For a party  $A$  and reverse firewall  $\mathcal{W}$ , we define  $\mathcal{W} \circ A$  as the “composed” party in which  $\mathcal{W}$  is applied to the messages that  $A$  receives before  $A$  “sees them” and the messages that  $A$  sends before they “leave the local network of  $A$ .”  $\mathcal{W}$  has access to all public parameters, but not to the private input of  $A$  or the output of  $A$ . (We can think of  $\mathcal{W}$  as an “active router” that sits at the boundary between Alice’s private network and the outside world and modifies Alice’s incoming and outgoing messages.) We repeat all relevant definitions from [MS14] below, and we add two new ones.

As in [MS14], we assume that a cryptographic protocol comes with some functionality or correctness requirements  $\mathcal{F}$  and security requirements  $\mathcal{S}$ . (For example, a functionality requirement  $\mathcal{F}$  might require that Alice and Bob output the same thing at the end of the protocol. A security requirement might ask that no efficient adversary can distinguish between the transcript of the protocol and uniformly random strings.) Throughout, we use  $\tilde{A}$  to represent arbitrary adversarial implementations of party  $A$  and  $\bar{A}$  to represent functionality-maintaining adversarial implementations of  $A$  (i.e., implementations of  $A$  that still satisfy the functionality requirements of the protocol.) For a protocol  $\mathcal{P}$  with party  $A$ , we write  $\mathcal{P}_{A \rightarrow \tilde{A}}$  to represent the protocol in which the role of party  $A$  is replaced by party  $\tilde{A}$ .

We are only interested in firewalls that themselves maintain functionality. In other words, the *composed party*  $\mathcal{W} \circ A$  should not break the correctness of the protocol. (Equivalently,  $\mathcal{P}_{A \rightarrow \mathcal{W} \circ A}$  should satisfy the same functionality requirements as the underlying protocol  $\mathcal{P}$ .) We follow [MS15] in requiring something slightly stronger—reverse firewalls should be “stackable”, so that many reverse firewalls composed in series  $\mathcal{W} \circ \dots \circ \mathcal{W} \circ A$  still do not break correctness. All of our firewalls will trivially satisfy this notion. Note as well that we are not interested in protocols whose functionality “depends on the presence of the reverse firewall,” so we require that the protocol without the reverse firewall is also functional.

**Definition 1 (Reverse firewall).** *A reverse firewall  $\mathcal{W}$  maintains functionality  $\mathcal{F}$  for party  $A$  in protocol  $\mathcal{P}$  if protocol  $\mathcal{P}$  satisfies  $\mathcal{F}$ , the protocol  $\mathcal{P}_{A \rightarrow \mathcal{W} \circ A}$  satisfies  $\mathcal{F}$ , and the protocol  $\mathcal{P}_{A \rightarrow \mathcal{W} \circ \dots \circ \mathcal{W} \circ A}$  also satisfies  $\mathcal{F}$ . (I.e., we can compose arbitrarily many reverse firewalls without breaking functionality.)*

Of course, a firewall is not interesting unless it provides some benefit. The most natural reason to deploy a reverse firewall is to *preserve* the security of a protocol, even in the presence of compromise.

The below definition (which again follows [MS14]) captures this notion by asking that the protocol obtained by replacing party  $A$  with  $\mathcal{W} \circ \tilde{A}$  for an arbitrary corrupted party  $\tilde{A}$  still achieves some notion of security. For example, when we consider message transmission, we will want the firewall to guarantee Alice’s privacy against some adversary, even when Alice’s own computer has been corrupted. (As before, we are only interested in protocols that are secure without the reverse firewall as well. I.e., the reverse firewall should *preserve* security of the protocol when a party has been compromised, but it should not be necessary when the party is behaving correctly.)

**Definition 2 (Security preservation).** *A reverse firewall strongly preserves security  $\mathcal{S}$  for party  $A$  in protocol  $\mathcal{P}$  if protocol  $\mathcal{P}$  satisfies  $\mathcal{S}$ , and for any polynomial-time algorithm  $\tilde{A}$ , the protocol  $\mathcal{P}_{A \rightarrow \mathcal{W} \circ \tilde{A}}$  satisfies  $\mathcal{S}$ . (I.e., the firewall can guarantee security even when an adversary has tampered with  $A$ .)*

*A reverse firewall preserves security  $\mathcal{S}$  for party  $A$  in protocol  $\mathcal{P}$  satisfying functionality requirements  $\mathcal{F}$  if protocol  $\mathcal{P}$  satisfies  $\mathcal{S}$ , and for any polynomial-time algorithm  $\tilde{A}$  such that  $\mathcal{P}_{A \rightarrow \tilde{A}}$  satisfies  $\mathcal{F}$ , the protocol  $\mathcal{P}_{A \rightarrow \mathcal{W} \circ \tilde{A}}$  satisfies  $\mathcal{S}$ . (I.e., the firewall can guarantee security even when an adversary has tampered with  $A$ , provided that the tampered implementation does not break the functionality of the protocol.)*

```

proc. LEAK( $\mathcal{P}, A, B, \mathcal{W}, \lambda$ )
 $(\tilde{A}, \tilde{B}, I) \leftarrow \mathcal{E}(1^\lambda)$ 
 $b \xleftarrow{\$} \{0, 1\}$ 
IF  $b = 1$ ,  $A^* \leftarrow \mathcal{W} \circ \tilde{A}$ 
ELSE,  $A^* \leftarrow \mathcal{W} \circ A$ 
 $\mathcal{T}^* \leftarrow \mathcal{P}_{P_j \rightarrow A^*, J \rightarrow \tilde{B}}(I)$ 
 $b^* \leftarrow \mathcal{E}(\mathcal{T}^*, S_{\tilde{B}})$ 
OUTPUT ( $b = b^*$ )

```

**Fig. 1:** LEAK( $\mathcal{P}, A, B, \mathcal{W}, \lambda$ ), the exfiltration resistance security game for a reverse firewall  $\mathcal{W}$  for party  $A$  in protocol  $\mathcal{P}$  against party  $B$  with input  $I$ .  $S_{\tilde{B}}$  is the state of  $\tilde{B}$  after the run of the protocol, and  $I$  may be any valid input for  $\mathcal{P}$ .

For technical reasons, we will also need a new definition not present in [MS14]. We wish to show composition theorems, allowing us to construct message-transmission protocols with secure reverse firewalls from key-agreement protocols with their own firewalls. In order to accomplish this, we will need the notion of *detectable failure*. Essentially, a protocol fails detectably if we can distinguish between transcripts of valid runs of the protocol and invalid transcripts. We will use this to make sure that a firewall of some larger message-transmission protocol can always detect if the key-agreement sub-protocol has failed in some way. We make this precise below. In order to do so, we need to carefully consider what it means for a transcript to be valid. (The reader may wish to skip these definitions on a first read.) For simplicity, we assume that honest parties always output  $\perp$  when they receive a malformed message (e.g., when a message that should be a pair of group elements is not a pair of group elements).

**Definition 3 (Valid transcripts).** *A sequence of bits  $r$  and private input  $I$  generate transcript  $\mathcal{T}$  in protocol  $\mathcal{P}$  if a run of the protocol  $\mathcal{P}$  with input  $I$  in which the parties’ coin flips are taken from  $r$  results in the transcript  $\mathcal{T}$ . A transcript  $\mathcal{T}$  is a valid transcript for protocol  $\mathcal{P}$  if there is a sequence  $r$  and private input  $I$  generating  $\mathcal{T}$  such that no party outputs  $\perp$  at the end of the run. (Here, we*



assume that the public input is part of the transcript.) A protocol has unambiguous transcripts if for any  $(r, I)$  generating  $\mathcal{T}$ , no party outputs  $\perp$  at the end of a protocol run with  $(r, I)$ . (In other words, a valid transcript never results in a failed run of the protocol.)

**Definition 4 (Detectable failure).** A reverse firewall  $\mathcal{W}$  detects failure for party  $A$  in protocol  $\mathcal{P}$  if

- if  $\mathcal{P}_{A \rightarrow \mathcal{W} \circ A}$  has unambiguous transcripts;
- the firewall  $\mathcal{W}$  outputs the special symbol  $\perp$  when run on any transcript that is not valid for  $\mathcal{P}_{A \rightarrow \mathcal{W} \circ A}$ ; and
- there is a polynomial-time deterministic algorithm that decides whether a transcript  $\mathcal{T}$  is valid for  $\mathcal{P}_{A \rightarrow \mathcal{W} \circ A}$ .

We will also need the notion of *exfiltration resistance*, introduced in [MS14]. We refer the reader to [MS14] for a much more thorough description of exfiltration resistance, but we provide a brief discussion below. Intuitively, a reverse firewall is exfiltration resistant if “no corrupt implementation of Alice can leak information through the firewall.” We say that it is exfiltration resistant for Alice against Bob if Alice cannot leak information to Bob through the firewall, and we say that it is exfiltration resistant against eavesdroppers (or just exfiltration resistant) if Alice cannot leak information through the firewall to an adversary that is only given access to the protocol transcript.

The relationship between security preservation and exfiltration resistance depends on the security notion of the cryptographic primitive in question. E.g., in a message-transmission protocol, a reverse firewall for Alice resists exfiltration if and only if it preserves semantic security. However, it is possible to construct a reverse firewall for a key-agreement protocol that preserves security but does not resist exfiltration (for example, we can add an arbitrary message to the beginning of any key-agreement protocol without changing its security properties, but clearly such a message can be used to leak information). It is also possible to construct a firewall that resists exfiltration but does not preserve key-agreement security. The second definition below (which uses the notion of valid transcripts) is new to this paper and is necessary for our composition theorems.

**Definition 5 (Exfiltration resistance).** A reverse firewall is exfiltration resistant for party  $A$  against party  $B$  in protocol  $\mathcal{P}$  satisfying functionality  $\mathcal{F}$  if no PPT algorithm  $\mathcal{E}$  with output circuits  $\tilde{A}$  and  $\tilde{B}$  such that  $\mathcal{P}_{A \rightarrow \tilde{A}}$  and  $\mathcal{P}_{B \rightarrow \tilde{B}}$  satisfy  $\mathcal{F}$  has non-negligible advantage in  $\text{LEAK}(\mathcal{P}, A, B, \mathcal{W}, \lambda)$ . If  $B$  is empty, then we simply say that the firewall is exfiltration resistant.

A reverse firewall is exfiltration resistant for party  $A$  against party  $B$  in protocol  $\mathcal{P}$  with valid transcripts if no PPT algorithm  $\mathcal{E}$  with output circuits  $\tilde{A}$  and  $\tilde{B}$  such that  $\mathcal{P}_{A \rightarrow \tilde{A}}$  and  $\mathcal{P}_{B \rightarrow \tilde{B}}$  produce valid transcripts for  $\mathcal{P}$  has non-negligible advantage in  $\text{LEAK}(\mathcal{P}, A, B, \mathcal{W}, \lambda)$ . If  $B$  is empty, then we simply say that the firewall is exfiltration resistant with valid transcripts.

A reverse firewall is strongly exfiltration resistant for party  $A$  against party  $B$  in protocol  $\mathcal{P}$  if no PPT adversary  $\mathcal{E}$  has non-negligible advantage in  $\text{LEAK}(\mathcal{P}, A, B, \mathcal{W}, \lambda)$ . If  $B$  is empty, then we simply say that the firewall is strongly exfiltration resistant.

## 2.2 Message-transmission protocols

A *message-transmission protocol* is a two-party protocol in which one party, Alice, is able to communicate a plaintext message to the other party, Bob. (For simplicity, we only formally model the case in which Alice wishes to send a single plaintext to Bob per run of the protocol, but this of course naturally extends to a more general case in which Alice and Bob wish to exchange many plaintext messages.) We consider two notions of security for such messages. First, we consider *CPA security*,

in which the adversary must distinguish between the transcript of a run of the protocol in which Alice communicates the plaintext  $m_0$  to Bob and the transcript with which Alice communicates  $m_1$  to Bob, where  $m_0$  and  $m_1$  are adversarially chosen plaintexts. (Even in this setting, we allow the adversary to start many concurrent runs of the protocol with adaptively chosen plaintexts.) Our strongest notion of security is *CCA security* in which the adversary may “feed” the parties any messages and has access to a decryption oracle. Our security definitions are similar in spirit to [DF14], but adapted for our setting.

**Session ids.** Throughout this paper, we consider protocols that may be run concurrently many times between the same two parties. In order to distinguish one run of a protocol from another, we therefore “label” each run with a unique session id, denoted  $\text{sid}$ . We view  $\text{sid}$  as an implicit part of every message, and we often ignore  $\text{sid}$  when it is not important. Our parties and firewalls are stateful, and we assume that the parties and the firewall maintain a list of the relevant session ids, together with any information that is relevant to continue the run of the protocol corresponding to  $\text{sid}$  (such as the number of messages sent so far, any values that need to be used later in the protocol, etc.). We typically suppress explicit reference to these states. In our security games, the adversary may choose the value  $\text{sid}$  for each run of the protocol, provided that each party has a unique run for each session  $\text{sid}$ . (In fact, it does not even make sense for the adversary to use the same  $\text{sid}$  for two different runs of the protocol with the same party, as this party will necessarily view any calls with the same  $\text{sid}$  as corresponding to a single run of the protocol. However, as is clear from our security games, an active adversary may maintain two separate runs of a protocol with two different parties but the same  $\text{sid}$ .) In practice,  $\text{sid}$  can be a simple counter or any other nonce (perhaps together with any practical information necessary for communication, such as IP addresses). We note in passing that, in the setting of reverse firewalls, a counter is preferable to, e.g., a random nonce to avoid providing a channel through  $\text{sid}$ , but such concerns are outside of our model and the scope of this paper.

The definition below makes the above formal and provides us with some useful terminology.

**Definition 6 (Message-transmission protocol).** *A message-transmission protocol is a two-party protocol in which one party, Alice, receives as input a plaintext  $m$  from some plaintext space  $\mathcal{M}$ . The protocol is correct if for any input  $m \in \mathcal{M}$ , Bob’s output is always  $m$ .*

*We represent the protocol by four algorithms  $\mathcal{P} = (\text{setup}, \text{next}_A, \text{next}_B, \text{return}_B)$ .  $\text{setup}$  takes as input  $1^\lambda$ , where  $\lambda$  is the security parameter, and returns the starting states for each party,  $S_A, S_B$ , which consist of both private input,  $\sigma_A$  and  $\sigma_B$  respectively, and public input  $\pi$ . Each party’s  $\text{next}$  procedure is a stateful algorithm that takes as input  $\text{sid}$  and an incoming message, updates the party’s state, and returns an outgoing message. The  $\text{return}_B$  procedure takes as input Bob’s state  $S_B$  and  $\text{sid}$  and returns Bob’s final output.*

*We say that a message-transmission protocol is*

- unkeyed if  $\text{setup}$  does not return any private input  $\sigma_A$  or  $\sigma_B$ ;
- singly keyed if  $\text{setup}$  returns private input  $\sigma_B$  for Bob but none for Alice;
- publicly keyed if  $\text{setup}$  returns private input for both parties  $\sigma_A$  and  $\sigma_B$ , but these private inputs are independently distributed; and
- privately keyed if  $\text{setup}$  returns private input for both parties whose distributions are dependent.

When we present protocols, we will often omit the formality of defining explicit functions  $\mathcal{P} = (\text{setup}, \text{next}_A, \text{next}_B, \text{return}_B)$  and states for the parties, preferring instead to use diagrams as in Figure 4. But, this formulation is convenient for our security definitions. In particular, we present the relevant subprocedures for our security games in Figure 2. An adversary plays the game depicted in Figure 2 by first calling  $\text{initialize}$  (receiving as output  $\pi$ ) and then making various calls to the

other subprocedures. Each time it calls a subprocedure, it receives any output from the procedure. The game ends when the adversary calls `finalize`, and the adversary wins if and only if the output of `finalize` is one.

<pre> <b>proc.</b> initialize(<math>1^\lambda</math>) <math>(\sigma_A, \sigma_B, \pi) \xleftarrow{\\$} \text{setup}(1^\lambda)</math> <math>S_A \leftarrow (\sigma_A, \pi); S_B \leftarrow (\sigma_B, \pi)</math> <math>\text{sid}^* \leftarrow \perp; \text{compromised} \leftarrow \text{false}</math> <math>b \xleftarrow{\\$} \{0, 1\}</math> <b>OUTPUT</b> <math>\pi</math>  <b>proc.</b> finalize(<math>b^*</math>) <b>IF</b> <math>b = b^*</math>, <b>RETURN</b> 1 <b>ELSE</b>, <b>RETURN</b> 0  <b>proc.</b> start-run(<math>\text{sid}, m</math>) <b>IF</b> <math>\text{sid} \notin S_A</math>, <math>S_A.\text{add}(\text{sid}, m)</math>  <b>proc.</b> start-challenge(<math>\text{sid}, m_0, m_1</math>) <b>IF</b> <math>\text{sid} \notin S_A</math> <b>AND</b> <math>\text{sid}^* = \perp</math>,    <math>\text{sid}^* \leftarrow \text{sid}</math>    <math>S_A.\text{add}(\text{sid}, m_b)</math> </pre>	<pre> <b>proc.</b> get-next<sub>A</sub>(<math>\text{sid}, M</math>) <b>IF</b> <b>compromised</b>,    <b>OUTPUT</b> <math>\perp</math> <b>OUTPUT</b> next<sub>A</sub>(<math>S_A, \text{sid}, M</math>)  <b>proc.</b> get-next<sub>B</sub>(<math>\text{sid}, M</math>) <b>IF</b> <b>compromised</b>,    <b>OUTPUT</b> <math>\perp</math> <b>OUTPUT</b> next<sub>B</sub>(<math>S_B, \text{sid}, M</math>)  <b>proc.</b> get-output<sub>B</sub>(<math>\text{sid}</math>) <b>IF</b> <math>\text{sid} = \text{sid}^*</math> <b>OR</b> <b>compromised</b>,    <b>OUTPUT</b> <math>\perp</math> <b>OUTPUT</b> return<sub>B</sub>(<math>S_B, \text{sid}</math>)  <b>proc.</b> get-secrets <b>compromised</b> <math>\leftarrow</math> <b>true</b> <b>OUTPUT</b> <math>(\sigma_A, \sigma_B)</math> </pre>
---	--

**Fig. 2:** Procedures used to define security for message-transmission protocol  $\mathcal{P} = (\text{setup}, \text{next}_A, \text{next}_B)$ . An adversary plays this game by first calling `initialize` and then making various oracle calls. The game ends when the adversary calls `finalize`, and the output of `finalize` is one if the adversary wins and zero otherwise.

