



Cryptographic reverse firewalls for interactive proof systems

DOI:

[10.1016/j.tcs.2020.11.043](https://doi.org/10.1016/j.tcs.2020.11.043)

Document Version

Accepted author manuscript

[Link to publication record in Manchester Research Explorer](#)

Citation for published version (APA):

Ganesh, C., Magri, B., & Venturi, D. (2021). Cryptographic reverse firewalls for interactive proof systems. *Theoretical Computer Science*, 855, 104-132. <https://doi.org/10.1016/j.tcs.2020.11.043>

Published in:

Theoretical Computer Science

Citing this paper

Please note that where the full-text provided on Manchester Research Explorer is the Author Accepted Manuscript or Proof version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version.

General rights

Copyright and moral rights for the publications made accessible in the Research Explorer are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Takedown policy

If you believe that this document breaches copyright please refer to the University of Manchester's Takedown Procedures [<http://man.ac.uk/04Y6Bo>] or contact uml.scholarlycommunications@manchester.ac.uk providing relevant details, so we can investigate your claim.



Cryptographic Reverse Firewalls for Interactive Proof Systems

Chaya Ganesh

Department of Computer Science and Automation, Indian Institute of Science, India

chaya@iisc.ac.in

Bernardo Magri

Department of Computer Science, Aarhus University, Denmark

magri@cs.au.dk

Daniele Venturi

Department of Computer Science, Sapienza University of Rome, Italy

venturi@di.uniroma1.it

Abstract

We study interactive proof systems (IPSeS) in a strong adversarial setting where the machines of *honest parties* might be corrupted and under control of the adversary. Our aim is to answer the following, seemingly paradoxical, questions:

- Can Peggy convince Vic of the veracity of an NP statement, without leaking any information about the witness even in case Vic is malicious and Peggy does not trust her computer?
- Can we avoid that Peggy fools Vic into accepting false statements, even if Peggy is malicious and Vic does not trust her computer?

At EUROCRYPT 2015, Mironov and Stephens-Davidowitz introduced cryptographic reverse firewalls (RFs) as an attractive approach to tackling such questions. Intuitively, a RF for Peggy/Vic is an external party that sits between Peggy/Vic and the outside world and whose scope is to sanitize Peggy's/Vic's incoming and outgoing messages in the face of subversion of her/his computer, *e.g.* in order to destroy subliminal channels.

In this paper, we put forward several natural security properties for RFs in the concrete setting of IPSeS. As our main contribution, we construct efficient RFs for different IPSeS derived from a large class of Sigma protocols that we call *malleable*.

A nice feature of our design is that it is completely transparent, in the sense that our RFs can be directly applied to already deployed IPSeS, without the need to re-implement them.

2012 ACM Subject Classification Security and privacy → Cryptography

Keywords and phrases Subversion, Algorithm substitution attacks, Cryptographic reverse firewalls, Interactive proofs, Zero knowledge

Digital Object Identifier [10.4230/LIPIcs.ICALP.2020.55](https://doi.org/10.4230/LIPIcs.ICALP.2020.55)

Related Version A full version of the paper is available at <https://eprint.iacr.org/2020/204>.

Funding *Bernardo Magri*: This work was supported by Concordium Blockchain Research Center, Aarhus University, Denmark



© Chaya Ganesh and Bernardo Magri and Daniele Venturi;
licensed under Creative Commons License CC-BY

47th International Colloquium on Automata, Languages, and Programming (ICALP 2020).

Editors: Artur Czumaj, Anuj Dawar, and Emanuela Merelli; Article No. 55; pp. 55:1–55:17

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

1 Introduction

An interactive proof system (IPS) $\Pi = (P, V)$ allows a prover P to convince a verifier V about the veracity of a public statement $x \in \mathcal{L}$, where \mathcal{L} is an NP language and where both P and V are modeled as interactive PPT machines. The prover is facilitated by possessing a witness w to the fact that, indeed, $x \in \mathcal{L}$, and the interaction with the verifier may consist of several rounds of communication, at the end of which the verifier outputs a verdict on the membership of x in \mathcal{L} .

In order to be useful, an IPS should satisfy the following properties:

- *Completeness*: If $x \in \mathcal{L}$, the honest prover (almost) always convinces the honest verifier.
- *Soundness*: If $x \notin \mathcal{L}$, no (computationally bounded) malicious prover can convince the honest verifier that $x \in \mathcal{L}$. An even stronger guarantee, known as *knowledge soundness* [9], is to require that the only way a prover can convince the honest verifier that $x \in \mathcal{L}$ is to “know” a valid witness w corresponding to x . Such proofs¹ are called *proofs of knowledge* (PoKs).
- *Zero Knowledge (ZK)*: A valid proof reveals nothing beyond the fact that $x \in \mathcal{L}$, and thus in particular it leaks no information about the witness w , even in case the proof is conducted in the presence of a (computationally bounded) malicious verifier [36]. A weaker guarantee, known as *witness indistinguishability* (WI) [24], is that, whenever there are multiple witnesses attesting that $x \in \mathcal{L}$, no (computationally bounded) malicious verifier can distinguish whether a proof is conducted using either of two witnesses.

One of the motivations for studying IPSes with the above properties is that they are ubiquitous in cryptography, with applications ranging from identification protocols [24], blind digital signatures [42], and electronic voting [16], to general-purpose maliciously secure multi-party computation [35].

1.1 Sigma Protocols

While WI/ZK PoKs exist for all of NP, based on minimal cryptographic assumptions [23, 34, 33], efficiency is a different story. Fortunately, it is possible to design practical interactive proofs for specific languages, typically in the form of so-called Sigma protocols. Briefly, a Sigma protocol is a special type of IPS consisting of just three rounds, where the prover sends a first message α (the commitment), the verifier sends a random string β (the challenge), and finally the prover forwards a last message γ (the response). Sigma protocols satisfy two main properties: The first one, known as *special soundness*, is a strong form of knowledge soundness; the second one, known as *honest-verifier zero knowledge* (HVZK), is a weak form of the zero knowledge property that only holds against honest-but-curious verifiers.

The applications of Sigma protocols to cryptographic constructions are countless (see, e.g., [25, 17, 48, 22, 43]). These results are perhaps surprising, as Sigma protocols only satisfy HVZK and thus guarantee no security in the presence of malicious verifiers. In some cases, the solution to this apparent paradox is due to a beautiful technique put forward by Cramer, Damgård, and Schoenmakers [15], which allows to add WI to any Sigma protocol. Moreover, it is relatively easy to transform any Sigma protocol into an interactive ZK PoK at the cost of adding a single round of interaction [33].

¹ Sometimes, the term “proof” is used to refer to statistically sound IPSes, while computationally sound IPSes are typically called “arguments”.

78 **1.2 Our Question**

79 The standard definitions of security for IPSes (implicitly) rely on the assumption that honest
 80 parties can fully trust their machines. In practice, however, such an assumption may just
 81 be too optimistic, as witnessed by the revelations of Edward Snowden about subversion of
 82 cryptographic standards [45, 7], and in light of the numerous (seemingly accidental) bugs in
 83 widespread pieces of cryptographic software [38, 1, 2].

84 Motivated by the above incidents, we ask the following question which constitutes the
 85 main source of inspiration for this work:

86 *Can we design practical interactive proofs that remain secure even if the machines of*
 87 *the honest parties running them have been tampered with?*

88 In order to see why the above question is well motivated and not trivial, let us analyze
 89 the dramatic consequences of subverting the prover of ZK IPSes. Clearly, the problem of
 90 subversion-resistant interactive zero knowledge is just impossible in its utmost generality, as
 91 a subverted prover could just reveal the witness to the verifier. However, one may argue that
 92 these kind of attacks are easily detectable, and thus can be avoided.

93 The problem becomes more interesting if we restrict the subversion to be *undetectable*,
 94 as suggested by Bellare, Paterson, and Rogaway [11] in their seminal work on subversion of
 95 symmetric encryption, where the authors show how to subvert any sufficiently randomized
 96 cipher in an undetectable manner, using rejection sampling. A moment of reflection shows
 97 that their attack can be adapted to the case of IPSes.² The solution proposed by [11] is to
 98 rely on deterministic symmetric encryption. Unfortunately, this approach is not viable for
 99 the case of IPSes, as it is well-known that interactive proofs with deterministic provers can
 100 be zero knowledge only for trivial languages [32, §4.5].

101 **Reverse firewalls**

102 The above described undetectable attacks show that the problem of designing IPSes that
 103 remain secure even when run on untrusted machines is simply impossible if we are not
 104 willing to make any further assumption. In this paper, we study how to tackle subversion
 105 attacks against interactive proofs in the framework of “cryptographic reverse firewalls (RFs)”,
 106 introduced by Mironov and Stephens-Davidowitz [40]. In such a setting, both the prover and
 107 the verifier are equipped with their own RF W , also modeled as an interactive PPT machine,
 108 whose scope is solely to sanitize the parties’ incoming and outgoing messages in the face of
 109 subversion.

110 Importantly, neither the prover nor the verifier put any trust in the RF, meaning that they
 111 are not allowed to share secrets with the firewall itself. The hope is that an uncorrupted³ RF
 112 can provide meaningful security guarantees even in case the honest prover’s and/or verifier’s
 113 machines have been tampered with. Note that a RF can never “create security”, as it does
 114 not even know the inputs to the protocol, but at best can preserve the security guarantees
 115 satisfied by the initial IPS. At the same time, the RF should not ruin the functionality of the

² In particular, a subverted prover with a hardwired secret key k for a pseudorandom function $F_k(\cdot)$, could sample the random coins $r^{(i)}$ needed to generate the honest prover’s message $m^{(i)}$ (for round $i \in \mathbb{N}$) multiple times, until $F_k(m^{(i)})$ leaks one bit of the witness. This attack works provided that at least one of the prover’s messages has high-enough min-entropy.

³ Clearly, if both the machine of the honest party and the firewall are corrupted, there is no hope for security. On the other hand, in case the machine is honest and the firewall is corrupt, the underlying protocol is still secure, since we can simply think of the RF as being part of the adversary [21].

116 underlying IPS, in the sense that the sanitized IPS should still work in case no subversion
117 takes place.

118 Mironov and Stephens-Davidowitz construct general-purpose RFs that can be used in
119 order to preserve both functionality and security of any two-party protocol. It is important
120 to note that since ZK/WI IPSes are a special case of secure two-party computation, their
121 RF constructions already seem to solve our problem.⁴ However, the solutions in [40] are not
122 practical. In particular, one of their RFs increases the round complexity of the initial IPS,
123 and, more importantly, it requires to carry out the underlying IPS in the encrypted domain,
124 thus requiring to completely change the original protocol. In contrast, we seek constructions
125 of RFs that can be applied directly to existing IPSes, without adding any overhead, and
126 without the need to re-implement them.

127 **2 Reverse Firewalls for Interactive Proofs**

128 In this section, we give security definitions for RFs applied to IPSes. Our notions can be
129 seen as special cases of the generic framework by Mironov and Stephens-Davidowitz [40],
130 who defined security of RFs for the more general case of arbitrary two-party protocols.

131 Let $\Pi = (P, V)$ be an IPS for a relation \mathcal{R} . A cryptographic reverse firewall is an external
132 party W that can be attached either to the prover P or to the verifier V , whose scope is
133 to sanitize incoming and outgoing messages in the face of parties' subversion. Importantly,
134 the RF is allowed to keep its own state but cannot share state with any of the parties.
135 Similarly to [40], we model an interactive Turing machine M as a triple of algorithms
136 $M := (M_{\text{next}}, M_{\text{rec}}, M_{\text{out}})$ specified as follows: (i) Algorithm M_{next} takes as input the current
137 state and outputs the next message to be sent; (ii) Algorithm M_{rec} takes as input an incoming
138 message, and updates the state; (iii) Algorithm M_{out} takes as input the final state at the
139 completion of the protocol, and returns a bit.