The below definitions capture formally the intuitive notions of security that we presented above. In particular, the CPA security definition allows the adversary to start arbitrarily many concurrent runs of the protocol with adversarial input, but it does not allow the adversary to change the messages sent by the two parties or to send its own messages. We also define forward secrecy, which requires that security hold even if the parties’ secret keys may be leaked to the adversary.

**Definition 7 (Message-transmission security).** *A message-transmission protocol is called*

- chosen-plaintext secure (*CPA-secure*) if no PPT adversary has non-negligible advantage in the game presented in Figure 2 when `get-nextA(sid, M)` and `get-nextB(sid, M)` output  $\perp$  unless this is the first `get-next` call with this `sid` or  $M$  is the output from the previous `get-nextA` call with the same `sid` or the previous `get-nextB` with the same `sid` respectively (i.e., the adversary is passive); and
- chosen-ciphertext secure (*CCA-secure*) if no PPT adversary has non-negligible advantage in the game presented in Figure 2 with access to all oracles.

We say that the protocol is chosen-plaintext (resp. chosen-ciphertext) secure without forward secrecy if the above holds without access to the `get-secrets` oracle.

We note that it does not make sense to consider chosen-ciphertext security when Bob may be corrupted. In this case, the output of `get-outputB` could be arbitrary. (Note that the firewall can potentially “sanitize” Bob’s *messages*, but it of course does not have access to his output *output*.) We therefore only consider firewalls that preserve CPA security for Bob.

### 2.3 Key agreement

Key-agreement protocols will play a central role in our constructions, so we now provide a definition of key agreement that suffices for our purposes. Our notion of key agreement closely mirrors the definitions from the previous section.

**Definition 8 (Key agreement).** A key-agreement protocol is represented by five algorithms,  $\mathcal{P} = (\text{setup}, \text{next}_A, \text{next}_B, \text{return}_A, \text{return}_B)$ .  $\text{setup}$  takes as input  $1^\lambda$ , where  $\lambda$  is the security parameter and returns the starting states for each party,  $S_A, S_B$ , which consists of public input  $\pi$  and the private input for each party  $\sigma_A$  and  $\sigma_B$ . Each party's  $\text{next}$  procedure is a stateful algorithm that takes as input  $\text{sid}$  and an incoming message, updates the party's state, and returns an outgoing message. Each party's  $\text{return}$  procedure takes as input the relevant party's state and  $\text{sid}$  and returns the party's final output from some key space  $\mathcal{K}$  or  $\perp$ . We also allow auxiliary input  $\text{aux}$  to be added to Alice's state before the first message of a protocol is sent.

The protocol is correct if Alice and Bob always output the same thing at the end of the run of a protocol for any random coins and auxiliary input  $\text{aux}$ .

We say that a key-agreement protocol is

- unkeyed if  $\text{setup}$  does not return any private input  $\sigma_A$  or  $\sigma_B$ ;
- singly keyed if  $\text{setup}$  returns private input  $\sigma_B$  for Bob but no private input  $\sigma_A$  for Alice; and
- publicly keyed if  $\text{setup}$  returns private input for both parties  $\sigma_A$  and  $\sigma_B$ .

<pre> <b>proc.</b> initialize(<math>1^\lambda</math>) (<math>\sigma_A, \sigma_B, \pi</math>) <math>\stackrel{\\$}{\leftarrow}</math> setup(<math>1^\lambda</math>) <math>S_A \leftarrow (\sigma_A, \pi)</math> <math>S_B \leftarrow (\sigma_B, \pi)</math> <math>\text{sid}^* \leftarrow \perp</math> <math>\text{compromised} \leftarrow \text{false}</math> <math>b \stackrel{\\$}{\leftarrow} \{0, 1\}</math> <b>OUTPUT</b> <math>\pi</math>  <b>proc.</b> finalize(<math>b^*</math>) <b>IF</b> <math>b = b^*</math>,     <b>RETURN</b> 1 <b>ELSE, RETURN</b> 0         </pre>	<pre> <b>proc.</b> get-next<math>_A</math>(<math>\text{sid}, M</math>) <b>IF NOT</b> compromised,     <b>OUTPUT</b> next<math>_A</math>(<math>S_A, \text{sid}, M</math>)  <b>proc.</b> get-next<math>_B</math>(<math>\text{sid}, M</math>) <b>IF NOT</b> compromised,     <b>OUTPUT</b> next<math>_B</math>(<math>S_B, \text{sid}, M</math>)  <b>proc.</b> get-output<math>_A</math>(<math>\text{sid}</math>) <b>IF</b> compromised, <b>OUTPUT</b> <math>\perp</math> <b>IF</b> <math>\text{sid} = \text{sid}^*</math> <b>AND</b> <math>b = 0</math>,     <b>IF</b> return<math>_A</math>(<math>S_A, \text{sid}</math>) = <math>\perp</math>, <b>OUTPUT</b> <math>\perp</math>     <b>ELSE, OUTPUT</b> <math>R_{\text{sid}}</math> <b>ELSE, OUTPUT</b> return<math>_A</math>(<math>S_A, \text{sid}</math>)  <b>proc.</b> get-output<math>_B</math>(<math>\text{sid}</math>) <b>IF</b> compromised, <b>OUTPUT</b> <math>\perp</math> <b>IF</b> <math>\text{sid} = \text{sid}^*</math> <b>AND</b> <math>b = 0</math>,     <b>IF</b> return<math>_B</math>(<math>S_B, \text{sid}</math>) = <math>\perp</math>, <b>OUTPUT</b> <math>\perp</math>     <b>ELSE, OUTPUT</b> <math>R_{\text{sid}}</math> <b>ELSE, OUTPUT</b> return<math>_B</math>(<math>S_B, \text{sid}</math>)  <b>proc.</b> get-secrets <b>compromised</b> <math>\leftarrow</math> true <b>OUTPUT</b> (<math>\sigma_A, \sigma_B</math>)         </pre>
<pre> <b>proc.</b> start-run(<math>\text{sid}, \text{aux}</math>) <b>IF</b> <math>\text{sid} \notin S_A</math>,     <math>S_A.\text{add}(\text{sid}, \text{aux})</math>  <b>proc.</b> start-challenge(<math>\text{sid}, \text{aux}</math>) <b>IF</b> <math>\text{sid}^* = \perp</math> <b>AND</b> <math>\text{sid} \notin S_A</math>,     <math>\text{sid}^* \leftarrow \text{sid}</math>     <math>\mathcal{R}_{\text{sid}^*} \stackrel{\\$}{\leftarrow} \mathcal{K}</math>     <math>S_A.\text{add}(\text{sid}, \text{aux})</math>         </pre>	

**Fig. 3:** Procedures used to define security for key-agreement protocol  $\mathcal{P} = (\text{setup}, \text{next}_A, \text{next}_B, \text{return}_A, \text{return}_B)$ . An adversary plays this game by first calling  $\text{initialize}$  and then making various oracle calls. The game ends when the adversary calls  $\text{finalize}$ , and the output of  $\text{finalize}$  is one if the adversary wins and zero otherwise. We suppress the auxiliary input  $\text{aux}$  when it is irrelevant.

**Definition 9 (Key-agreement security).** *A key-agreement protocol is*

- secure against passive adversaries *if no probabilistic polynomial-time adversary has non-negligible advantage in the game presented in Figure 3 when  $\text{get-next}_A(\text{sid}, M)$  and  $\text{get-next}_B(\text{sid}, M)$  output  $\perp$  unless this is the first  $\text{get-next}$  call with this  $\text{sid}$  or  $M$  is the output from the previous  $\text{get-next}_B$  call with the same  $\text{sid}$  or the previous  $\text{get-next}_A$  call with the same  $\text{sid}$  respectively (i.e., the adversary is passive);*
- secure against active adversaries for Alice *if no probabilistic polynomial-time algorithm has non-negligible advantage in the game presented in Figure 3 without access to the  $\text{get-output}_A$  oracle;*
- secure against active adversaries for Bob *if no probabilistic polynomial-time algorithm has non-negligible advantage in the game presented in Figure 3 without access to the  $\text{get-output}_B$  oracle;*  
*and*
- secure against active adversaries *if it is secure against active adversaries for both Bob and Alice;*  
*and*
- authenticated for Bob *if no probabilistic polynomial-time algorithm playing the game presented in Figure 3 can output a valid transcript with corresponding session id  $\text{sid}$  unless  $\text{return}_B(S_B, \text{sid}) \neq \perp$  or  $\text{compromised} = \text{true}$ . (I.e., it is hard to find a valid transcript unless Bob returns a key.) Furthermore, if the transcript is valid and  $\text{get-output}_A(\text{sid}) \neq \perp$  then  $\text{get-output}_A(\text{sid}) = \text{return}_B(\text{sid})$ . (I.e., if the transcript is valid and Alice outputs a key, then Bob outputs the same key.)*

Note that these definitions are far from standard. In particular, in the case of active adversaries, we define security for Alice in terms of the keys that *Bob* outputs and security for Bob in terms of the keys that *Alice* outputs. This may seem quite counterintuitive. But, in our setting, we are worried that Alice may be corrupted. In this case, we cannot hope to restrict Alice’s output after she receives invalid messages. (The firewall can modify Alice’s *messages*, but not her *output*.) So, the best we can hope for is that the firewall prevents a tampered implementation of Alice (together with an active adversary) from “tricking” Bob into returning an insecure key.

## 2.4 Reverse firewalls and message transmission

Now that we have formal definitions, we discuss the relationship between message-transmission protocols, reverse firewalls, and functionality-maintaining adversaries. We start with an obvious observation: a firewall cannot allow Alice to non-trivially communicate with Bob while simultaneously preventing an arbitrarily corrupted implementation of Alice from communicating something nefarious to Bob! In the language of [MS15] and Section 2.1, a reverse firewall for Alice in a message-transmission protocol cannot *strongly* resist exfiltration against Bob. Such basic impossibility results are bound to arise when we consider notions of security for parties that may be corrupted!

However, if we only consider corrupted implementations of Alice that *maintain the functionality* of this message-transmission protocol, then it is technically possible for a firewall to resist exfiltration against Bob. ([MS15] refer to this as exfiltration resistance against Bob, but not *strong* exfiltration resistance.) The functionality of a message-transmission protocol requires Bob to output (i.e., “receive”) the plaintext that Alice is given as input. So, a functionality-maintaining corruption of Alice must “send” *this* message. Intuitively, a firewall that resists exfiltration against Bob for this restricted class of corruptions “allows Alice to send her message to Bob, but nothing else.” In other words, these firewalls prevent steganography. All of our reverse firewalls easily achieve this, and in some sense this is the “best possible” notion of exfiltration resistance against Bob.

Of course, for message-transmission protocols, we are primarily concerned with security against a third party adversary. We say that a firewall *preserves security* if it protects Alice from such an

adversary in some way.<sup>4</sup> Ideally, we would like to build reverse firewalls for Alice that protect her from such an adversary even if Alice has been arbitrarily corrupted.

In order to do this, we have to worry about an “attack by refusal.” E.g., a corrupted implementation of Alice might simply refuse to communicate at all (in our model, she can send the special “abort” symbol  $\perp$ ) when the first bit of her bank account number is zero and communicates normally otherwise. Clearly, a firewall can only protect against such attacks if it can “spoofer” messages from Alice, so that a third party adversary cannot tell the difference between a legitimate message sent from Alice and a “spoofed” message from the firewall. Now, suppose Alice has a secret key that is used to publicly sign something at some point in the protocol. It should be clear that the strong notion of security preservation discussed above (which [MS15] simply call strong security preservation) is not possible for any such protocols. So, our firewalls that work in the setting in which Alice has a public-key/private-key pair preserve Alice’s security, but they do not *strongly* preserve it. (See Section 4.) In the “mixed” setting in which Bob has a public-key/private-key pair but Alice does not, we show that strong security preservation is equivalent to a strengthened notion of rerandomizable RCCA-secure encryption, a primitive that is currently being actively studied. (See Section 5.) In the unkeyed setting, we do in fact achieve strong security with a simple protocol. (See Section 3.)

Finally, while the barrier discussed above seems inherent, one might hope that this “attack by refusal” is in some sense the only attack possible. Perhaps, for example, Alice can only communicate a single bit to the adversary per run of the message-transmission protocol. We could try to formalize a security notion that captures this, but it is clear that our more advanced protocols would not meet it. In particular, our more advanced protocols use the symmetric-key protocol from [BPR14a] as a subprotocol. The protocol from [BPR14a] is completely deterministic and suffers from an attack in which a corrupted implementation of Alice can simply try to encrypt many different messages until she finds one whose corresponding ciphertext has a specific form (say, one that starts with many zeros). Obviously, this channel can be used to communicate with an eavesdropper, and both of our key-agreement-based protocols inherit this weakness, including our strongest result from Section 4.1. (We note, however, that this attack only works against our protocols if the corrupted implementation chooses which message to send *after* the key agreement is complete.) We leave it to future work to find a stronger notion of security that is still achievable (but, we stress that we already achieve a remarkably strong notion of security with surprisingly efficient protocols). It seems that something might be possible here, but we might need to sacrifice efficiency by applying (typically very slow) public-key operations to (potentially very long) plaintexts.

### 3 The unkeyed setting

In this section we investigate the scenario in which Alice and Bob have neither a shared secret key nor each other’s public keys. Note that we of course can only achieve CPA security here because neither Alice nor Bob can be authenticated. As such, this section’s primary purpose is to illustrate core concepts and to act as a useful stepping stone toward the more complicated and stronger protocols that we present in Section 4. Two approaches are feasible: a simple two-round solution in which Bob (the receiver) sends a fresh public key to Alice and receives an encryption of her plaintext under this public key in reply, and a more flexible and efficient protocol based on key agreement.

---

<sup>4</sup> It is easy to see that a firewall for Alice in a message-transmission protocol is exfiltration resistant against an eavesdropper if and only if it preserves semantic security—the weakest notion of security that interests us. So, we may safely ignore exfiltration resistance against eavesdroppers in this context.

### 3.1 A two-round protocol from rerandomizable encryption

We first consider the simple case of CPA-secure two-round schemes in which the first message is a public key chosen randomly by Bob and the second message is an encryption of Alice’s plaintext under this public key. Figure 4 shows the protocol.

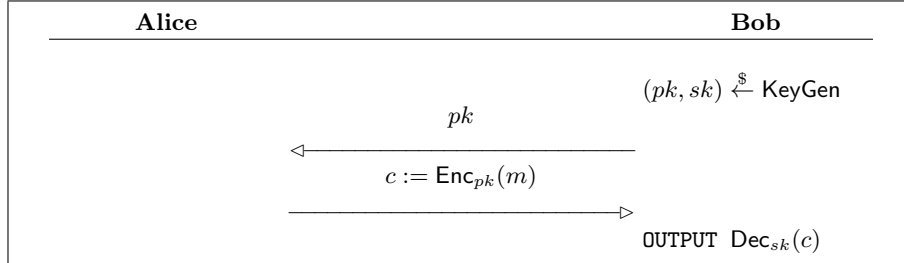


Fig. 4: Two round message-transmission protocol using the public-key encryption scheme (KeyGen, Enc, Dec).

In order to provide a reverse firewall for Alice in this protocol, the encryption scheme must be *rerandomizable*. In order to provide a reverse firewall for Bob, the scheme must be *key malleable*. Intuitively, a scheme is key malleable if a third party can “rerandomize” a public key and map ciphertexts under the “rerandomized” public key to ciphertexts under the original public key. For completeness, we include formal definitions in Appendix B, and we observe there that ElGamal encryption suffices.

**Definition 10 (Public-key encryption).** A public-key encryption scheme is a triple of efficient algorithms (KeyGen, Enc, Dec). KeyGen takes as input  $1^\lambda$  where  $\lambda$  is the security parameter and outputs a public-key/private-key pair,  $(pk, sk)$ . Enc takes as input the public key and a plaintext  $m$  from some plaintext space  $\mathcal{M}$  and outputs a ciphertext  $c$  from some ciphertext space  $\mathcal{C}$ . Dec takes as input a ciphertext and the private key and outputs a plaintext or the special symbol  $\perp$ . We sometimes omit the keys from the input to Enc and Dec and the security parameter input to KeyGen. The scheme is correct if  $\text{Dec}(\text{Enc}(m)) = m$  for all  $m \in \mathcal{M}$ . The scheme is semantically secure if for any adversarially chosen pair of plaintexts  $(m_0, m_1)$ ,  $\text{Enc}(m_0)$  is computationally indistinguishable from  $\text{Enc}(m_1)$ .

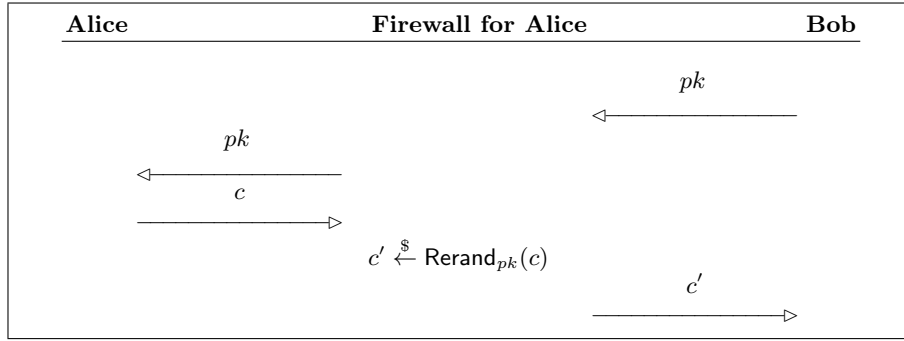
**Definition 11 (Rerandomizable encryption).** A semantically secure public-key encryption scheme is rerandomizable if there is an efficient algorithm Rerand (with access to the public key) such that for any ciphertext  $c$  such that  $\text{Dec}(c) \neq \perp$ , we have  $\text{Rerand}(\text{Dec}(c)) = \text{Dec}(c)$ , and the pair  $(c, \text{Rerand}(c))$  is computationally indistinguishable from  $(c, \text{Rerand}(\text{Enc}(0)))$ . We say that it is strongly rerandomizable if the previous property holds even when  $\text{Dec}(c) = \perp$ .

**Definition 12 (Key malleability).** A public-key encryption scheme is key malleable if (1) the output of KeyGen is distributed uniformly over the space of valid keys; (2) for each public key  $pk$  there is a unique associated private key  $sk$ ; and (3) there is a pair of efficient algorithms KeyMaul and CKeyMaul that behave as follows. KeyMaul is a randomized algorithm that takes as input a public key  $pk$  and returns a new public key  $pk'$  whose distribution is uniformly random over the public key space and independent of  $pk$ . Let  $(sk, pk)$  be a private key/public key pair. Let  $(sk', pk')$  be the unique pair associated with randomness  $r$  such that  $pk' = \text{KeyMaul}(pk; r)$ . Then, CKeyMaul takes as input a ciphertext  $c$  and randomness  $r$  and returns  $c'$  such that  $\text{Dec}_{sk'}(c) = \text{Dec}_{sk}(c')$ . We suppress the

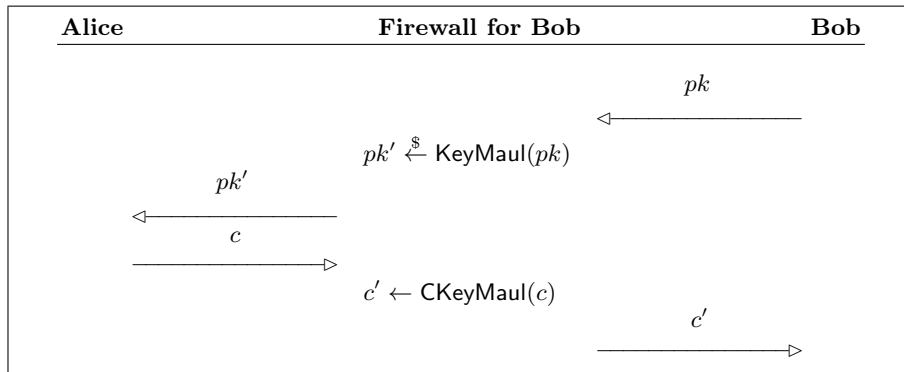
input  $r$  when it is understood. Furthermore, we require that **KeyMaul** outputs a uniformly random key  $pk'$  if called on input that is not in the public-key space.

**Example.** It is well-known that ElGamal encryption [ELG85] is both key malleable and strongly rerandomizable. In particular, given an ElGamal public-key  $(g, h)$  over a group of order  $p$ , a ciphertext  $(u, v)$  can be rerandomized by applying the operation  $(u, v) \rightarrow (g^r u, h^r v)$  where  $r \xleftarrow{\$} \mathbb{Z}_p^*$  is chosen uniformly at random. The public key can be mauled by applying the operation  $(g, h) \rightarrow (g^\alpha, h^\beta)$  where  $(\alpha, \beta) \xleftarrow{\$} (\mathbb{Z}_p^*)^2$  are chosen uniformly and independently at random. Finally, a ciphertext  $(u, v)$  under key  $(g^\alpha, h^\beta)$  can be converted into a ciphertext under  $(g, h)$  by applying the operation  $(u, v) \rightarrow (u^{\beta/\alpha}, v)$ .

If the underlying encryption scheme in Figure 4 is rerandomizable, then we can build a reverse firewall for Alice as in Figure 5. If it is key malleable, then we can build a reverse firewall for Bob as in Figure 6. The following theorem shows that this protocol and its reverse firewalls are secure.



**Fig. 5:** Reverse firewall for Alice for the protocol shown in Figure 5 that works if the encryption scheme is rerandomizable.



**Fig. 6:** Reverse firewall for Bob for the protocol shown in Figure 4 that works if the encryption scheme is key malleable. We suppress the randomness  $r$  used as input to **KeyMaul** and **CKeyMaul**.

**Theorem 1.** *The unkeyed message-transmission protocol shown in Figure 4 is CPA-secure if the underlying encryption scheme is semantically secure. If the scheme is also (strongly) rerandomizable,*



then the reverse firewall shown in Figure 5 (strongly) preserves security for Alice and (strongly) resists exfiltration for Alice. If the scheme is key malleable, then the reverse firewall shown in Figure 6 maintains functionality for Bob, strongly preserves Bob’s security, and strongly resists exfiltration for Bob against Alice.

*Proof.* It is a common folklore result that the underlying protocol (i.e., the protocol without reverse firewalls) is CPA-secure. It is clear that the two firewalls maintain functionality.

Security and exfiltration resistance of Bob’s firewall follows from the definition of key malleability. In particular, for any tampered implementation  $\tilde{B}$  of Bob, after the post-processing by the reverse firewall, the key  $pk'$  is uniformly random by the definition of key malleability, regardless of the behavior of  $\tilde{B}$ . This implies exfiltration resistance, and security then follows from the fact that the underlying protocol is secure when the key is chosen legitimately.

Consider a tampered implementation of Alice  $\tilde{A}$  that maintains functionality, and let  $c$  be the output of Alice. Since Alice maintains functionality,  $\text{Dec}(c) \neq \perp$ . So, by the definition of rerandomizability, the output of  $\text{Rerand}(c)$  is indistinguishable from  $\text{Enc}_{pk}(0)$ . The security preservation and exfiltration resistance of Alice’s firewall follows. A similar argument shows that strong rerandomizability implies strong security preservation and strong exfiltration resistance.  $\square$

### 3.2 Hybrid encryption fails.

A major drawback of the above scheme is that it requires public-key operations of potentially very long plaintexts, which can be very inefficient in practice. A common solution in the classical setting is to use *hybrid encryption*, in which  $\text{Enc}_{pk}(m)$  is replaced by  $(\text{Enc}_{pk}(rk), \text{SEnc}_{rk}(m))$ , where  $\text{SEnc}$  is some suitable symmetric-key encryption scheme and  $rk$  is a freshly chosen uniformly random key for  $\text{SEnc}$ . However, if we simply replace the public-key encryption in Figure 4 with the corresponding hybrid-key encryption scheme, then this fails spectacularly. For example, a tampered version of Alice  $\tilde{A}$  can choose some *fixed* secret key  $rk^*$  and send the message  $(\text{Enc}_{pk}(rk^*), \text{SEnc}_{rk^*}(m))$ . If  $rk^*$  is a valid key, then  $\tilde{A}$  maintains functionality, but an adversary that knows  $rk^*$  can of course read any messages that Alice sends.

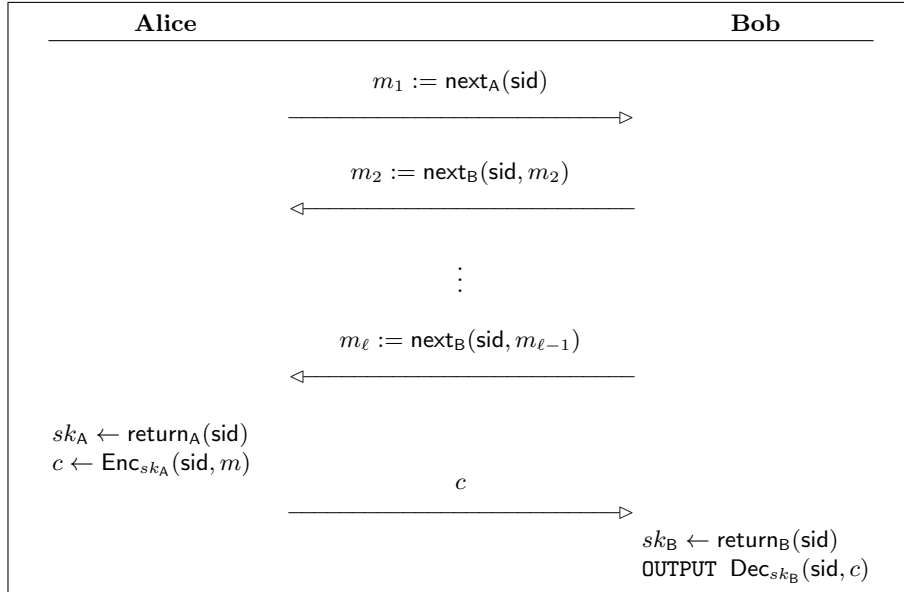
### 3.3 A solution using key agreement

Recall that we are in the unkeyed setting, so we are still interested in CPA security. (We address CCA security in the next section.) The protocol from Section 3.1 works, but it requires a public-key operation on the plaintext, which may be very long. In practice, this can be very inefficient. Above, we show that one common solution to this problem in the classical setting, hybrid encryption, seems hopeless with reverse firewalls because it allows Alice to *choose* a key  $rk$  that will be used to encrypt the plaintext—thus allowing a tampered version of Alice to “choose a bad key.”

So, we instead consider an alternative common solution to this efficiency problem: key agreement followed by symmetric-key encryption. (See Figure 7.) As in Appendix A, we use a nonce-based encryption scheme with unique ciphertexts. We can view this as a modification of hybrid encryption in which “Alice and Bob together choose the key  $rk$ ” that will be used to encrypt the plaintext. More importantly from our perspective, the messages that determine the key will go through the firewall. As an added benefit, once a key is established, Alice can use it to efficiently send multiple messages, not just one, without any additional public-key operations (though we do not model this here).<sup>5</sup>

<sup>5</sup> To keep our definitions relatively simple, we only formally model the case in which Alice wishes to send a single message to Bob per key agreement, and we only formally prove security in this model. We note in passing that this weaker model actually only requires semantically secure symmetric-key encryption. The same is true of Theorem 4.

The composition theorem below shows that this protocol can in fact have a reverse firewall for both parties, provided that the key-agreement protocol itself has a reverse firewall that satisfies some suitable security requirements. See Appendix C for the proof. In the next section, we construct such a protocol.



**Fig. 7:** The message-transfer protocol obtained by combining a key-agreement scheme ( $\text{setup}$ ,  $\text{next}_A$ ,  $\text{next}_B$ ,  $\text{return}_A$ ,  $\text{return}_B$ ) and a nonce-based encryption scheme,  $(\text{Enc}, \text{Dec})$ .

**Theorem 2 (Composition theorem for CPA security).** *Let  $\mathcal{W}_A$  and  $\mathcal{W}_B$  be reverse firewalls in the underlying key-agreement protocol in Figure 7 for Alice and Bob respectively. Let  $\mathcal{W}_A^*$  be the firewall for Alice in the full protocol in Figure 7 obtained by applying  $\mathcal{W}_A$  to the key-agreement messages and then letting the last message through if  $\mathcal{W}_A$  does not output  $\perp$  and replacing the last message by  $\perp$  otherwise. Let  $\mathcal{W}_B^*$  be the firewall for Bob in the full protocol in Figure 7 obtained by applying  $\mathcal{W}_B$  to the key-agreement messages and simply letting the last message through. Then,*

1. *the protocol in Figure 7 is CPA-secure if the underlying key-agreement protocol is secure against passive adversaries and the underlying nonce-based encryption scheme is CPA-secure;*
2.  *$\mathcal{W}_B^*$  preserves CPA security if  $\mathcal{W}_B$  preserves security of the key-agreement protocol; and*
3.  *$\mathcal{W}_A^*$  preserves CPA security if the encryption scheme has unique ciphertexts and  $\mathcal{W}_A$  preserves semantic security and is exfiltration resistant against Bob.*

Finally, we note that strong security preservation is not possible for this protocol (at least for Alice).

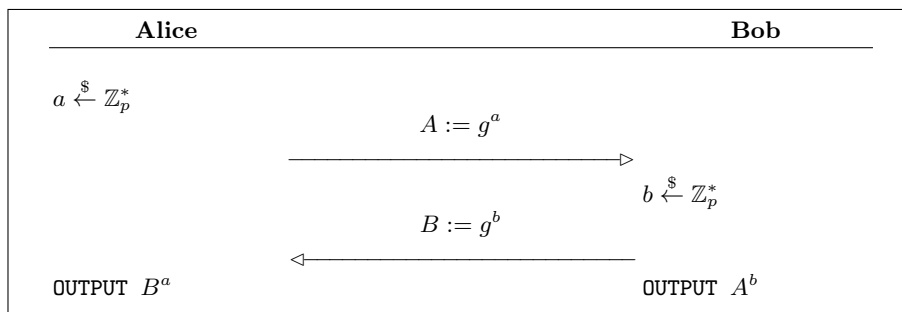
*Remark 1 (Informal).* There is no reverse firewall for Alice in the protocol illustrated in Figure 7 that maintains functionality and strongly preserves Alice’s security.

*Proof.* Consider the tampered implementation of Alice  $\bar{A}$  that goes through the key-agreement protocol as normal and then sends as its last message  $c := \text{Enc}_{sk'}(\text{sid}, m)$  for some fixed key  $sk'$  chosen by the adversary by simulating a run of the key-agreement protocol. Since the firewall

maintains functionality, it must be the case that the message sent by the firewall to Bob  $c'$  satisfies  $\text{Dec}_{sk'}(\text{sid}, c') = m$ . So, the adversary can decrypt the message itself. Clearly, this protocol is not secure (by any reasonable definition).  $\square$

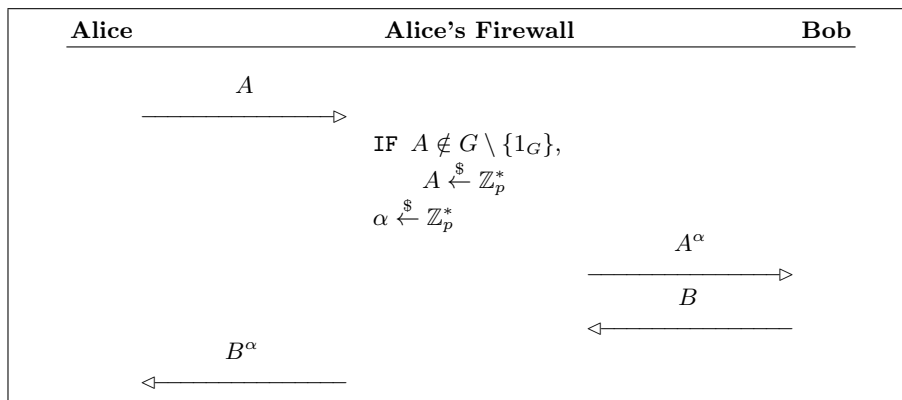
### 3.4 Key agreement secure against passive adversaries

Theorem 2 motivates the study of unkeyed key-agreement protocols with reverse firewalls that preserve security against passive adversaries. In the classical setting (i.e., without reverse firewalls), the canonical example is the elegant key-agreement protocol of Diffie and Hellman [DH06], shown in Figure 8, whose security follows immediately from the hardness of DDH over the base group  $G$ . We use this as an example to illustrate the basic idea of a reverse firewall in the key-agreement setting.