140 **► Definition 1 (RF for IPSes).** *Let $\Pi = (P, V)$ be an IPS for a relation \mathcal{R} . A cryptographic*
141 *reverse firewall (RF) for Π is a stateful algorithm W that takes as input a message, its state,*
142 *and outputs a sanitized message, together with an updated state. For an interactive Turing*
143 *machine $M = (M_{\text{next}}, M_{\text{rec}}, M_{\text{out}}) \in \{P, V\}$, and RF W , the sanitized machine $W \circ M := \hat{M} =$*
144 *$(\hat{M}_{\text{next}}, \hat{M}_{\text{rec}}, \hat{M}_{\text{out}})$ is specified as follows:*

$$145 \quad \hat{M}_{\text{next}}(\sigma) := W(M_{\text{next}}(\sigma))$$

$$146 \quad \hat{M}_{\text{rec}}(\sigma, m) := M_{\text{rec}}(\sigma, W(m))$$

$$147 \quad \hat{M}_{\text{out}}(\sigma) := M_{\text{out}}(\sigma).$$

149 As our first contribution, we put forward several natural properties that a RF for an IPS
150 might satisfy. In particular, we consider the following notions (see the full version [29] for
151 more formal definitions).

- 152 **■ Completeness preservation:** The sanitized IPS (*i.e.*, the IPS obtained by sanitizing both
153 the honest prover's and the honest verifier's messages) still satisfies completeness.
- 154 **■ Strong soundness preservation:** Whenever $x \notin \mathcal{L}$, no malicious prover can convince the
155 verifier that $x \in \mathcal{L}$, even if the verifier's implementation has been arbitrarily subverted.

⁴ At least to some extent, since, strictly speaking, their results for IPSes are incomparable to ours. We refer the reader to §5.1 for more details.

- 156 ■ *Strong ZK preservation*: A valid proof reveals nothing beyond the fact that $x \in \mathcal{L}$, even
 157 in case the proof is conducted in the presence of a malicious verifier talking to a prover
 158 whose implementation has been arbitrarily subverted.
 - 159 ■ *Strong WI preservation*: Whenever there are multiple witnesses attesting that $x \in \mathcal{L}$, no
 160 malicious verifier talking to a prover whose implementation has been arbitrarily subverted
 161 can distinguish whether a proof is conducted using either of two witnesses.
 - 162 ■ *Strong exfiltration resistance for the prover (resp. verifier)*: Transcripts produced by
 163 running the sanitized IPS in the presence of a malicious verifier (resp. prover) talking
 164 to a prover (resp. verifier) whose implementation has been arbitrarily subverted are
 165 indistinguishable to transcripts produced by running the sanitized IPS in the presence of
 166 a malicious verifier (resp. prover) talking to the honest prover (resp. verifier).
- 167 For each of the above properties (except for completeness), we also consider a weak variant
 168 which only holds w.r.t. *functionality-maintaining* provers/verifiers. Intuitively, a prover is
 169 functionality maintaining if, upon input a valid statement/witness pair, it still convinces the
 170 honest verifier with overwhelming probability. Similarly, a verifier is functionality maintaining
 171 if, upon input a valid statement, it still accepts with overwhelming probability in a protocol
 172 run with the honest prover.

173 What is possible and what is impossible

174 A moment of reflection shows that soundness preservation is impossible to achieve. In fact,
 175 an arbitrarily subverted verifier might always⁵ output one, thus automatically accepting
 176 both true and false statements. Such a verifier is still functionality maintaining,⁶ and thus
 177 this simple attack even rules out *weak* soundness preservation. One way to circumvent
 178 this impossibility would be to only consider *partial subversion*, *i.e.* split the verifier into
 179 two components, one for computing the next messages in the protocol, and the other one
 180 for determining the final verdict on the veracity of a statement; hence, assume the latter
 181 component to be untamperable.

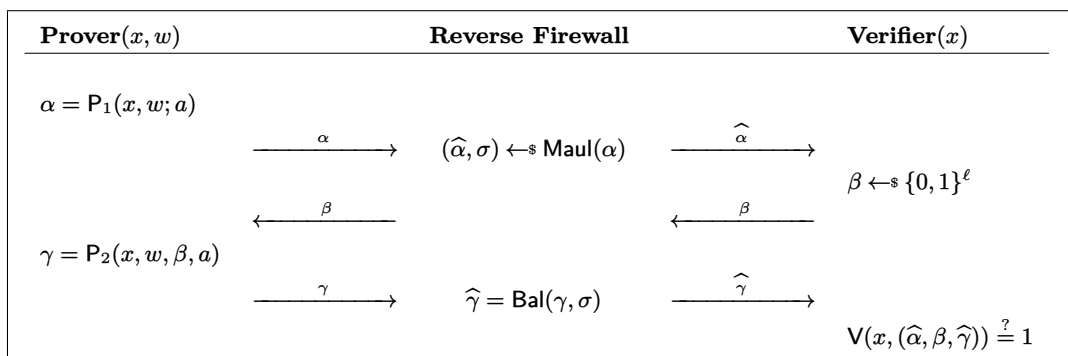
182 Turning to subversion of the prover, consider the subverted prover that always outputs
 183 the all-zero string. The soundness property of the underlying IPS implies that, for any RF
 184 and for any *false* statement $x \notin \mathcal{L}$, a sanitized transcript in this case can never be accepting.
 185 Moreover, assuming the language \mathcal{L} is non-trivial, the latter holds true even in case x is a
 186 *true* statement, which in turn rules out strong exfiltration resistance. For similar reasons,
 187 strong ZK/WI preservation are also impossible to achieve.

188 **3** Firewall Constructions from Malleable Sigma Protocols

189 As our second contribution, we formalize a class of Sigma protocols which admit simple, and
 190 very efficient, RFs for the prover. (See the full version [29] for similar constructions dealing
 191 with functionality-maintaining subversion of the verifier.) The main idea is to use the RF
 192 to re-randomize the prover's messages, in order to destroy any potential subliminal channel
 193 signaling information about the witness. The difficulty, though, is that such re-randomization
 194 must be carried out without knowing a witness, and while at the same time preserving the
 195 completeness property of the underlying IPS.