**Fig. 8:** Diffie-Hellman key agreement over a group  $G$  of prime order  $p$  with generator  $g$ .

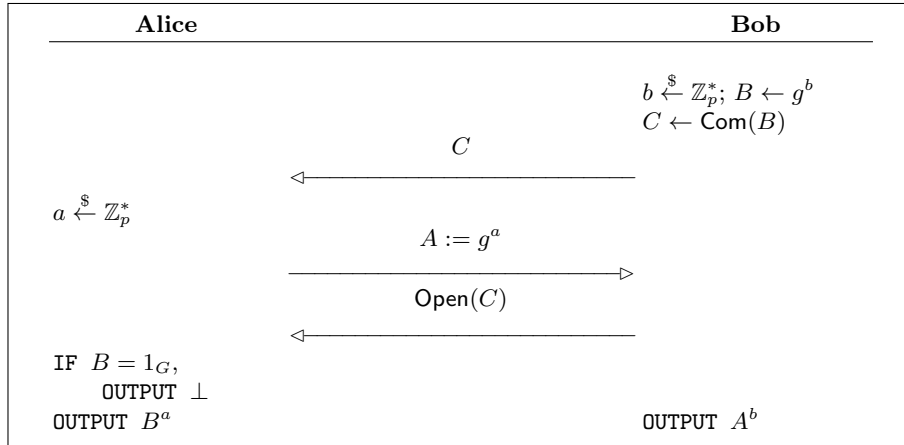
Diffie-Hellman key agreement has a simple reverse firewall for Alice, which raises both messages to a single random power,  $\alpha \in \mathbb{Z}_p^*$ . We present this reverse firewall in Figure 9. Note that this firewall effectively replaces Alice's message with a uniformly random message. Security then follows from the security of the underlying protocol, since the transcript and resulting key in the two cases are distributed identically. This very simple idea of rerandomizing Diffie-Hellman key agreement is part of all of our key-agreement protocols.



**Fig. 9:** Reverse firewall for Alice in the protocol from Figure 8.

But, this protocol cannot have a reverse firewall that maintains correctness and preserves security for Bob. Consider the tampered implementation  $\tilde{B}$  that repeatedly samples  $b \xleftarrow{\$} \mathbb{Z}_p^*$  and computes  $A^b$  until the first bit of the result is zero.  $\tilde{B}$  then sends  $g^b$  to Alice as normal. If the firewall maintains correctness, then Alice’s output must be  $A^b$ , but this output is clearly distinguishable from random—its first bit is zero! This is a problem that we must overcome if we want to build protocols that have secure firewalls for *both* parties. In particular, we run the risk that one party has the ability to “accept or reject” any key.

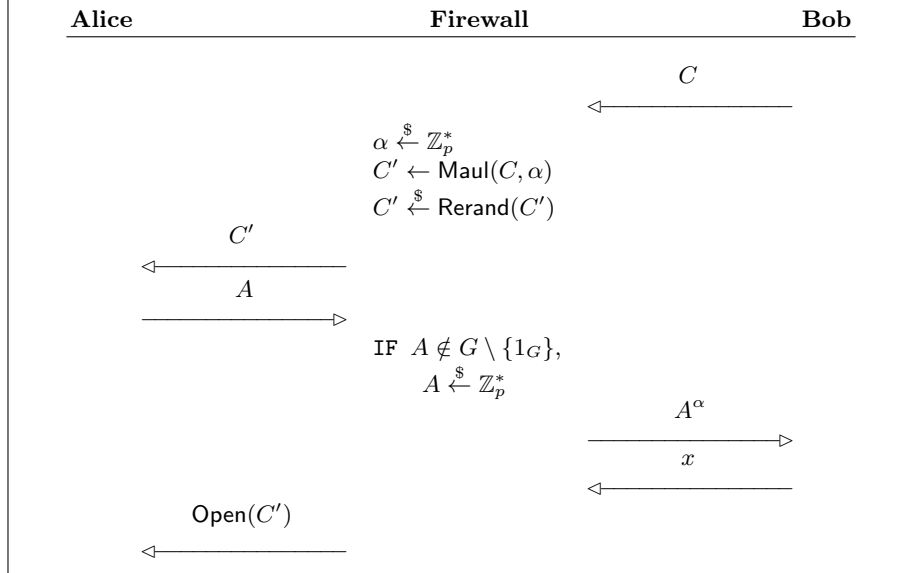
To solve this problem, we add an additional message to the beginning of the protocol in which Bob commits to the message that he will send later. Of course, in order to permit a secure firewall, the commitment scheme itself must be both malleable (so that the firewall can rerandomize the underlying message that Bob has committed to, mapping a commitment of  $B$  to a commitment of  $B^\alpha$ ) and rerandomizable (so that the randomness used by Bob to commit and open will not leak any information about his message). To achieve our strongest level of security, we also need the scheme to be statistically hiding and for each commitment to have a unique opening for a given message. (These requirements are easily met in practice. For example, a simple variant of the Pedersen commitment suffices [Ped92]. For completeness, we present such a scheme in Appendix D.) The protocol is shown in Figure 10. In Figure 11, we present a single reverse firewall for this protocol that happens to work for either party. (Each party would need to deploy its own version of this firewall to guarantee its own security. It just happens that each party’s firewall would have the same “code.”)



**Fig. 10:** A variant of Diffie-Hellman key agreement over a group  $G$  of prime order  $p$  with public generator  $g$ . ( $\text{Com}, \text{Open}$ ) is a commitment scheme.

**Theorem 3.** *The protocol in Figure 10 is secure against passive adversaries if DDH is hard in  $G$ . The reverse firewall  $\mathcal{W}$  in Figure 11 is functionality maintaining. If the commitment scheme is statistically hiding, then  $\mathcal{W}$  preserves security for Alice and is strongly exfiltration resistant against Bob. If the commitment scheme is computationally binding, then  $\mathcal{W}$  is exfiltration resistant for Bob against Alice and preserves security for Bob.  $\mathcal{W}$  also detects failure for both parties.*

*Proof.* It is clear that the underlying protocol is secure provided that DDH is hard in  $G$ . It is also clear that  $\mathcal{W}$  maintains functionality and fails detectably for both parties. (Here, we assume that the firewall outputs  $\perp$  if it ever receives a malformed message.)



**Fig. 11:** Reverse firewall for either Alice or Bob in the protocol from Figure 8.  $\text{Maul}(C, \alpha)$  takes a commitment  $C = \text{Com}(B)$  and converts it into a commitment of  $B^\alpha$ .  $\text{Rerand}(C)$  takes a commitment  $C = \text{Com}(B)$  and converts it into a uniformly random commitment of  $B$ . We assume that a rerandomized and mauled commitment can be opened with access to an opening of the original commitment (and the randomness used in the rerandomization and mauling).

Note that, after rerandomization and mauling, the commitment  $C'$  is a uniformly random commitment of a uniformly random group element, independent of the original commitment  $C$ . Since Bob is functionality maintaining, his second message is fixed unless he can find an alternative opening for the commitment, which by assumption is computationally hard. It follows that  $\mathcal{W}$  is exfiltration resistant for Bob against Alice and preserves security for Bob.

To prove strong exfiltration resistance for Alice against Bob and security preservation, we again note that Bob's first message is a uniformly random commitment of a uniformly random group element. Since the commitment is statistically binding, it is statistically close to independent from  $\alpha$ , regardless of Bob's choice of  $C$ . Therefore, Alice's message  $A$  is statistically close to independent from  $\alpha$ , and  $A^\alpha$  is statistically close to uniform. The result follows.  $\square$

## 4 The publicly keyed setting

We now consider the publicly keyed setting. I.e., we assume that both Bob and Alice have a public key and a private key (though their private keys must be independent). In the setting of the previous section, with no public-key infrastructure, it is trivially impossible to achieve CCA-security. (An adversary can simply "pretend to be Bob" and read Alice's plaintext.) In this section, we show that a CCA-secure message-transmission protocol with reverse firewalls does in fact exist in the publicly keyed setting. In particular, in Section 4.1, we give the CCA analogue of Theorem 2, showing that a key-agreement protocol that is secure against active adversaries and has sufficiently secure reverse firewalls together with a symmetric-key encryption scheme suffices. As in the previous section, this key-agreement-based protocol has the additional benefit that it is efficient, in the sense that it does not apply public-key operations to the plaintext. In Section 4.2, we construct a key-agreement protocol that suffices.

## 4.1 A solution using key agreement

The theorem below shows that we can build an efficient CCA-secure message-transmission protocol with reverse firewalls in this setting, provided that we have a sufficiently secure key-agreement protocol with reverse firewalls. (Recall that Bob’s reverse firewall can only preserve CPA security. Such a firewall is already given by Theorem 2, so we do not repeat this here.) See Appendix E for the proof.

**Theorem 4 (Composition theorem for CCA security).** *Define  $\mathcal{W}_A$  and  $\mathcal{W}_A^*$  as in Theorem 2. Then,*

1. *the protocol in Figure 7 is CCA-secure if the underlying key-agreement protocol is secure against active adversaries for Alice and the underlying nonce-based encryption scheme is CCA-secure; and*
2.  *$\mathcal{W}_A^*$  preserves CCA-security if the encryption scheme has unique ciphertexts, the key-agreement protocol is authenticated for Bob, and  $\mathcal{W}_A$  preserves security for Alice, is exfiltration resistant against Bob with valid transcripts, and detects failure.*

## 4.2 Key agreement secure against active adversaries

Theorem 4 motivates the study of key-agreement protocols with reverse firewalls that preserve security against active adversaries. In the classical setting, the common solution is essentially for each of the parties to sign the transcript of this run of the protocol. Intuitively, this lets Alice know that she has received unaltered messages from Bob, and vice versa. However, this solution does not work in our setting. In particular, it is important for us that messages *can* be altered without breaking functionality, so that the firewall can rerandomize messages when necessary.

Of course, while we want to allow for the possibility that Alice and Bob see a different transcript but still output a key, we still want them to agree on the key itself. This leads to the idea of signing some deterministic function of the key, so that the signatures can be used to verify that the parties share the same key without necessarily requiring them to share the same transcript. This is the heart of our solution.

However, we also have to worry that the signatures themselves can provide channels, allowing tampered versions of the parties to leak some information. We solve this by using a unique signature scheme, as defined by [MRV99]. These guarantee that there is a unique signature that verifies for each plaintext, so that functionality-maintaining implementations of Alice and Bob have only one option for their signature, given a fixed key. (See [AMV15] for a thorough analysis of signatures in the context of reverse firewalls and corrupted implementations, including alternative ways to implement signatures that would suffice for our purposes.)

Furthermore, in order for our firewall to fail detectably, it has to be able to check the signature itself—so that it can distinguish a valid transcript from an invalid one. So, we would like the parties to sign a deterministic function of  $g^{ab}$  that is efficiently computable given only access to  $g^a$  and  $g^b$ . This leads naturally to the use of a symmetric bilinear map  $e : G \times G \rightarrow G_T$ . The parties then sign  $e(g^a, g^b)$ . Of course,  $g^{ab}$  is no longer indistinguishable from random in the presence of a bilinear map. But, it can be hard to compute. So, we apply a cryptographic hash function  $H$  to  $g^{ab}$  in order to extract the final key  $H(g^{ab})$ .

We now provide two definitions to make this precise.

**Definition 13 (Unique Signatures).** *A unique signature scheme is a triple of algorithms (KeyGen, USig, Ver). KeyGen takes as input  $1^\lambda$  where  $\lambda$  is the security parameter and outputs a public key  $pk$  and a private*

key  $sk$ .  $\text{Sig}$  takes as input the secret key  $sk$  and a plaintext  $m$  and outputs a signature  $\tau$ .  $\text{Ver}$  takes as input the public key  $pk$ , a signature  $\tau$  and a message  $m$  and outputs either `true` or `false`. A signature scheme is correct if  $\text{Ver}_{pk}(\text{USig}_{sk}(m), m) = \text{true}$ . It is unique if for each plaintext  $m$  and public key  $pk$ , there is a unique signature  $\tau$  such that  $\text{Ver}_{pk}(\tau, m) = \text{true}$ .

A signature scheme is secure against adaptive chosen-message existential-forgery attacks if no adversary with access to the public key and a signature oracle can produce a valid signature not returned by the oracle.

We will need to use a group with a symmetric bilinear map in which the following variant of the computational Diffie-Hellman assumption holds.

**Definition 14 (Inverse CDH).** For a group  $G$  of order  $p$ , we say that inverse CDH is hard in  $G$  if no probabilistic polynomial-time adversary taking input  $(g, g^a, g^b)$  where  $g \xleftarrow{\$} G$  and  $(a, b) \xleftarrow{\$} \mathbb{Z}_p^2$  has non-negligible probability of returning  $(h, h^{1/a}, h^{1/b})$  for some element  $h \in G \setminus \{1_G\}$ .

Note that inverse CDH is stronger than the standard CDH assumption, as  $(g^{ab}, g^b, g^a)$  is a solution to inverse CDH.

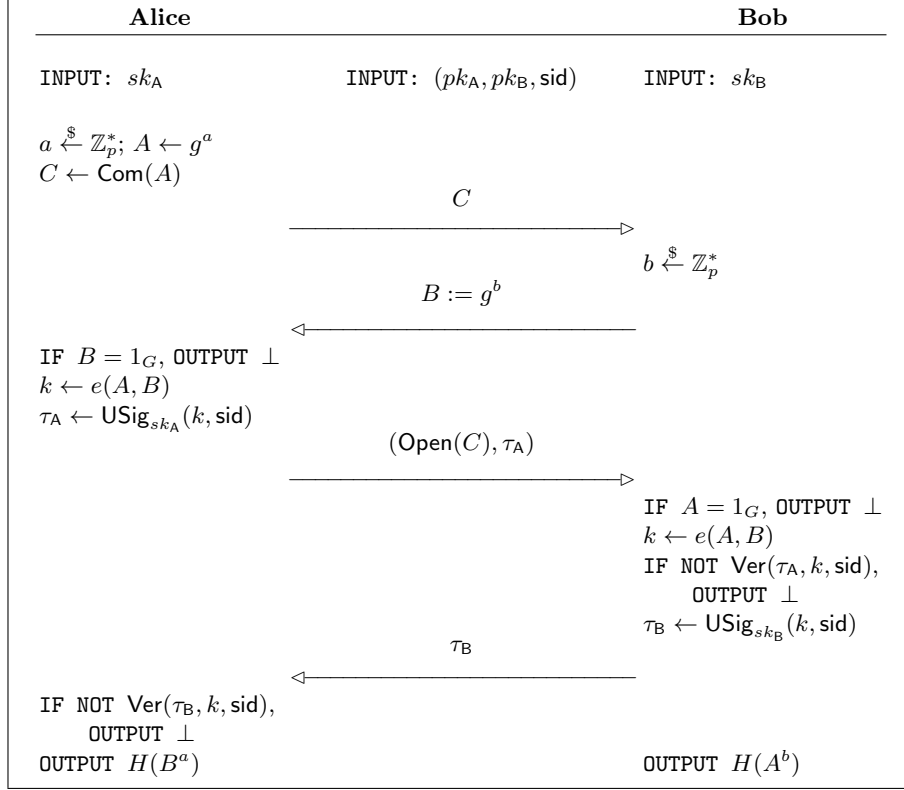
We now present our protocol in Figure 12 with a reverse firewall for both parties in Figure 13. It requires a unique signature scheme  $(\text{USig}, \text{Ver})$  with public keys  $pk_A$  for Alice and  $pk_B$  for Bob and corresponding secret keys  $sk_A$  and  $sk_B$  respectively, a base group  $G$  with generator  $g$  in which inverse CDH is hard, a target group  $G_T$ , and a non-trivial bilinear map between the two groups  $e : G \times G \rightarrow G_T$ . We also need a function  $H : G \rightarrow \{0, 1\}^\ell$  for some polynomially large  $\ell$  that extracts hardcore bits from CDH. Presumably a standard cryptographic hash function will work. For simplicity, we model  $H$  as a random oracle, but we note that the proof can be modified to apply to any function  $H$  such that  $(g^a, g^b, H(g^{ab}))$  is indistinguishable from random. (See [Kil07] for a discussion of such functions.) We stress again that this protocol is remarkably efficient, and we think that it can and should be used in practice.

**Theorem 5.** *The protocol shown in Figure 12 is authenticated for Bob and secure against active adversaries if the signature scheme is secure and inverse CDH is hard in  $G$ . The reverse firewall  $\mathcal{W}$  shown in Figure 13 preserves security against active adversaries for Alice, preserves authenticity, is exfiltration resistant for Alice against Bob with valid transcripts, and detects failure for Alice.  $\mathcal{W}$  also preserves security against active adversaries for Bob, is exfiltration resistant for Bob against Alice with valid transcripts, and detects failure for Bob.*

*Proof.* The fact that the protocol is authenticated for Bob is clear. Informally, any valid transcript must have a valid signature from Bob. By the security of the signature scheme, except with negligible probability, the adversary cannot produce a signature with a given `sid` without Bob returning a key. The same proof shows that the firewall preserves authenticity. The fact that the firewall detects failure for both parties is also clear.

We prove that  $\mathcal{W}$  preserves security against active adversaries for Alice. Let  $\tilde{A}$  be a functionality-maintaining tampered implementation of Alice. Consider the following sequence of games.

- **Game 1** is the key-agreement security game against active adversaries for Alice in the protocol shown in Figure 12 with  $A$  replaced by  $\mathcal{W} \circ \tilde{A}$ .
- **Game 2** is **Game 1** in which the adversary never provides a signature  $\tau_A$  or  $\tau_B$  unless it came from a response from Alice or Bob respectively.
- **Game 3** is **Game 2** in which Alice’s commitment  $C$  is replaced by  $A$  (the value that her commitment opens to), and the firewall’s commitment  $C'$  is replaced by  $A^\alpha$  in all runs of the protocol. We also remove Alice’s second message is simply removed and join Bob’s last two messages into a single message,  $(B, \tau_B)$ .



**Fig. 12:** Authenticated key agreement with a firewall for both parties. USig is a unique signature.

- **Game 4** is **Game 3** in which Alice’s message  $A$  is replaced by a uniformly random group element, and the firewall takes  $\alpha := 1$  (i.e., the firewall leaves all messages alone, besides checking signatures) in all runs of the protocol.

The following claim follows from the definition of adaptive chosen-message existential-forgery security.

**Claim 5.1.** *If the signature scheme is secure, then for any PPT adversary  $\mathcal{E}$ ,  $|\text{Adv}^{\text{Game } 1}(\mathcal{E}) - \text{Adv}^{\text{Game } 2}(\mathcal{E})|$  is negligible.*

It is trivial to take any adversary in **Game 2** and convert it to an adversary in **Game 3** with the same advantage. Similarly, **Game 4** only differs from **Game 3** syntactically. Indeed,  $A' := A^\alpha$  is a uniformly random group element independent of everything else, and  $B' := B^\alpha$  is uniformly random and independent of  $A'$  and everything else. (Recall that Bob is honest.) So, since Alice’s next message is deterministic (because the signature is unique and Alice maintains functionality) and no output from Alice is given to the adversary other than these two messages (recall that the  $\text{get-output}_A$  oracle is not allowed in this game), this is merely a change of variables.

**Claim 5.2.** *No PPT adversary  $\mathcal{E}$  has non-negligible advantage in **Game 4** if inverse CDH is hard in  $G$ .*

*Proof.* We assume without loss of generality that  $\mathcal{E}$  makes no “trivial” oracle calls, whose output can be easily predicted based on prior calls, such as calls to  $\text{get-output}_B$  before the relevant run of the protocol is finished. We build an adversary  $\mathcal{E}'$  in the inverse CDH game as follows. On input





if  $\mathcal{E}$  has non-negligible advantage, it must call the random oracle on the key. Therefore,  $\mathcal{E}'$  can search through its random oracle queries until it finds  $h_3$  satisfying  $e(g_1, h_3) = e(h_1, g_3)$ . It then returns  $(h_3, h_2, h_1)$ . If it finds nothing, it returns  $\perp$ . The result follows from noting that  $(h_3, h_2, h_1)$  is a valid solution to CDH. (5.2) ■

It follows that  $\mathcal{W}$  preserves security against active adversaries for Alice. The proof for Bob is essentially identical. And, essentially the same proof shows exfiltration resistance. □

## 5 The singly keyed setting

Finally, we consider the singly keyed setting, in which Bob has a public-key/private-key pair, but Alice does not. As in the classical case, we show that there exist one-round singly keyed protocols that are not forward secret. These are essentially the single-round analogue of the two-round protocol presented in Figure 4 in Section 3.1. We show that the existence of such a protocol is equivalent to the existence of rerandomizable encryption, and we show how to achieve CCA-security (though not forward secrecy).

### 5.1 One-round CPA-secure protocols

The next theorem shows that one-round CPA-secure protocols with reverse firewalls are equivalent to rerandomizable public-key encryption.

**Theorem 6.** *Any (strongly) rerandomizable semantically secure public-key encryption scheme implies a one-round CPA-secure singly keyed message-transmission protocol without forward security with a reverse firewall that (strongly) preserves security and (strongly) resists exfiltration. Conversely, any one-round CPA-secure message-transmission protocol with a reverse firewall that (strongly) preserves security implies a (strongly) rerandomizable semantically secure public-key encryption scheme.*

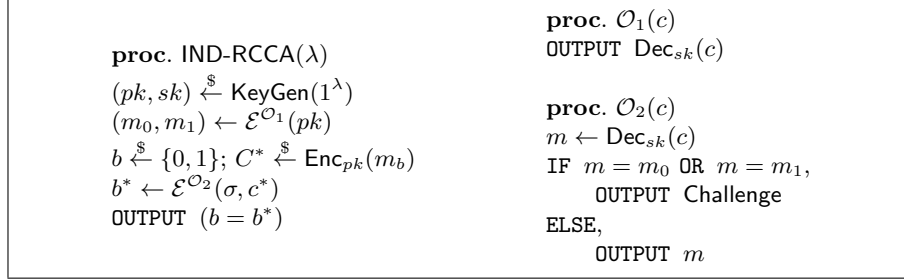
*Proof.* To prove the first statement, we consider the protocol in which Alice simply sends Bob an encryption of the plaintext under Bob’s public key. Alice’s firewall applies the **Rerand** algorithm to the plaintext. (Bob does not need a firewall, since he does not send any messages.) Security of this protocol follows immediately from the security of the encryption scheme, and the fact that the firewall preserves security follows immediately from the rerandomizability of the encryption scheme.

To prove the second statement, consider the following encryption scheme. The key generation algorithm runs the **setup** algorithm of the underlying MTP protocol, receiving as output  $\sigma_A$ ,  $\sigma_B$ , and  $\pi$ . The public key is then  $\pi$  and  $\sigma_A$ , and the private key is  $\sigma_B$ . The encryption algorithm first uses  $\sigma_A$  and an arbitrarily chosen **sid** to compute Alice’s single message in the protocol, given the plaintext. It then applies the reverse firewall to this message; the result is the ciphertext. The rerandomization algorithm simply applies the reverse firewall to this message. The security and rerandomizability of the scheme are immediate from the security of the underlying MTP and the security of the firewall respectively. □

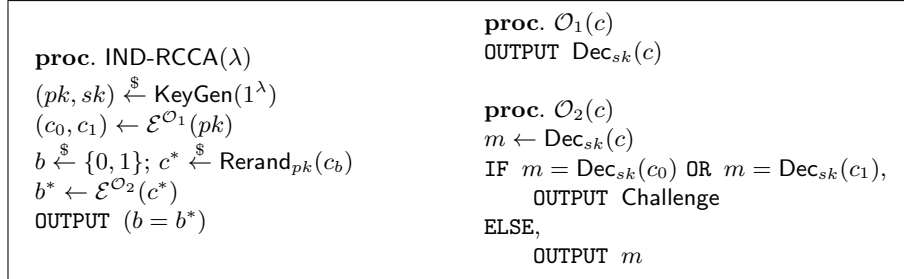
### 5.2 A one-round CCA-secure protocol

To extend this idea to stronger notions of security, we need the underlying encryption scheme to satisfy stronger notions of security. A natural candidate is CCA security. However, CCA-secure encryption schemes cannot be rerandomizable, so we need a slightly weaker notion. RCCA security, as defined by [CKN03], suffices, and rerandomizable RCCA-secure schemes do exist (see,

e.g., [Gro04, PR07]), though they are relatively inefficient. (They are not *strongly* rerandomizable; their rerandomization procedures do not work on invalid ciphertexts.) We present the RCCA security game in Figure 14. In addition, we need a rerandomized ciphertext to be indistinguishable from a valid encryption *even with access to a decryption oracle*. Figure 15 and the definition below makes this precise.



**Fig. 14:** The RCCA security game.



**Fig. 15:** The RCCA rerandomization game.

**Definition 15.** An encryption scheme is RCCA secure if no probabilistic polynomial-time adversary  $\mathcal{E}$  has non-negligible advantage in the game presented in Figure 14. It is RCCA rerandomizable if there exists an algorithm  $\text{Rerand}$  with access to the public key such that for any ciphertext  $c$  with  $\text{Dec}(c) \neq \perp$ ,  $\text{Dec}(\text{Rerand}(c)) = \text{Dec}(c)$  and no probabilistic polynomial-time adversary  $\mathcal{E}$  has non-negligible advantage in the game presented in Figure 15 when we require that  $\text{Dec}(c_i) \neq \perp$ . It is strongly RCCA-rerandomizable if the previous statement holds even if  $\text{Dec}(c_i) = \perp$ .

The below theorem is the CCA analogue of Theorem 6.

**Theorem 7.** Any (strongly) RCCA-rerandomizable, RCCA-secure encryption scheme implies a one-round CCA-secure singly keyed message-transmission protocol without forward security with a reverse firewall that (strongly) preserves security and (strongly) resists exfiltration.

*Proof.* Consider the protocol in which Alice, given as input a plaintext  $m$  and a session id  $\text{sid}$ , simply sends the message  $\text{Enc}_{pk}(\text{sid}, m)$ . On input  $c$ , Bob's return function computes  $(\text{sid}^\dagger, m) \leftarrow \text{Dec}_{sk}(c)$ . If  $\text{sid}^\dagger = \text{sid}$ , it outputs  $m$ . Otherwise, it outputs  $\perp$ . Alice's firewall simply applies the  $\text{Rerand}$  algorithm to Alice's message.

For any PPT adversary  $\mathcal{E}$  in the CCA-security game against this message-transmission protocol (without forward security), we construct  $\mathcal{E}'$  in the RCCA-security game against the underlying encryption scheme with  $\text{Adv}(\mathcal{E}') = \text{Adv}(\mathcal{E})$ . We assume without loss of generality that  $\mathcal{E}$  never makes a “useless” oracle call. I.e.,  $\mathcal{E}$  never calls `start-run` with an already used `sid`, never calls `get-nextA` with a never used `sid`, never calls `get-output` on the challenge `sid`, etc. The adversary  $\mathcal{E}'$  behaves as follows in response to the oracle calls of  $\mathcal{E}$ .

- When  $\mathcal{E}$  calls the `start-run` oracle with input  $(\text{sid}, m)$ ,  $\mathcal{E}'$  sets  $c_{\text{sid},A} \leftarrow \text{Enc}_{pk}(\text{sid}, m)$ .
- When  $\mathcal{E}$  calls `start-challenge`(`sid`,  $m_0, m_1$ ),  $\mathcal{E}'$  sets  $\text{sid}^* \leftarrow \text{sid}$ . It then returns  $(\text{sid}^*, m_0^*)$  and  $(\text{sid}^*, m_1^*)$  as its challenge plaintexts, receiving in response the challenge ciphertext  $c^*$ . It sets  $c_{\text{sid}^*} = c^*$ .
- When  $\mathcal{E}$  calls `get-nextA`(`sid`),  $\mathcal{E}'$  replies with  $c_{\text{sid},A}$ .
- When  $\mathcal{E}$  calls `get-nextB`(`sid`,  $c$ ),  $\mathcal{E}'$  sets  $c_{\text{sid},B} \leftarrow c$ .
- When  $\mathcal{E}$  calls `get-outputB`(`sid`),  $\mathcal{E}'$  calls  $\mathcal{O}_1(c_{\text{sid},B})$ . If the output is not of the form  $(\text{sid}, m)$ ,  $\mathcal{E}'$  responds with  $\perp$ . Otherwise, it responds with  $m$ .
- When  $\mathcal{E}$  calls `finalize`( $b^*$ ),  $\mathcal{E}'$  simply returns  $b^*$ .

Note that the view of  $\mathcal{E}$  is identical to its view in the CCA-security game against the message-transmission protocol. Furthermore,  $\mathcal{E}'$  wins the RCCA-security game if and only if  $\mathcal{E}$  wins its simulated game. The result follows.

It follows that the protocol is CCA secure. The proof that the firewall preserves security is nearly identical to the above proof.  $\square$

### 5.3 Achieving forward secrecy and efficiency?

Note that the protocols described above suffer from two problems: they do not have forward secrecy, and they are inefficient (i.e., they require public-key operations on the entire plaintext). Ideally, we would like a composition theorem in this setting to match Theorems 2 and 4. Such a theorem would solve both problems, allowing us to achieve an efficient CCA-secure and forward-secret message-transmission protocol in the singly keyed setting. We leave this as an open question.

As a potential alternative direction to achieving forward secrecy, we note that the protocol from Theorem 7 can be converted into a two-round CCA-secure and forward-secret protocol with a reverse firewall for Alice *but no reverse firewall for Bob*. In particular, in the first round of the protocol, Bob generates a fresh pair of keys  $(pk^\dagger, sk^\dagger)$  for an RCCA-secure RCCA-rerandomizable encryption scheme and sends Alice  $pk^\dagger$  together with a signature  $\tau$  of  $(\text{sid}, pk^\dagger)$ , where the signature is under Bob’s signature key. Alice checks the signature and, if it is valid, sends Bob an encryption of her plaintext under the secret key. This is essentially the CCA analogue of the protocol in Section 3.1. As in that case, Alice’s reverse firewall can simply check the signature and rerandomize the ciphertext. This protocol is in fact CCA secure, and Alice’s firewall does preserve this security. However, we do not know how to construct a reverse firewall for Bob in this setting. In analogy with Section 3.1, a possible method would be to find an RCCA-rerandomizable encryption scheme that is also key malleable. We know of no such scheme, but even if such a scheme were constructed, the signature scheme would have to be similarly malleable, while still achieving an appropriate notion of security. So, we also leave this as an open question.

## 6 Conclusion and open questions

We consider the problem of message-transmission protocols in the cryptographic reverse firewalls framework of [MS15]. We show that this problem has a rich structure, in analogy with the classical

setting, with a variety of solutions that require different setup assumptions, achieve different levels of security, and provide different levels of efficiency. Perhaps surprisingly, we show that it is possible to achieve concurrent, interactive CCA security in the presence of compromise against functionality-maintaining adversaries and CPA security against arbitrary adversaries. Many of our protocols (including those that provide the strongest notions of security) are very efficient and relatively simple, so that they can (and should) be implemented in practice.

Therefore, the most important work that we leave open is the implementation of our protocols. The most important theoretical question that we leave open is whether there is a non-trivial composition theorem in the singly keyed case, in analogy with Theorem 4. Note that Theorem 2 naturally extends to the singly keyed case, but we see no inherent reason why CCA-security should not be achievable from key agreement followed by symmetric-key encryption in this setting.

In addition, our work brings new attention to the question of rerandomizable RCCA-secure schemes. In particular, in Section 5, we show that such schemes give a one-round CCA-secure message-transmission protocol with a reverse firewall (without forward security). However, we do not know of such schemes that are “strongly rerandomizable” (as defined in Section 5). If such schemes existed, then we show that they would immediately imply one-round CCA-secure message transmission with a reverse firewall that *strongly* preserves security, a very powerful primitive. (Even just more efficient versions of current schemes would be quite interesting.)

## References

- [AMV15] Giuseppe Ateniese, Bernardo Magri, and Daniele Venturi. Subversion-resilient signature schemes. Cryptology ePrint Archive, Report 2015/517, 2015. <http://eprint.iacr.org/>.
- [BBG13] James Ball, Julian Borger, and Glenn Greenwald. Revealed: how US and UK spy agencies defeat internet privacy and security. *Guardian Weekly*, September 2013.
- [BH15] Mihir Bellare and Viet Tung Hoang. Resisting randomness subversion: Fast deterministic and hedged public-key encryption in the standard model. In Elisabeth Oswald and Marc Fischlin, editors, *Advances in Cryptology - EUROCRYPT 2015*, volume 9057 of *Lecture Notes in Computer Science*, pages 627–656. Springer Berlin Heidelberg, 2015.
- [BJK15] Mihir Bellare, Joseph Jaeger, and Daniel Kane. Mass-surveillance without the state: Strongly undetectable algorithm-substitution attacks. Cryptology ePrint Archive, Report 2015/808, 2015. <http://eprint.iacr.org/2015/808>.
- [BPR14a] Mihir Bellare, Kenneth G. Paterson, and Phillip Rogaway. Security of symmetric encryption against mass surveillance. In *CRYPTO 2014*, volume 8616 of *LNCS*, pages 1–19. Springer, 2014. Full version: [BPR14b].
- [BPR14b] Mihir Bellare, Kenneth G. Paterson, and Phillip Rogaway. Security of symmetric encryption against mass surveillance. Cryptology ePrint Archive, Report 2014/438, 2014. <http://eprint.iacr.org/>.
- [CKN03] Ran Canetti, Hugo Krawczyk, and Jesper Buus Nielsen. Relaxing chosen-ciphertext security. In Dan Boneh, editor, *CRYPTO*, volume 2729 of *LNCS*, pages 565–582. Springer, 2003.
- [CVE14a] Vulnerability summary for CVE-2014-1260 (‘Heartbleed’). <http://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2014-1260>, April 2014.
- [CVE14b] Vulnerability summary for CVE-2014-1266 (‘goto fail’). <http://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2014-1266>, February 2014.
- [CVE14c] Vulnerability summary for CVE-2014-6271 (‘Shellshock’). <http://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2014-6271>, September 2014.
- [DF14] Yevgeniy Dodis and Dario Fiore. Interactive encryption and message authentication. In Michel Abdalla and Roberto De Prisco, editors, *Security and Cryptography for Networks*, volume 8642 of *Lecture Notes in Computer Science*, pages 494–513. Springer International Publishing, 2014.
- [DFP15] Jean Paul Degabriele, Pooya Farshim, and Bertram Poettering. A more cautious approach to security against mass surveillance. In Gregor Leander, editor, *Fast Software Encryption*, volume 9054 of *Lecture Notes in Computer Science*, pages 579–598. Springer Berlin Heidelberg, 2015.
- [DGG<sup>+</sup>15] Yevgeniy Dodis, Chaya Ganesh, Alexander Golovnev, Ari Juels, and Thomas Ristenpart. A formal treatment of backdoored pseudorandom generators. In Elisabeth Oswald and Marc Fischlin, editors, *Advances in Cryptology - EUROCRYPT 2015*, volume 9056 of *Lecture Notes in Computer Science*, pages 101–126. Springer Berlin Heidelberg, 2015.