⁵ If one insists on undetectability, the subverted verifier may output 1 upon some hard-wired, randomly chosen, false statement $\bar{x} \notin \mathcal{L}$.

⁶ The latter is because completeness is a guarantee that only concerns true statements.



■ **Figure 1** Cryptographic reverse firewall for a malleable Sigma protocol

196 For the sake of concreteness, let us describe our firewall applied to the classical Sigma
 197 protocol for proving knowledge of a discrete logarithm [49]. Here, the statement consists of a
 198 description of a cyclic group \mathbb{G} with generator g and prime order q , together with a value
 199 $x \in \mathbb{G}$ such that $x = g^w$ for some $w \in \mathbb{Z}_q$. The prover's first message is a random group
 200 element $\alpha = g^a \in \mathbb{G}$. Finally, the prover's last message is $\gamma = a - w \cdot \beta$, where $\beta \in \mathbb{Z}_q$ is the
 201 verifier's challenge; the verifier accepts (α, β, γ) if and only if $g^\gamma = \alpha \cdot x^{-\beta}$. Our RF sanitizes
 202 the messages α and γ from a possibly subverted implementation of the prover as follows:

$$203 \quad \hat{\alpha} = \alpha \cdot g^\sigma$$

$$204 \quad \hat{\gamma} = \gamma + \sigma,$$

206 for random $\sigma \in \mathbb{Z}_q$. Note that $g^{\hat{\gamma}} = g^a \cdot g^\sigma \cdot x^{-\beta} = \hat{\alpha} \cdot x^{-\beta}$, and thus the RF preserves
 207 completeness.

208 We now sketch the proof of weak HVZK preservation. Observe that for any $\tilde{\alpha} = g^{\tilde{a}}$ sent
 209 by a functionality-maintaining subverted prover, the distribution of $\hat{\alpha} = g^{a+\sigma}$ is uniform
 210 over \mathbb{G} and independent of $\tilde{\alpha}, \tilde{a}$, and in fact it is identical to the distribution of α in an
 211 honest run of the original Sigma protocol (without the firewall). As for $\hat{\gamma}$, note that if there
 212 would be two possible values γ, γ' which make both $\tau = (\alpha, \beta, \gamma)$ and $\tau' = (\alpha, \beta, \gamma')$ valid
 213 transcripts, the choice of which response to pick could be used by a functionality-maintaining
 214 subverted prover as a subliminal channel signaling information about the witness. Hence,
 215 we exploit the fact that for any prefix α, β , there exists a unique response γ such that the
 216 verifier accepts upon input x and (α, β, γ) .

217 It follows that the distribution of $\hat{\gamma}$ is identical to that of γ in an honest run of the
 218 original Sigma protocol (without the firewall). Putting it all together, we have shown that
 219 the distribution of a sanitized transcript $\hat{\tau} = (\hat{\alpha}, \beta, \hat{\gamma})$ is identical to the distribution of an
 220 honest transcript $\tau = (\alpha, \beta, \gamma)$. Thus, weak HVZK preservation follows by the fact that
 221 Schnorr's Sigma protocol is HVZK.

222 3.1 HVZK Preservation

223 Let us now explain how to generalize the above idea to a large class of Sigma protocols
 224 that we call *malleable*. In what follows, given a Sigma protocol $\Sigma = (P, V)$, we denote by
 225 P_1 and P_2 the algorithms that compute, respectively, the first prover's message α , and the

226 last prover’s message (response) γ . The challenge space is represented⁷ as $\{0, 1\}^\ell$, so that
 227 there are 2^ℓ possible challenges, and we write V for the algorithm that the verifier runs upon
 228 statement x and transcript τ to make its final decision. Let \mathcal{A} be the space of all possible
 229 prover’s first messages; we assume that membership in \mathcal{A} can be tested efficiently, so that V
 230 always outputs \perp whenever $\alpha \notin \mathcal{A}$.

231 As for the case of Schnorr’s Sigma protocol, an additional requirement that we need is
 232 that the prover’s responses are unique, meaning that for all $x \in \mathcal{L}$, and for any $\alpha \in \mathcal{A}$ and
 233 $\beta \in \{0, 1\}^\ell$, there exists at most one⁸ value γ such that $V(x, (\alpha, \beta, \gamma)) = 1$.

234 Intuitively, a Sigma protocol is malleable if there exists an efficient algorithm Maul for
 235 randomizing the prover’s first message α into a value $\hat{\alpha}$ which is distributed identically to
 236 the first message of an honest prover. Moreover, for any challenge β , given the coins used
 237 to randomize α and any response γ yielding a valid transcript $\tau = (\alpha, \beta, \gamma)$, there exists an
 238 efficient algorithm Bal for computing a balanced response $\hat{\gamma}$ such that $(\hat{\alpha}, \beta, \hat{\gamma})$ is also valid.
 239 As we show in the full version [29], many natural Sigma protocols are already malleable.
 240 In particular, the latter holds true for Maurer’s unifying protocol [39], which includes the
 241 protocols by Fiat-Shamir [25], Guillou-Quisquater [37], Schnorr [49], Okamoto [41], and many
 242 others as special cases.

243 Our RF construction is depicted in Fig. 1. Intuitively, the firewall uses the malleability
 244 property of the underlying Sigma protocol in order to re-randomize the prover’s first and
 245 last messages, in such a way that a functionality-maintaining subverted prover cannot signal
 246 information about the witness through them. The theorem below, whose proof appears in
 247 the full version [29], establishes its security.

248 **► Theorem 2.** *Let $\Sigma = (P = (P_1, P_2), V)$ be a malleable Sigma protocol with unique responses,*
 249 *for a relation \mathcal{R} . The RF W of Fig. 1 preserves completeness, and is weakly HVZK preserving*
 250 *for the prover.*

251 3.2 ZK Preservation

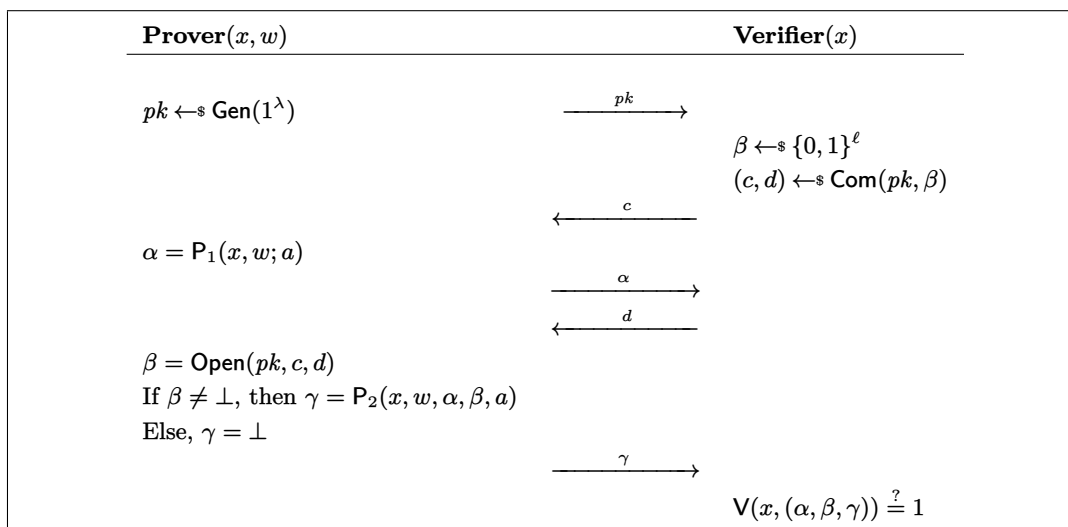
252 As Sigma protocols are not in general zero knowledge, there is no hope to prove that the
 253 above firewall weakly preserves ZK. However, a standard trick [33] allows to transform any
 254 Sigma protocol into a 5-round IPS satisfying ZK. The idea is to let the prover send the public
 255 key pk of a commitment scheme ($\text{Gen}, \text{Com}, \text{Open}$) during the first round. Then, during the
 256 second round, the verifier forwards to the prover a commitment c to the challenge β . Finally,
 257 the Sigma protocol is executed as before with the difference that the verifier also needs to
 258 open the commitment, with the prover aborting if the opening is invalid. We depict such a
 259 modified protocol in Fig. 2.

260 In order to build a RF for this IPS, we need to sanitize the additional messages from the
 261 (possibly subverted, but functionality-maintaining) prover.⁹ We do so by relying on a special
 262 type of *key-malleable* commitment, which intuitively allows to *maul* any public key pk (via
 263 an algorithm MaulKey) into a uniformly random public key \widehat{pk} , in such a way that, given
 264 a commitment c with opening d w.r.t. \widehat{pk} , it is possible to map (c, d) into a commitment
 265 \widehat{c} with opening \widehat{d} w.r.t. pk , without changing the message inside the commitment. We

⁷ In the case of Schnorr’s Sigma protocol, the challenge space is a cyclic group. However, we can embed such group in $\{0, 1\}^\ell$ for some $\ell \in \mathbb{N}$.

⁸ This property is met by many natural Sigma protocols, and was already considered in several previous works [26, 22, 51].

⁹ The other messages are sanitized as before, *i.e.* we still exploit the fact that the underlying Sigma protocol is malleable.



■ **Figure 2** Sigma protocol compiled with standard techniques to obtain full zero knowledge

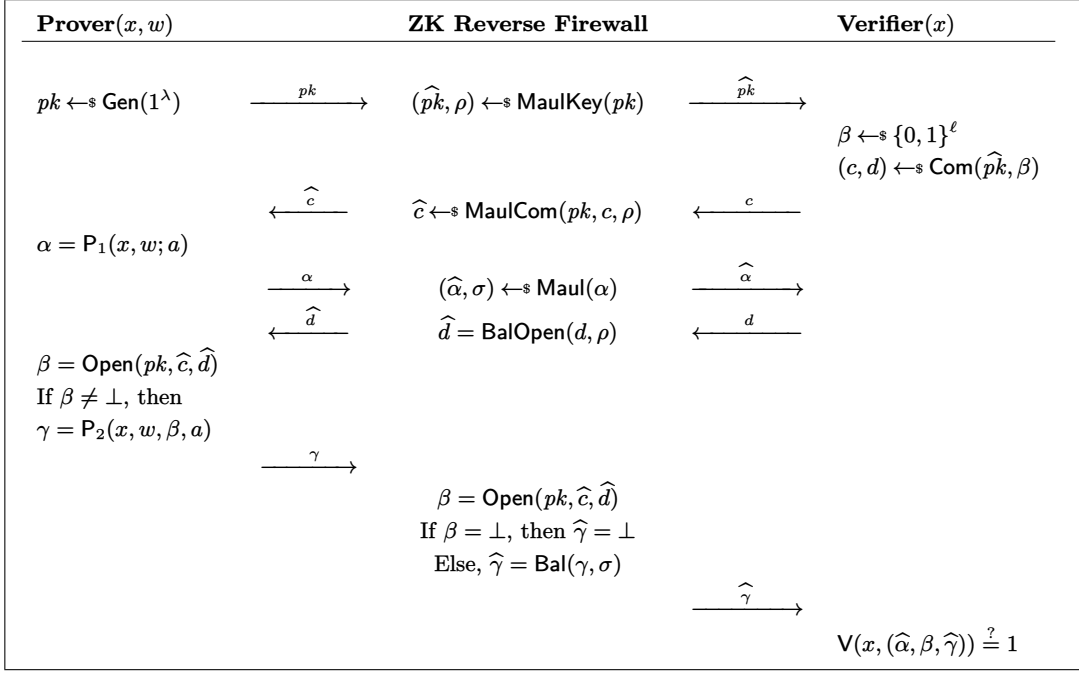
266 denote by MaulCom and BalOpen , respectively, the algorithms for mauling the commitment
 267 c and the opening d , and additionally require that the distribution of mauled public keys
 268 and commitments is identical, respectively, to that of honestly computed public keys and
 269 commitments. As we show in the full version [29], the standard Pedersen's commitment [44]
 270 is easily seen to be key malleable, thus yielding a concrete instantiation under the Discrete
 271 Logarithm assumption.