- [DH06] W. Diffie and M. Hellman. New directions in cryptography. *IEEE Trans. Inf. Theor.*, 22(6):644–654, September 2006.
- [DPSW06] Yvo Desmedt, Josef Pieprzyk, Ron Steinfeld, and Huaxiong Wang. A non-malleable group key exchange protocol robust against active insiders. In *Information Security*, volume 4176 of *Lecture Notes in Computer Science*, pages 459–475. Springer Berlin Heidelberg, 2006.
- [ElG85] Taher ElGamal. A public key cryptosystem and a signature scheme based on discrete logarithms. In *Proceedings of CRYPTO 84 on Advances in Cryptology*, pages 10–18, New York, NY, USA, 1985. Springer-Verlag New York, Inc.
- [Gre14] Glenn Greenwald. *No Place to Hide: Edward Snowden, the NSA, and the U.S. Surveillance State*. Metropolitan Books, May 2014.
- [Gro04] Jens Groth. Rerandomizable and replayable adaptive chosen ciphertext attack secure cryptosystems. In Moni Naor, editor, *TCC*, volume 2951 of *LNCS*, pages 152–170. Springer, 2004.
- [Kil07] Eike Kiltz. Chosen-ciphertext secure key-encapsulation based on Gap Hashed Diffie-Hellman. In Tatsuaki Okamoto and Xiaoyun Wang, editors, *Public Key Cryptography - PKC 2007, 10th International Conference on Practice and Theory in Public-Key Cryptography, Beijing, China, April 16-20, 2007, Proceedings*, volume 4450 of *Lecture Notes in Computer Science*, pages 282–297. Springer, 2007.
- [KS05] Jonathan Katz and Ji Sun Shin. Modeling insider attacks on group key-exchange protocols. In *Proceedings of the 12th ACM Conference on Computer and Communications Security, CCS '05*, pages 180–189, New York, NY, USA, 2005. ACM.
- [LHA<sup>+</sup>12] Arjen K. Lenstra, James P. Hughes, Maxime Augier, Joppe W. Bos, Thorsten Kleinjung, and Christophe Wachter. Public keys. In Reihaneh Safavi-Naini and Ran Canetti, editors, *Advances in Cryptology — CRYPTO 2012*, volume 7417 of *LNCS*, pages 626–642. Springer Berlin Heidelberg, 2012.
- [MRV99] Silvio Micali, Michael Rabin, and Salil Vadhan. Verifiable random functions. In *40th Annual Symposium on Foundations of Computer Science (New York, 1999)*, pages 120–130. IEEE Computer Soc., Los Alamitos, CA, 1999.
- [MS14] Ilya Mironov and Noah Stephens-Davidowitz. Cryptographic reverse firewalls. Cryptology ePrint Archive, Report 2014/758, 2014. <http://eprint.iacr.org/2014/758>.
- [MS15] Ilya Mironov and Noah Stephens-Davidowitz. Cryptographic reverse firewalls. In *EUROCRYPT 2015*, 2015. Available from [MS14].
- [Ped92] TorbenPryds Pedersen. Non-interactive and information-theoretic secure verifiable secret sharing. In Joan Feigenbaum, editor, *Advances in Cryptology — CRYPTO'91*, volume 576 of *Lecture Notes in Computer Science*, pages 129–140. Springer Berlin Heidelberg, 1992.
- [PLS13] Nicole Perlroth, Jeff Larson, and Scott Shane. N.S.A. able to foil basic safeguards of privacy on Web. *The New York Times*, September 2013.
- [PR07] Manoj Prabhakaran and Mike Rosulek. Rerandomizable RCCA encryption. In Alfred Menezes, editor, *CRYPTO 2007*, volume 4622 of *LNCS*, pages 517–534. Springer, 2007.
- [PW03] Josef Pieprzyk and Huaxiong Wang. Key control in multi-party key agreement protocols. Workshop on Coding, Cryptography and Combinatorics, 2003.
- [RTYZ15] Alexander Russell, Qiang Tang, Moti Yung, and Hong-Sheng Zhou. Cliptography: Clipping the power of kleptographic attacks. Cryptology ePrint Archive, Report 2015/695, 2015. <https://eprint.iacr.org/2015/695>.
- [SFKR15] Bruce Schneier, Matthew Fredrikson, Tadayoshi Kohno, and Thomas Ristenpart. Surreptitiously weakening cryptographic systems. Technical report, IACR Cryptology ePrint Archive, 2015: 97, 2015. <http://eprint.iacr.org/2015/97>.
- [Sim84] Gustavus Simmons. The prisoners’ problem and the subliminal channel. In David Chaum, editor, *Advances in Cryptology*, pages 51–67. Springer US, 1984.
- [Sup15] <https://www.us-cert.gov/ncas/alerts/TA15-051A>, February 2015.

## A The symmetric-key setting

As a warmup, we first consider the setting in which Alice and Bob share a private key. We observe that an elegant one-round protocol due to Bellare, Paterson, and Rogaway provides a solution that does not even need a reverse firewall [BPR14a]. We will also use this scheme in the sequel to build protocols that do not rely on shared private keys. We first define nonce-based encryption.

**Definition 16 (Nonce-based encryption).** A nonce-based symmetric-key encryption scheme is a pair of deterministic algorithms (Enc, Dec). Enc takes as input a key from a key space  $\mathcal{K}$ , a nonce

from a nonce space  $\mathcal{N}$ , and a plaintext from a plaintext space  $\mathcal{M}$  and outputs a ciphertext from a ciphertext space  $\mathcal{C}$ .  $\text{Dec}$  takes as input a key, a nonce, and a ciphertext and returns a plaintext or the special symbol  $\perp$ . The scheme is correct if for any key  $sk$ , nonce  $r$ , and plaintext  $m$ ,  $\text{Dec}(r, \text{Enc}_{sk}(r, m)) = m$ .

Such a scheme is CPA secure if no probabilistic polynomial-time adversary can distinguish between  $\text{Enc}_{sk}(r^*, m_0)$  and  $\text{Enc}_{sk}(r^*, m_1)$  with non-negligible advantage where  $r^*$ ,  $m_0$ , and  $m_1$  are adversarially chosen when given access to an encryption oracle that outputs  $\perp$  unless given a unique nonce  $r$ . It is CCA-secure if no probabilistic polynomial-time has non-negligible advantage when also given access to a decryption oracle that outputs  $\perp$  if  $r = r^*$ .

Such a scheme scheme  $(\text{Enc}, \text{Dec})$  has unique ciphertexts if for any key  $sk$ , message  $m$ , and nonce  $r$ , there is exactly one ciphertext  $c$  such that  $\text{Dec}(r, c) = m$ .

**Theorem 8.** *Let  $(\text{Enc}, \text{Dec})$  be a nonce-based symmetric-key encryption scheme. Then, if the encryption scheme is CPA-secure (resp. CCA-secure) the one-round protocol in which Alice sends  $\text{Enc}_{sk}(\text{sid}, m)$  and Bob returns  $\text{Dec}_{sk}(\text{sid}, m)$  is a CPA-secure (resp. CCA-secure) one-round privately keyed message-transmission protocol without forward secrecy. If the encryption scheme has unique ciphertexts, then the “trivial” reverse firewall that simply passes Alice’s messages to Bob unchanged preserves security and is exfiltration resistant against Bob.*

See, e.g., [BPR14a] for formal analysis and the construction of such a scheme. The key thing to note from our perspective is that, as Bellare et al. observe, the fact that the encryption scheme has unique ciphertexts implies that any tampered version of Alice that maintains functionality necessarily behaves identically to honest Alice. The next theorem shows that we essentially cannot do any better for a one-round protocol without using public-key primitives. The proof is in appendix F.

**Theorem 9.** *There is a black-box reduction from semantically secure public-key encryption to CPA-secure symmetric-key encryption with at least four possible plaintexts and a reverse firewall that strongly preserves CPA security.*

Of course, the primary drawbacks of this approach are that it requires Alice and Bob to share a secret key and that it does not offer forward secrecy.

## B Rerandomizable and key-malleable public-key encryption

**Definition 17 (Public-key encryption).** *A public-key encryption scheme is a triple of efficient algorithms  $(\text{KeyGen}, \text{Enc}, \text{Dec})$ .  $\text{KeyGen}$  takes as input  $1^\lambda$  where  $\lambda$  is the security parameter and outputs a public-key/private-key pair,  $(pk, sk)$ .  $\text{Enc}$  takes as input the public key and a plaintext  $m$  from some plaintext space  $\mathcal{M}$  and outputs a ciphertext  $c$  from some ciphertext space  $\mathcal{C}$ .  $\text{Dec}$  takes as input a ciphertext and the private key and outputs a plaintext or the special symbol  $\perp$ . We sometimes omit the keys from the input to  $\text{Enc}$  and  $\text{Dec}$  and the security parameter input to  $\text{KeyGen}$ . The scheme is correct if  $\text{Dec}(\text{Enc}(m)) = m$  for all  $m \in \mathcal{M}$ . The scheme is semantically secure if for any adversarially chosen pair of plaintexts  $(m_0, m_1)$ ,  $\text{Enc}(m_0)$  is computationally indistinguishable from  $\text{Enc}(m_1)$ .*

**Definition 18 (Rerandomizable encryption).** *A semantically secure public-key encryption scheme is rerandomizable if there is an efficient algorithm  $\text{Rerand}$  (with access to the public key) such that for any ciphertext  $c$  such that  $\text{Dec}(c) \neq \perp$ , we have  $\text{Rerand}(\text{Dec}(c)) = \text{Dec}(c)$ , and the pair  $(c, \text{Rerand}(c))$  is computationally indistinguishable from  $(c, \text{Rerand}(\text{Enc}(0)))$ . We say that it is strongly rerandomizable if the previous property holds even when  $\text{Dec}(c) = \perp$ .*

**Definition 19 (Key malleability).** A public-key encryption scheme is key malleable if (1) the output of `KeyGen` is distributed uniformly over the space of valid keys; (2) for each public key  $pk$  there is a unique associated private key  $sk$ ; and (3) there is a pair of efficient algorithms `KeyMaul` and `CKeyMaul` that behave as follows. `KeyMaul` is a randomized algorithm that takes as input a public key  $pk$  and returns a new public key  $pk'$  whose distribution is uniformly random over the public key space and independent of  $pk$ . Let  $(sk, pk)$  be a private key/public key pair. Let  $(sk', pk')$  be the unique pair associated with randomness  $r$  such that  $pk' = \text{KeyMaul}(pk; r)$ . Then, `CKeyMaul` takes as input a ciphertext  $c$  and randomness  $r$  and returns  $c'$  such that  $\text{Dec}_{sk'}(c) = \text{Dec}_{sk}(c')$ . We suppress the input  $r$  when it is understood. Furthermore, we require that `KeyMaul` outputs a uniformly random key  $pk'$  if called on input that is not in the public-key space.

**Example.** It is well-known that ElGamal encryption [ELG85] is both key malleable and strongly rerandomizable. In particular, given an ElGamal public-key  $(g, h)$  over a group of order  $p$ , a ciphertext  $(u, v)$  can be rerandomized by applying the operation  $(u, v) \rightarrow (g^r u, h^r v)$  where  $r \xleftarrow{\$} \mathbb{Z}_p^*$  is chosen uniformly at random. The public key can be mauled by applying the operation  $(g, h) \rightarrow (g^\alpha, h^\beta)$  where  $(\alpha, \beta) \xleftarrow{\$} (\mathbb{Z}_p^*)^2$  are chosen uniformly and independently at random. Finally, a ciphertext  $(u, v)$  under key  $(g^\alpha, h^\beta)$  can be converted into a ciphertext under  $(g, h)$  by applying the operation  $(u, v) \rightarrow (u^{\beta/\alpha}, v)$ .

## C Proof of Theorem 2

*Proof.* The security of the underlying protocol (i.e., without firewalls) follows by a folklore composition theorem. We assume that the last message of the key-agreement protocol is sent by Bob. (This is without loss of generality, as we can always add an empty message from Bob at the end of the protocol.) We prove Item 2. The proof of Item 3 is a simplified version of the proof of Theorem 4.

Let  $\tilde{B}$  be some functionality-maintaining tampered implementation of Bob. Let  $\tilde{B}^*$  be the tampered implementation of Bob in the key-agreement protocol obtained by simply “truncating”  $\tilde{B}$ . (This is merely a syntactic change.) We assume without loss of generality that the adversary makes no “trivial” calls whose output can be trivially predicted from its previous oracle calls. E.g., it makes no calls to `get-output` (which always either returns  $\perp$  or the plaintext that the adversary provided for the corresponding `sid` in the CPA-security game), no `get-next` calls with modified messages, etc. Consider the following sequence of games.

- **Game 1** is the CPA-security game against the message-transfer protocol with Bob replaced by  $\mathcal{W}_{\tilde{B}}^* \circ \tilde{B}$ .
- **Game 2** is **Game 1** in which the final message of the challenge run of the protocol is replaced by an encryption of the challenge plaintext  $m_b$  under a uniformly random key  $sk$ .
- **Game 3** is **Game 2** in which the final message of the challenge run is replaced by an encryption of 0 under a uniformly random key.

Note that no adversary can have non-zero advantage in the last game, as none of the messages depend on the challenge bit.

**Claim 2.1.** *If the encryption scheme has unique ciphertexts, the key-agreement protocol is secure against passive adversaries, and  $\mathcal{W}_{\tilde{B}}$  preserves this security, then for any PPT adversary  $\mathcal{E}$ ,  $|\text{Adv}^{\text{Game 1}}(\mathcal{E}) - \text{Adv}^{\text{Game 2}}(\mathcal{E})|$  is negligible.*



*Proof.* We construct a passive adversary  $\mathcal{E}'$  in the security game against the key-agreement protocol with Bob replaced by  $\mathcal{W}_B \circ \tilde{\mathbf{B}}^*$  as follows.  $\mathcal{E}'$  receives as input the public parameters  $\pi$  and passes them to  $\mathcal{E}$ .  $\mathcal{E}'$  then selects  $b^\dagger \xleftarrow{\$} \{0, 1\}$  uniformly at random and sets  $S_A^\dagger \leftarrow \emptyset$ . It responds to the (non-trivial) oracle calls of  $\mathcal{E}$  as follows.

- When  $\mathcal{E}$  calls  $\text{start-run}(\text{sid}, m)$ ,  $\mathcal{E}'$  adds  $(\text{sid}, m)$  to  $S_A^\dagger$  and calls its own oracle  $\text{start-run}(\text{sid})$ .
- When  $\mathcal{E}$  calls  $\text{start-challenge}(\text{sid}^*, m_0, m_1)$ ,  $\mathcal{E}'$  adds  $(\text{sid}^*, m_{b^\dagger})$  to  $S_A^\dagger$  and calls its own oracle  $\text{start-challenge}(\text{sid}^*)$ .
- When  $\mathcal{E}$  calls  $\text{get-next}_A(\text{sid}, M)$ ,  $\mathcal{E}'$  calls its own oracle  $\text{get-next}_A(\text{sid}, M)$ . If this is not the last message of the protocol,  $\mathcal{E}'$  then simply passes the resulting message to  $\mathcal{E}$ . If it is the last message, it  $\mathcal{E}'$  calls its own oracle  $\text{get-output}_A(\text{sid})$ , receiving as output some key  $sk$ . It responds with  $\text{Enc}_{sk}(\text{sid}, m)$  where  $m$  is the unique plaintext such that  $(\text{sid}, m) \in S_A^\dagger$ . (Note that since the oracle calls of  $\mathcal{E}$  are non-trivial, it must have made a unique call to either  $\text{start-run}$  or  $\text{start-challenge}$  with this  $\text{sid}$ .)
- When  $\mathcal{E}$  calls  $\text{get-next}_B(\text{sid}, M)$ , if this is not the last message of the protocol,  $\mathcal{E}'$  calls its own oracle  $\text{get-next}_B(\text{sid}, M)$  and passes the resulting message to  $\mathcal{E}$ . Otherwise,  $\mathcal{E}'$  does nothing.
- When  $\mathcal{E}$  calls  $\text{get-secrets}$ ,  $\mathcal{E}'$  calls its own oracle  $\text{get-secrets}$  and passes the result to  $\mathcal{E}$ .
- When  $\mathcal{E}$  calls  $\text{finalize}(b^*)$ ,  $\mathcal{E}'$  returns 1 if  $b^\dagger = b^*$  and 0 otherwise.