272 Our RF for the protocol of Fig. 2 is depicted in Fig. 3. The theorem below, whose proof
 273 appears in the full version [29], establishes its security.

274 ► **Theorem 3.** *Let $\Sigma = (P = (P_1, P_2), V)$ be a malleable Sigma protocol with unique responses,
 275 for a relation \mathcal{R} . Let $\Gamma = (\text{Gen}, \text{Com}, \text{Open})$ be a key-malleable commitment scheme with
 276 message space $\{0, 1\}^\ell$. The RF W of Fig. 3 preserves completeness, and moreover is weakly
 277 exfiltration resistant and weakly zero-knowledge preserving for the prover.*

278 ► **Remark 4 (On knowledge soundness).** The IPS of Fig. 2 satisfies soundness, but is not in
 279 general a proof of knowledge. However, we would like to note that the prover's firewall still
 280 works for the standard transformation of a Sigma protocol into a zero-knowledge proof of
 281 knowledge. In such a transformation, a *trapdoor* commitment scheme is used to commit to
 282 the verifier's challenge. Then, after the verifier decommits, the prover sends the trapdoor to
 283 the verifier. This allows an extractor to learn the trapdoor, rewind the prover, and open the
 284 commitment to a different challenge, thus learning the response for two different challenges,
 285 which allows it to obtain a witness using special soundness.

286 The prover's RF for this protocol stays the same, except that it additionally needs to
 287 provide a trapdoor for the mauled public key \widehat{pk} given a trapdoor for the original public key
 288 pk . This is possible, for instance, using Pedersen's commitment, where given a public key
 289 $pk = (g, h = g^k)$ with trapdoor k , we can maul the key to $(\widehat{g} = g^{t_1}, \widehat{h} = h^{t_2})$ for random t_1, t_2 .
 290 Given the trapdoor k for the key pk , the trapdoor for the mauled key \widehat{pk} can be computed
 291 as $t_2 t_1^{-1} k$.



■ **Figure 3** Prover's RF for the protocol in Fig. 2

292 4 Firewalls for Proving Compound Statements

293 In this section, we show how to construct firewalls for Sigma protocols that prove compound
294 statements.

295 Given two Sigma protocols Σ_0 and Σ_1 for NP languages \mathcal{L}_0 and \mathcal{L}_1 , it is easy to obtain a
296 Sigma protocol Σ_{AND} for the NP language $\mathcal{L}_{\text{AND}} = \{(x_0, x_1) : x_0 \in \mathcal{L}_0 \wedge x_1 \in \mathcal{L}_1\}$ by simply
297 running Σ_0 and Σ_1 in parallel, with the verifier sending a single challenge. In a similar vein,
298 the OR technique by Cramer, Damgård, and Schoenmakers [15] allows to obtain a Sigma
299 protocol Σ_{OR} for the NP language $\mathcal{L}_{\text{OR}} = \{(x_0, x_1) : x_0 \in \mathcal{L}_0 \vee x_1 \in \mathcal{L}_1\}$. Importantly, if
300 Σ_0 and Σ_1 are both perfect HVZK, Σ_{OR} satisfies perfect WI. On the other hand, Garay *et*
301 *al.* [30] showed that if Σ_0 and Σ_1 are computational HVZK, Σ_{OR} satisfies computational WI,
302 although the latter holds only in case both statements x_0, x_1 in the definition of language
303 \mathcal{L}_{OR} are true (but the prover knows either a witness for x_0 or for x_1).

304 As long as Σ_0 and Σ_1 are malleable, it is easy to build RFs for Σ_{AND} and Σ_{OR} using
305 our techniques. The RF for Σ_{AND} weakly preserves HVZK, whereas the RF for Σ_{OR} weakly
306 preserves both HVZK and WI.

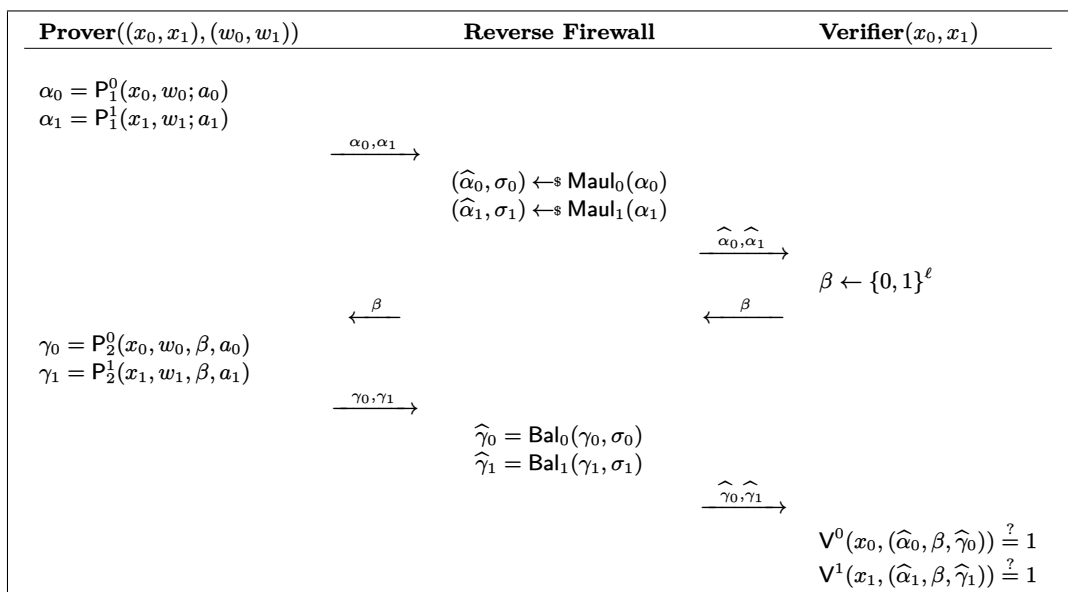
307 4.1 AND Composition

308 Given x_0, x_1 , a prover wishes to prove to a verifier that $x_0 \in \mathcal{L}_0$ and $x_1 \in \mathcal{L}_1$. More precisely,
309 consider the derived relation:

$$310 \quad \mathcal{R}_{\text{AND}} = \{((x_0, x_1), (w_0, w_1)) : (x_0, w_0) \in \mathcal{R}_0 \wedge (x_1, w_1) \in \mathcal{R}_1\}.$$

311 Let $\Sigma_0 = ((P_1^0, P_2^0), V^0)$ (resp. $\Sigma_1 = ((P_1^1, P_2^1), V^1)$) be a Sigma protocol for language \mathcal{L}_0
312 (resp. \mathcal{L}_1). A Sigma protocol Σ_{AND} for the relation \mathcal{R}_{AND} can be obtained by running the
313 two provers of Σ_0 and Σ_1 in parallel, with the verifier sending a single challenge for both

314 executions. Fig. 4 shows a RF for the prover of Σ_{AND} , assuming that both Σ_0 and Σ_1 are
 315 malleable. We prove the following result, whose proof appears in the full version [29].



316 **Figure 4** Reverse firewall for the AND composition of Sigma protocols

317

318 **Theorem 5.** Let $\Sigma_0 = (P^0 = (P_1^0, P_2^0), V^0)$ and $\Sigma_1 = (P^1 = (P_1^1, P_2^1), V^1)$ be malleable Sigma
 319 protocols with unique responses, for relations \mathcal{R}_0 and \mathcal{R}_1 . The RF W of Fig. 4 preserves
 320 completeness, and is weakly HVZK preserving for the prover of the Sigma protocol Σ_{AND} for
 321 relation \mathcal{R}_{AND} .

322 4.2 OR Composition

323

324 Given x_0, x_1 , a prover wishes to prove to a verifier that either $x_0 \in \mathcal{L}_0$ or $x_1 \in \mathcal{L}_1$ (without
 325 revealing which one is the case). More precisely, consider the derived relation

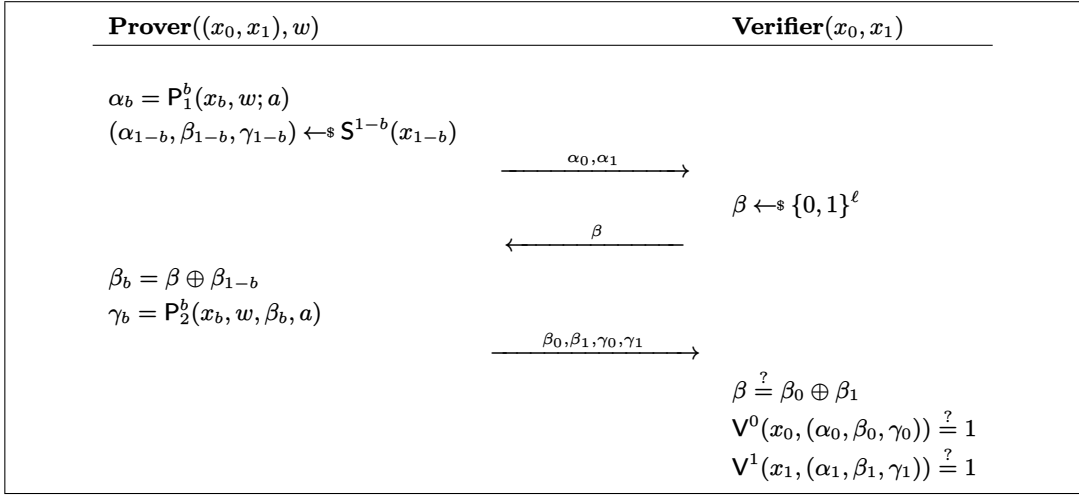
$$326 \mathcal{R}_{\text{OR}} = \{((x_0, x_1), w) : (x_0, w) \in \mathcal{R}_0 \vee (x_1, w) \in \mathcal{R}_1\}.$$

327 Let $\Sigma_0 = ((P_1^0, P_2^0), V^0)$ (resp. $\Sigma_1 = ((P_1^1, P_2^1), V^1)$) be a Sigma protocol for language \mathcal{L}_0 (resp.
 328 \mathcal{L}_1); we denote by S^0 (resp. S^1) the HVZK simulator for Σ_0 (resp. Σ_1). A Sigma protocol
 329 Σ_{OR} for the relation \mathcal{R}_{OR} has been constructed for the first time in [15], where the authors
 330 showed that Σ_{OR} satisfies both (perfect) special HVZK and (perfect) WI. We describe the
 331 protocol Σ_{OR} in Fig. 5, and depict our RF for the prover in Fig. 6.