Suppose the challenge bit  $b$  in the key-agreement security game is 0, so that the challenge key  $sk^*$  corresponding to the challenge session id  $\text{sid}^*$  was selected uniformly at random. Then, clearly the view of  $\mathcal{E}$  is identical to its view in **Game 2**, and  $\mathcal{E}'$  correctly outputs 0 if and only if  $\mathcal{E}$  “loses” its simulated game by returning  $b^\dagger = b^*$ . If, on the other hand, the challenge bit  $b$  is 1, then the view of  $\mathcal{E}$  is identical to its view in **Game 1**, and  $\mathcal{E}'$  correctly outputs 1 if and only if  $\mathcal{E}$  “wins” its simulated game. The result follows. (2.1) ■

**Claim 2.2.** *If the encryption scheme is CPA-secure, then for any PPT adversary  $\mathcal{E}$ ,  $|\text{Adv}^{\text{Game } 2}(\mathcal{E}) - \text{Adv}^{\text{Game } 3}(\mathcal{E})|$  is negligible.*

*Proof.* We construct an adversary  $\mathcal{E}'$  in the CPA-security game against the encryption scheme as follows.  $\mathcal{E}'$  first runs the  $\text{setup}$  procedure of the key-agreement protocol, receiving as output  $\sigma_A$ ,  $\sigma_B$ , and  $\pi$ .  $\mathcal{E}'$  will simulate many runs of the key-agreement protocol with Bob replaced by  $\mathcal{W}_B \circ \tilde{\mathbf{B}}^*$  and input  $(\sigma_A, \sigma_B, \pi)$ . For convenience, we give  $\mathcal{E}'$  an “oracle interface” to these simulated runs with “oracle” calls  $\text{start-run}$ ,  $\text{get-next}_A$ ,  $\text{get-next}_B$ , and  $\text{get-output}_A$  as in the key-agreement security game (Figure 3).  $\mathcal{E}'$  sends  $\pi$  to  $\mathcal{E}$ , selects  $b^\dagger \xleftarrow{\$} \{0, 1\}$  uniformly at random, and sets  $S_A^\dagger \leftarrow \emptyset$ . It then responds to the oracle queries of  $\mathcal{E}$  as follows.

- When  $\mathcal{E}$  calls  $\text{start-run}(\text{sid}, m)$ ,  $\mathcal{E}'$  adds  $(\text{sid}, m)$  to  $S_A^\dagger$ . It calls its own “oracle”  $\text{start-run}(\text{sid})$ .
- When  $\mathcal{E}$  calls  $\text{start-challenge}(\text{sid}^*, m_0, m_1)$ ,  $\mathcal{E}'$ , if  $b^\dagger = 0$ , it sends the challenge  $(\text{sid}^*, m_0, 0)$  to its challenger. Otherwise, it sends the challenge  $(\text{sid}^*, 0, m_1)$ . It stores the resulting challenge ciphertext  $c^*$  and calls its “oracle”  $\text{start-run}(\text{sid})$ .
- When  $\mathcal{E}$  calls  $\text{get-next}_A(\text{sid}, M)$ ,  $\mathcal{E}'$  calls its own “oracle”  $\text{get-next}_A(\text{sid}, M)$ . If this call does not correspond to the last message of the relevant run of the message-transfer protocol,  $\mathcal{E}'$  simply passes the response to  $\mathcal{E}$ . If this is the last message and  $\text{sid} = \text{sid}^*$ ,  $\mathcal{E}'$  sends  $c^*$  to  $\mathcal{E}$ . Otherwise,  $\mathcal{E}'$  sets  $sk \leftarrow \text{get-output}_A(\text{sid})$ , finds the unique  $m$  such that  $(\text{sid}, m) \in S_A^\dagger$ , and sends  $\text{Enc}_{sk}(\text{sid}, m)$  to  $\mathcal{E}$ .
- When  $\mathcal{E}$  calls  $\text{get-next}_B(\text{sid}, M)$ , if this is not the last message of the protocol,  $\mathcal{E}'$  calls its own “oracle”  $\text{get-next}_B(\text{sid}, M)$  and passes the response to  $\mathcal{E}$ . If this is the last message,  $\mathcal{E}'$  does nothing.

- When  $\mathcal{E}$  calls `get-secrets`,  $\mathcal{E}'$  responds with  $(\sigma_A, \sigma_B)$ .
- When  $\mathcal{E}$  calls `finalize( $b^*$ )`,  $\mathcal{E}'$  returns  $b^*$ .

Let  $b$  be the challenge bit of the CPA-security game. If  $b = b^\dagger$ , then the challenge ciphertext is an encryption of  $m_b$ , as in **Game 2** and  $\mathcal{E}'$  correctly outputs  $b$  if and only if  $\mathcal{E}$  “wins” and  $b^* = b^\dagger$ . Otherwise, the challenge ciphertext is an encryption of 0 as in **Game 3** and  $\mathcal{E}'$  correctly outputs  $b$  if and only if  $\mathcal{E}$  “loses.” The result follows. (2.2) ■

The result follows. □

## D A suitable malleable commitment scheme

We briefly describe a simple commitment scheme that is statistically hiding, computationally binding, and malleable in the way that we need. It is a basic variant of the Pedersen commitment [Ped92]. We first provide the relevant definitions. For our application, we require that the committed plaintext is a group element and that the commitment can be mauled in such a way that a commitment to group element  $B$  can be converted into a commitment to group element  $B^\alpha$ . As such, we define a malleable commitment scheme in this specific setting.

**Definition 20 (Commitment scheme).** A commitment scheme is a collection of four efficient algorithms (`KeyGen`, `Com`, `Open`, `Ver`). `KeyGen` takes as input  $1^\lambda$ , where  $\lambda$  is the security parameter, and outputs public parameters  $\rho$ . `Com` takes as input the public parameters  $\rho$ , a plaintext  $m$ , and randomness  $r$ , and outputs a commitment  $C$ . `Open` takes as input the public parameters  $\rho$ , a commitment  $C$ , and randomness  $r$  and outputs an opening  $x$ . Finally, `Ver` takes as input a commitment  $C$  and opening  $x$  and outputs either a plaintext  $m$  or the special symbol  $\perp$ . For simplicity, we often omit references to the public parameters and the randomness.

The scheme is correct if a commitment opens to the committed message, i.e.,  $\text{Ver}(C, x) = m$  whenever  $C = \text{Com}(m)$  and  $x = \text{Open}(C)$ . The scheme is perfectly hiding if for any two plaintexts  $m_1, m_2$ , the distribution of  $\text{Com}(m_1)$  is identical to that of  $\text{Com}(m_2)$ . The scheme is computationally binding if no efficient adversary can produce a commitment  $C$  and two openings  $x, x'$  that verify to different plaintexts (and not  $\perp$ ).

**Definition 21 (Additional commitment properties).** A commitment scheme is tight if for each commitment  $C$  and plaintext  $m$ , there is a unique opening  $x$  such that  $\text{Ver}(C, x) = m$ .

A commitment scheme is rerandomizable if there exists a pair of efficient algorithms (`Rerand`, `OpenRerand`) such that

1. for any commitment  $C$  and uniformly random  $r$ ,  $\text{Rerand}(C, r)$  is uniformly random over the commitment space; and
2. for any commitment  $C$  and opening  $x$ ,  $\text{Ver}(C, x) = \text{Ver}(C', x')$ , where  $C' = \text{Rerand}(C, r)$  and  $x' = \text{OpenRerand}(x', r)$ .

Similarly, a commitment scheme is malleable if the plaintext space is a group of order  $p$  (written multiplicatively) and there exists a pair of efficient algorithms (`Maul`, `OpenMaul`) such that for any commitment  $C$  and opening  $x$ ,  $\text{Ver}(C, x)^\alpha = \text{Ver}(C', x')$ , where  $C' = \text{Maul}(C, \alpha)$  and  $x' = \text{OpenMaul}(x', \alpha)$  for  $\alpha \in \mathbb{Z}_p$ . In the main body, we omit explicit reference to the functions `OpenRerand` and `OpenMaul`.

We now describe a commitment scheme that suffices for our purposes.

- **KeyGen** takes as input a group  $G$  of order  $p$  and returns two uniformly random non-identity group elements  $g, h$ .
- **Com** takes as input a group element  $m \in G$  and randomness  $(r, s) \in \mathbb{Z}_p^2$  and returns  $C := (g^r h^s, h^s m)$ .
- **Open** simply outputs the randomness  $(r, s)$ .
- **Ver** takes as input  $C = (c_1, c_2)$  and  $(r, s)$ . It first checks that  $c_1 = g^r h^s$ . If so, it returns  $h^{-s} c_2$ .
- **Rerand** takes as input  $C = (c_1, c_2)$  and randomness  $(r', s')$  and returns  $C' := (g^{r'} h^{s'} c_1, h^{s'} c_2)$ .
- **OpenRerand** takes as input an opening  $(r, s)$  and randomness  $(r', s')$  and returns  $(r + r', s + s')$ .
- **Maul** takes as input  $C = (c_1, c_2)$  and  $\alpha \in \mathbb{Z}_p$  and returns  $C' := (c_1^\alpha, c_2^\alpha)$ .
- **OpenMaul** takes as input  $(r, s)$  and  $\alpha$  and returns  $(\alpha r, \alpha s)$ .

The following proposition is immediate from inspection and the security of the standard Pedersen commitment scheme [Ped92].

**Proposition 1.** *The above commitment scheme is correct, statistically hiding, tight, rerandomizable, and malleable. If the discrete log is hard over  $G$ , it is also computationally binding.*

## E Proof of Theorem 4

*Proof.* Item 1 (the security of the underlying protocol without reverse firewalls) follows by a folklore composition theorem. We assume that the last message of the key-agreement protocol is sent by Bob. (This is without loss of generality, as we can always add an empty message from Bob at the end of the protocol.)

Let  $\tilde{A}$  be some functionality-maintaining tampered implementation of Alice in the protocol from Figure 7. Let  $\tilde{A}^*$  be the “truncation” of  $\tilde{A}$  to the key-agreement protocol. Note that  $\tilde{A}^*$  produces valid transcripts (though it may not preserve functionality). Let  $q$  be some polynomial bound on the number of oracle calls made by the adversary in the CCA-security game. We assume without loss of generality that the adversary makes no “trivial” calls whose output can be trivially predicted from its previous oracle calls. E.g., it makes no calls to  $\text{get-output}_B(\text{sid})$  where  $\text{sid}$  does not correspond to a completed run of the protocol, no  $\text{get-next}_A(\text{sid}, \cdot)$  without first calling  $\text{start-run}(\text{sid}, m)$  or  $\text{start-challenge}(\text{sid}, m_0, m_1)$ , etc. Consider the following sequence of games.

- **Game 1** is the CCA-security game against the message-transmission protocol with Alice replaced by  $\mathcal{W}_A^* \circ \tilde{A}$ .
- **Game 2** is **Game 1** in which Alice never sends the final message in a run of the protocol unless Bob’s key  $sk_B$  is well-defined. (I.e., a call to the  $\text{return}_B(\text{sid})$  procedure of the underlying key-agreement protocol for the relevant  $\text{sid}$  does not return  $\perp$ .)
- For  $i = 1, \dots, q$ , **Game  $i + 2$**  is **Game  $i + 1$**  in which the final message of the  $i$ th run of the protocol is replaced by an encryption of the relevant plaintext  $m$  under a uniformly random key.
- **Game  $q + 3$**  is **Game  $q + 2$**  in which the final message of the challenge run is replaced by an encryption of 0 under a uniformly random key.
- **Game  $q + 4$**  is **Game  $q + 3$**  in which the final message is removed from each run of the protocol and the oracle  $\text{get-output}_B$  is removed.
- For  $i = 1, \dots, q$ , **Game  $q + 4 + i$**  is **Game  $q + 3 + i$**  in which  $\mathcal{W}_A^* \circ \tilde{A}$  is replaced by  $\mathcal{W}_A^* \circ A$  (the honest implementation of Alice composed with the firewall) in the  $i$ th run of the protocol.

Note that no adversary can have any advantage in the last game because none of the responses to any of the adversary’s queries depend on the challenge bit.

**Claim 4.1.** *If the encryption scheme has unique ciphertexts, the key-agreement protocol is authenticated for Bob, and  $\mathcal{W}_A$  fails detectably and preserves authenticity, then for any PPT adversary  $\mathcal{E}$ ,  $|\text{Adv}^{\text{Game } 1}(\mathcal{E}) - \text{Adv}^{\text{Game } 2}(\mathcal{E})|$  is negligible.*

*Proof.* We construct an adversary  $\mathcal{E}'$  in the authentication game against the key-agreement protocol with Alice replaced by  $\mathcal{W}_A \circ \tilde{A}^*$  as follows.  $\mathcal{E}'$  receives as input the public parameters  $\pi$  and passes them to  $\mathcal{E}$ .  $\mathcal{E}'$  then selects  $b^\dagger \xleftarrow{\$} \{0, 1\}$  uniformly at random and sets  $S_A^\dagger, S_B^\dagger, \text{keys} \leftarrow \emptyset$ .  $\mathcal{E}'$  then responds to oracle calls as follows.

- When  $\mathcal{E}$  calls  $\text{start-run}(\text{sid}, m)$ ,  $\mathcal{E}'$  adds  $(\text{sid}, m)$  to  $S_A^\dagger$  and calls its own oracle  $\text{start-run}(\text{sid}, m)$ .
- When  $\mathcal{E}$  calls  $\text{start-challenge}(\text{sid}, m_0, m_1)$ ,  $\mathcal{E}'$  adds  $(\text{sid}, m_{b^\dagger})$  to  $S_A^\dagger$  and calls its own oracle  $\text{start-run}(\text{sid}, m_{b^\dagger})$ .
- When  $\mathcal{E}$  calls  $\text{get-next}_A(\text{sid}, M)$ ,  $\mathcal{E}'$  calls its own oracle  $\text{get-next}_A(\text{sid}, M)$ . If this is not the last message of this run of the protocol,  $\mathcal{E}'$  then simply passes the resulting message to  $\mathcal{E}$ . If it is the last message, it checks if the transcript of the underlying key-agreement protocol is valid for  $\mathcal{W}_A \circ A$  (using the efficient algorithm guaranteed by detectable failure). If it is invalid,  $\mathcal{E}'$  responds to  $\mathcal{E}$  with the special symbol  $\perp$ . If it is valid, let  $m$  be the unique plaintext and index such that  $(\text{sid}, m) \in S_A^\dagger$  and let  $sk \leftarrow \text{get-output}_B(\text{sid})$ . If  $sk = \perp$ , then  $\mathcal{E}'$  returns the transcript of this run of the key-agreement protocol (and wins the authentication game). Otherwise, it adds  $(\text{sid}, sk)$  to  $\text{keys}$  and responds with  $\text{Enc}_{sk}(\text{sid}, m)$ .
- When  $\mathcal{E}$  calls  $\text{get-next}_B(\text{sid}, M)$ , if this is not the last message of the protocol,  $\mathcal{E}'$  calls its own oracle  $\text{get-next}_B(\text{sid}, M)$  and passes the resulting message to  $\mathcal{E}$ . If this is the last message of the protocol, it adds  $(\text{sid}, M)$  to  $S_B^\dagger$  and sends nothing to  $\mathcal{E}$ .
- When  $\mathcal{E}$  calls  $\text{get-output}_B(\text{sid})$ ,  $\mathcal{E}'$  finds the unique  $M$  and  $sk$  such that  $(\text{sid}, M) \in S_B^\dagger$  and  $(\text{sid}, sk) \in \text{keys}$ . It computes  $\text{Dec}_{sk}(\text{sid}, M)$  and responds with the result.
- When  $\mathcal{E}$  calls  $\text{get-secrets}$ ,  $\mathcal{E}'$  calls its own oracle  $\text{get-secrets}$  and passes the result to  $\mathcal{E}$ .
- When  $\mathcal{E}$  calls  $\text{finalize}(b^*)$ ,  $\mathcal{E}'$  simply terminates.

Note that the view of  $\mathcal{E}$  is identical to its view in both **Game 1** and **Game 2** unless at some point it constructs a valid transcript such that  $sk_B$  is not well-defined. If it does construct such a transcript, then  $\mathcal{E}'$  wins the authentication game. The result follows. (4.1) ■

**Claim 4.2.** *If the key-agreement protocol is secure against active adversaries for Alice, the encryption scheme has unique ciphertexts, and  $\mathcal{W}_A$  preserves Alice's security and fails detectably, then for any PPT adversary  $\mathcal{E}$ ,  $|\text{Adv}^{\text{Game } i+1}(\mathcal{E}) - \text{Adv}^{\text{Game } i+2}(\mathcal{E})|$  is negligible.*

*Proof.* We construct an adversary  $\mathcal{E}'$  in the security game against the key-agreement protocol with Alice replaced by  $\mathcal{W}_A \circ \tilde{A}^*$  as follows.  $\mathcal{E}'$  receives as input the public parameters  $\pi$  and passes them to  $\mathcal{E}$ .  $\mathcal{E}'$  then selects  $b^\dagger \xleftarrow{\$} \{0, 1\}$  uniformly at random and sets  $S_A^\dagger, S_B^\dagger, \text{keys} \leftarrow \emptyset$  and  $j \leftarrow 0$ . It responds to the oracle calls of  $\mathcal{E}$  as follows.

- When  $\mathcal{E}$  calls  $\text{start-run}(\text{sid}, m)$ ,  $\mathcal{E}'$  adds  $(\text{sid}, j, m)$  to  $S_A^\dagger$ . If  $j = i$ , it calls its own oracle  $\text{start-challenge}(\text{sid}^*, m)$ ; otherwise it calls  $\text{start-run}(\text{sid}, m)$ . Finally, it increments  $j$ .
- When  $\mathcal{E}$  calls  $\text{start-challenge}(\text{sid}, m_0, m_1)$ ,  $\mathcal{E}'$  adds  $(\text{sid}, m_{b^\dagger})$  to  $S_A^\dagger$  and  $(\text{sid}, j)$  to  $S_B^\dagger$ . If  $j = i$ , it calls its own oracle  $\text{start-challenge}(\text{sid}^*, m_{b^\dagger})$ ; otherwise it calls  $\text{start-run}(\text{sid}, m_{b^\dagger})$ . Finally, it increments  $j$ .
- When  $\mathcal{E}$  calls  $\text{get-next}_A(\text{sid}, M)$ ,  $\mathcal{E}'$  calls its own oracle  $\text{get-next}_A(\text{sid}, M)$ . if this is not the last message of this run of the protocol,  $\mathcal{E}'$  then simply passes the resulting message to  $\mathcal{E}$ . If it is the last message, it checks if the transcript of the underlying key-agreement protocol is valid

for  $\mathcal{W}_A \circ A$  (using the efficient algorithm guaranteed by detectable failure). If it is invalid,  $\mathcal{E}'$  responds to  $\mathcal{E}$  with the special symbol  $\perp$ . If it is valid, let  $m, k$  be the unique plaintext and index such that  $(\text{sid}, k, m) \in S_A^\dagger$ . If  $k < i$ ,  $\mathcal{E}'$  selects  $sk \xleftarrow{\$} \mathcal{K}$ . If  $k \geq i$ ,  $\mathcal{E}'$  sets  $sk \leftarrow \text{get-output}_B(\text{sid})$ . Finally, it responds to  $\mathcal{E}$  with the message  $\text{Enc}_{sk}(\text{sid}, m)$ .

- When  $\mathcal{E}$  calls  $\text{get-next}_B(\text{sid}, M)$ ,  $\mathcal{E}'$  calls its own oracle  $\text{get-next}_B(\text{sid}, M)$  and passes the resulting message to  $\mathcal{E}$ . If this is the last message of the protocol, it also adds  $(\text{sid}, M)$  to  $S_B^\dagger$ .
- When  $\mathcal{E}$  calls  $\text{get-output}_B(\text{sid})$ ,  $\mathcal{E}'$  finds the unique  $M$  and  $sk$  such that  $(\text{sid}, M) \in S_B^\dagger$  and  $(\text{sid}, sk) \in \text{keys}$ . It computes  $\text{Dec}_{sk}(\text{sid}, M)$  and responds with the result.
- When  $\mathcal{E}$  calls  $\text{get-secrets}$ ,  $\mathcal{E}'$  calls its own oracle  $\text{get-secrets}$  and passes the result to  $\mathcal{E}$ .
- When  $\mathcal{E}$  calls  $\text{finalize}(b^*)$ ,  $\mathcal{E}'$  returns 1 if  $b^\dagger = b^*$  and 0 otherwise.

Suppose the challenge bit  $b$  in the key-agreement security game is 0 so that the challenge key  $sk^*$  corresponding to the  $i$ th run of the protocol was selected uniformly at random or is the special symbol  $\perp$ . Then, clearly the view of  $\mathcal{E}$  is identical to its view in **Game**  $i + 2$ , and  $\mathcal{E}'$  correctly outputs 0 if and only if  $\mathcal{E}$  “loses” its simulated game by returning  $b^\dagger = b^*$ . If, on the other hand, the challenge bit  $b$  is 1, then the view of  $\mathcal{E}$  is identical to its view in **Game**  $i + 1$ , and  $\mathcal{E}'$  correctly outputs 1 if and only if  $\mathcal{E}$  “wins” its simulated game. The result follows.

(4.2) ■

**Claim 4.3.** *If the encryption scheme is CCA-secure, then for any PPT adversary  $\mathcal{E}$ ,  $|\text{Adv}^{\text{Game } q+2}(\mathcal{E}) - \text{Adv}^{\text{Game } q+3}(\mathcal{E})|$  is negligible.*

*Proof.* We construct an adversary  $\mathcal{E}'$  in the CCA-security game against the encryption scheme as follows.  $\mathcal{E}'$  first runs the **setup** procedure of the key-agreement protocol, receiving as output  $\sigma_A, \sigma_B$ , and  $\pi$ .  $\mathcal{E}'$  will simulate many runs of the key-agreement protocol with Alice replaced by  $\mathcal{W}_A \circ A^*$  and input  $(\sigma_A, \sigma_B, \pi)$ . For convenience, we give  $\mathcal{E}'$  an “oracle interface” to these simulated runs with “oracle” calls **start-run**, **get-next<sub>A</sub>**, and **get-next<sub>B</sub>** as in the key-agreement security game (Figure 3).  $\mathcal{E}'$  sends  $\pi$  to  $\mathcal{E}$ , selects  $b^\dagger \xleftarrow{\$} \{0, 1\}$  uniformly at random, and sets  $S_A, S_B, \text{keys} \leftarrow \emptyset$ . It then responds to the oracle queries of  $\mathcal{E}$  as follows.

- When  $\mathcal{E}$  calls **start-run**( $\text{sid}, m$ ),  $\mathcal{E}'$  adds  $(\text{sid}, m)$  to  $S_A$ . It calls its own “oracle” **start-run**( $\text{sid}, m$ ).
- When  $\mathcal{E}$  calls **start-challenge**( $\text{sid}^*, m_0, m_1$ ), if  $b^\dagger = 0$ , it  $\mathcal{E}'$  sends the challenge  $(\text{sid}^*, m_0, 0)$  to its challenger. Otherwise, it sends the challenge  $(\text{sid}^*, 0, m_1)$ . It stores the resulting challenge ciphertext  $c^*$ . and calls its “oracle” **start-run**( $\text{sid}, m_{b^\dagger}$ ).
- When  $\mathcal{E}$  calls **get-next<sub>A</sub>**( $\text{sid}, M$ ), if this call does not correspond to the last message of the relevant run of the message-transmission protocol,  $\mathcal{E}'$  calls its own “oracle” **get-next<sub>A</sub>**( $\text{sid}, M$ ) and passes the response to  $\mathcal{E}$ . If the transcript of the underlying key-agreement protocol is invalid, then  $\mathcal{E}'$  sends  $\perp$  to  $\mathcal{E}$ . If it is valid and  $\text{sid} = \text{sid}^*$ , it sends  $c^*$  to  $\mathcal{E}$ . Otherwise,  $\mathcal{E}'$  selects a key  $sk \xleftarrow{\$} \mathcal{K}$  uniformly at random, adds  $(\text{sid}, sk)$  to  $\text{keys}$ , finds the unique  $m$  such that  $(\text{sid}, m) \in S_A$ , and sends  $\text{Enc}_{sk}(\text{sid}, m)$  to  $\mathcal{E}$ .
- When  $\mathcal{E}$  calls **get-next<sub>B</sub>**( $\text{sid}, M$ ), if this is not the last message of the protocol  $\mathcal{E}'$  calls its own “oracle” **get-next<sub>B</sub>**( $\text{sid}, M$ ) and passes the response to  $\mathcal{E}$ . If this is the last message of the protocol, it adds  $(\text{sid}, M)$  to  $S_B$ .
- When  $\mathcal{E}$  calls **get-output<sub>B</sub>**( $\text{sid}$ ),  $\mathcal{E}'$  finds the unique  $M$  and  $sk$  such that  $(\text{sid}, M) \in S_B$  and  $(\text{sid}, sk) \in \text{keys}$ . It computes  $\text{Dec}_{sk}(\text{sid}, M)$  and responds with the result.
- When  $\mathcal{E}$  calls **get-secrets**,  $\mathcal{E}'$  responds with  $(\sigma_A, \sigma_B)$ .
- When  $\mathcal{E}$  calls **finalize**( $b^*$ ),  $\mathcal{E}'$  returns  $b^*$ .

Let  $b$  be the challenge bit in the CPA-security game against  $\mathcal{E}'$ . If  $b^\dagger = b$ , the view of  $\mathcal{E}$  is identical to its view in **Game 2**. In this case,  $\mathcal{E}'$  wins if and only if  $\mathcal{E}$  wins the simulated **Game 2**. If  $b^\dagger \neq b$ , the view of  $\mathcal{E}$  is identical to its view in **Game 3** and  $\mathcal{E}'$  wins if and only if  $\mathcal{E}$  loses the simulated game. The result follows. (4.3) ■

It should be clear that any adversary in game **Game  $q + 3$**  can be easily converted into an adversary in **Game  $q + 4$**  with same advantage.

**Claim 4.4.** *If  $\mathcal{W}_A$  is exfiltration resistant against  $\mathbf{B}$  with valid transcripts then for any PPT adversary  $\mathcal{E}$ ,  $|\text{Adv}^{\text{Game } q+3+i}(\mathcal{E}) - \text{Adv}^{\text{Game } q+4+i}(\mathcal{E})|$  is negligible.*

*Proof.* Let  $\tilde{\mathbf{A}}^{(m)}$  be  $\tilde{\mathbf{A}}^*$  with input plaintext fixed to  $m$ .

We construct an adversary  $\mathcal{E}'$  in LEAK (Figure 1) as follows.  $\mathcal{E}'$  first runs the setup procedure of the key-agreement protocol, receiving as output  $\sigma_A, \sigma_B$ , and  $\pi$ . As above,  $\mathcal{E}'$  will simulate many runs of the key-agreement protocol with Alice replaced by  $\mathcal{W}_A \circ \mathbf{A}$  and input  $(\sigma_A, \sigma_B, \pi)$ . So, for convenience, we give  $\mathcal{E}'$  an “oracle interface” to these simulated runs with “oracle” calls `start-run`, `get-nextA` and `get-nextB` as in the key-agreement security game (Figure 3).  $\mathcal{E}'$  selects  $b^\dagger \xleftarrow{\$} \{0, 1\}$  uniformly at random, sets  $j \leftarrow 1$  and  $\text{ids} \leftarrow \emptyset$ , and simulates a run of  $\mathcal{E}$ , responding to oracle calls as follows.

- When  $\mathcal{E}$  calls `start-run(sid, m)`,  $\mathcal{E}'$  adds  $(\text{sid}, j)$  to  $\text{ids}$ . If  $j \neq i$ , it calls its own “oracle” `start-run(sid, m)` and increments  $j$ . If  $j = i$ ,  $\mathcal{E}'$  increments  $j$  and constructs the circuit  $\tilde{B}$  described below with  $\text{sid}^*$ ,  $(\sigma_A, \sigma_B, \pi)$ ,  $b^\dagger$ ,  $j$ ,  $\text{ids}$ , the state of  $\mathcal{E}$ , and the state of the various “oracles” hard-coded into it. It then returns  $(\tilde{\mathbf{A}}^{(m_{b^\dagger})}, \tilde{B}, (\sigma_A, \sigma_B, \pi))$ .
- When  $\mathcal{E}$  calls `get-nextA(sid, M)`,  $\mathcal{E}'$  calls its own “oracle” `get-nextA(sid, M)` and passes the response to  $\mathcal{E}$ .
- When  $\mathcal{E}$  calls `get-nextB(sid, M)`,  $\mathcal{E}'$  calls its own “oracle” `get-nextB(sid, M)` and passes the response to  $\mathcal{E}$ .
- When  $\mathcal{E}$  calls `start-challenge(sid*, m0, m1)`,
- When  $\mathcal{E}$  calls `get-secrets`,  $\mathcal{E}'$  responds with  $(\sigma_A, \sigma_B)$ .

$\tilde{B}$  will play the role of Bob in the key-agreement protocol, and it has the state of  $\mathcal{E}$  and the “oracles” hard-coded into it. It can make “oracle” calls to simulated protocols with Alice replaced by  $\mathcal{W}_A \circ \mathbf{A}$ . It also starts its own “oracle” simulations with Alice replaced by  $\mathcal{W}_A \circ \tilde{\mathbf{A}}^*$  instead. To distinguish these oracles, we write, e.g., `get-next $\mathcal{W}_A \circ \tilde{\mathbf{A}}^*$`  and `get-next $\mathcal{W}_A \circ \mathbf{A}$` . Note that  $\tilde{B}$  is itself playing a game in which it exchanges its own messages with the challenge party  $\mathbf{A}^*$  in LEAK (Figure 1).  $\tilde{B}$  continues to simulate  $\mathcal{E}$  from its current state, responding to oracle calls as follows.

- When  $\mathcal{E}$  calls `start-run(sid, m)`,  $\tilde{B}$  adds  $(\text{sid}, j)$  to  $\text{ids}$ , increments  $j$ , and calls its own “oracle” `start-run(sid, m)`.
- When  $\mathcal{E}$  calls `get-nextA(sid, M)`,  $\tilde{B}$  finds the unique  $k$  such that  $(\text{sid}, k) \in \text{ids}$ . If  $k < i$ ,  $\tilde{B}$  calls its “oracle” `get-next $\mathcal{W}_A \circ \mathbf{A}$ (sid, M)` and passes the response to  $\mathcal{E}$ . If  $k = i$ , then  $\tilde{B}$  sends the message  $M$  to the challenge party  $\mathbf{A}^*$  and passes the response to  $\mathcal{E}$ . If  $k > i$ , then  $\tilde{B}$  calls its “oracle” `get-next $\mathcal{W}_A \circ \tilde{\mathbf{A}}^*$ (sid, M)` and passes the response to  $\mathcal{E}$ .
- When  $\mathcal{E}$  calls `get-nextB(sid, M)`,  $\tilde{B}$  calls its own “oracle” `get-nextB(sid, M)` and passes the response to  $\mathcal{E}$ .
- When  $\mathcal{E}$  calls `finalize(b*)`,  $\tilde{B}$  sets its state to 0 if  $b^* = b^\dagger$  and to 1 otherwise.
- When  $\mathcal{E}$  calls `get-secrets`,  $\tilde{B}$  responds with  $(\sigma_A, \sigma_B)$ .

Finally,  $\mathcal{E}'$  receives the state of  $\tilde{B}$  and simply returns its value.

Let  $b$  be the challenge bit in LEAK. Then, if  $b = 0$  so that the challenge party is honest, the view of  $\mathcal{E}$  is identical to its view in **Game**  $q + 5$ . Then, the final state of  $\tilde{B}$  matches  $b$  if and only if  $\mathcal{E}$  wins this simulated game. If, on the other hand,  $b = 1$ , then the view of  $\mathcal{E}$  is identical to its view in **Game**  $q + 4$ , and the final state of  $\tilde{B}$  matches  $b$  if and only if  $\mathcal{E}$  loses this simulated game. The result follows. (4.4) ■

□

## F Proof of Theorem 9

*Proof of Theorem 9.* Let  $(\text{KeyGen}, \text{Enc}, \text{Dec})$  be a CPA-secure encryption scheme with some reverse firewall  $\mathcal{W}$ . Note that we can view  $\mathcal{W}$  as a map between ciphertexts.

We present a one-bit PKE scheme as follows. Let  $m_0, m_1$  be distinct plaintexts. The public key is then  $(e_0 = \mathcal{W}(\text{Enc}_{sk}(m_0)), e_1 = \mathcal{W}(\text{Enc}_{sk}(m_1)))$ , and the secret key is just the secret key of the underlying scheme. To encrypt a bit  $b$ , we compute  $\mathcal{W}(e_b)$ . The decryption algorithm of the public-key scheme runs the decryption algorithm of the symmetric-key scheme  $\text{Dec}$  and outputs 0 if the result is  $m_0$ , 1 if it is  $m_1$ , and  $\perp$  otherwise.

Let  $\mathcal{E}$  be a PPT adversary in the semantic-security game against the above scheme. We assume without loss of generality that  $\mathcal{E}$  never outputs a pair of identical challenge plaintexts. We construct an efficient tampered encryption algorithm  $\bar{\text{Enc}}$  and an efficient adversary  $\mathcal{E}'$  in the CPA-security game against  $\mathcal{W} \circ \bar{\text{Enc}}$ . Choose  $m_0^\dagger$  and  $m_1^\dagger$  uniformly at random from the plaintext space. Simulate  $\mathcal{E}$  polynomially many times and let  $i$  such in at least polynomial many of these runs, the challenge plaintexts chosen by  $\mathcal{E}$  differ in the  $i$ th bit. Fix  $c_0 = \text{Enc}_{sk}(m_0)$  and  $c_1 = \text{Enc}_{sk}(m_1)$ . Then, we define  $\bar{\text{Enc}}_{sk}(m_b^\dagger) = c_b$  and for all other plaintexts  $m$ ,  $\bar{\text{Enc}}(m) = \mathcal{W}(c_b)$  where  $b$  is the  $i$ th bit of  $m$ . Then,  $\mathcal{E}'$  behaves as follows.

- It calls the encryption oracle on input  $m_0^\dagger$  and  $m_1^\dagger$ . Call the results  $e_0$  and  $e_1$ .
- It runs  $\mathcal{E}$  with input  $pk = (e_0, e_1)$ , receiving as output two challenge plaintexts,  $(m_0^*, m_1^*)$ .  $\mathcal{E}'$  then outputs these as its own challenge plaintexts.
- On input  $c^*$ , a challenge ciphertext,  $\mathcal{E}'$  passes  $c^*$  to  $\mathcal{E}$ , receiving as output a bit  $b$ .
- If  $m_0^*$  and  $m_1^*$  differ in their  $i$ th bit and are distinct from  $m_0^\dagger$  and  $m_1^\dagger$ , output the bit corresponding to the plaintext whose  $i$ th bit is  $b$ . Otherwise, return a uniformly random bit.

Note that the view of  $\mathcal{E}$  is identical to its view in the semantic security game against the public-key scheme. Furthermore, with non-negligible probability, we have that  $m_0^\dagger, m_1^\dagger, m_0^*$ , and  $m_1^*$  are distinct and  $m_0^*$  and  $m_1^*$  differ in their  $i$ th bit. When both of these conditions are satisfied,  $\mathcal{E}'$  guesses correctly if and only if  $\mathcal{E}$  guesses correctly. The result follows. □