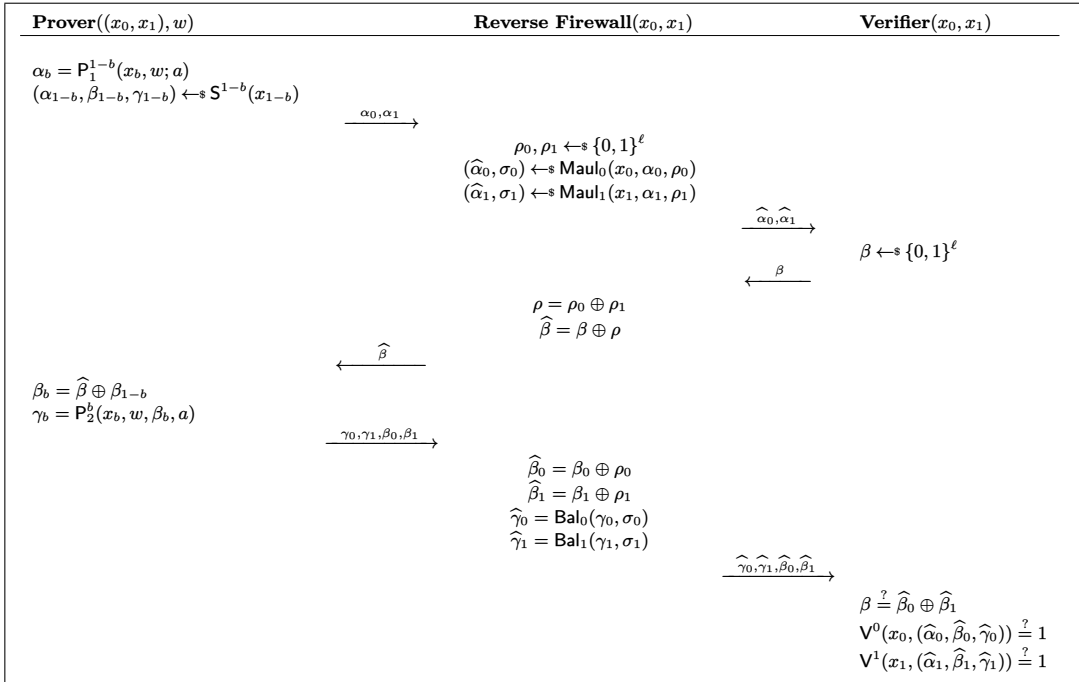
332 As in the case of AND composition, we still rely on the fact that the input Sigma
 333 protocols Σ_0, Σ_1 are malleable. An additional difficulty, however, stems from the fact that a
 334 functionality maintaining prover could now try to change the distribution of the challenges
 335 β_0, β_1 in such a way that, even if $\beta_0 \oplus \beta_1 = \beta$, the pair (β_0, β_1) signals some information
 336 about the witness w or about the hidden bit b . Intuitively, the RF in Fig. 6 tackles this attack
 337 by randomizing the challenges β, β_0, β_1 . The latter requires a different form of malleability
 338 from the underlying Sigma protocols, which we dub *instance-dependent malleability*, where it
 339 should be possible to maul the prover's first message in such a way that we can later balance
 the prover's last message as well as the verifier's challenge.

For the RF in Fig. 6, we prove the following result, whose proof appears in the full
 version [29] of this paper.

55:10 Cryptographic Reverse Firewalls for Interactive Proof Systems



■ **Figure 5** OR composition of Sigma protocols, where $b \in \{0, 1\}$ is s.t. $(x_b, w) \in \mathcal{R}_b$.



■ **Figure 6** Reverse Firewall for the basic OR composition of Sigma protocols, where $b \in \{0, 1\}$ is s.t. $(x_b, w) \in \mathcal{R}_b$.

340 ► **Theorem 6.** Let $\Sigma_0 = (P^0 = (P_1^0, P_2^0), V^0)$ and $\Sigma_1 = (P^1 = (P_1^1, P_2^1), V^1)$ be instance-
 341 dependent malleable Sigma protocols with unique responses, for relations \mathcal{R}_0 and \mathcal{R}_1 . The
 342 RFW of Fig. 6 preserves completeness, and is weakly HVZK/WI preserving for the prover
 343 of the Sigma protocol Σ_{OR} for relation \mathcal{R}_{OR} .

5 Previous Work

5.1 Comparison with Mironov and Stephens-Davidowitz

In their original paper, Mironov and Stephens-Davidowitz [40] build RFs for arbitrary two-party protocols. While their results are related to ours, since IPSes are just a special case of two-party computation, there are some crucial differences which we highlight below.

The first RF construction sanitizes a specific combination of re-randomizable garbled circuits and oblivious transfer, for obtaining general-purpose private function evaluation. The second RF construction sanitizes any two-party protocol, at the price of encrypting the full transcript under public keys that are broadcast at the beginning of the protocol. Both constructions can be instantiated based on (variants of) the DDH assumption. When cast to IPSes, their results yield:

- (i) A RF for the prover that weakly preserves ZK. This is comparable to our RF achieving weak ZK preservation using malleable Sigma protocols and key-malleable commitments. However, our constructions have the advantage that we do not need to change the initial IPS, and thus our RF can be applied directly to already existing implementations in a fully transparent manner (and without introducing any overhead).
- (ii) A RF for the prover satisfying a property called strong exfiltration resistance *against an eavesdropper*, which means that exfiltration resistance holds w.r.t. an arbitrarily subverted prover talking to the *honest verifier*. Note that the latter does not contradict our impossibility result ruling out strong ZK preservation, as our attacks crucially rely on the fact that the distinguisher can (passively) corrupt the verifier.
- (iii) A RF for the verifier satisfying both strong exfiltration resistance and the following weak guarantee: No malicious prover can find statements x_0, x_1 such that it can distinguish transcripts obtained by talking to an arbitrarily subverted verifier holding either input x_0 or input x_1 . Note that the latter does not contradict our impossibility result that rules out weak soundness preservation, since none of the above guarantees imply soundness preservation.

We observe that the above results have at least one of the following drawbacks: (i) The RF is not transparent, *i.e.* it cannot be applied to the initial protocol as is; (ii) The resulting sanitized protocol is not efficient, as we first need to encode the function being computed as a circuit.

Our techniques allow to overcome these limitations in the concrete case of IPSes, as our RFs are both transparent (*i.e.* they can be applied directly to already deployed protocols) and efficient (*i.e.* the sanitized IPSes have exactly the same efficiency as the original, both in terms of round and communication complexity). We see this as the main novelty of our work.

5.2 Additional Related Works

Besides the already mentioned constructions, RFs have also been realized in other settings including digital signatures [5], secure message transmission and key exchange [21, 12], and oblivious transfer [40, 12].

Moreover, a few other lines of research recently¹⁰ emerged to tackle the challenge of protecting cryptographic algorithms against (different forms of) subversion. We review the

¹⁰ All these research directions have their roots in the seminal works of Young and Yung [52] and Simmons [50], in the settings of kleptography and subliminal channels.

386 main ones below.

387 **Algorithm substitution attacks**

388 Bellare, Patterson, and Rogaway [11] studied subversion of symmetric encryption schemes in
 389 the form of algorithm substitution attacks (ASAs). In particular, they show that *undetectable*
 390 subversion of the encryption algorithm is possible, and may lead to severe security breaches;
 391 moreover, they prove that deterministic, stateful, ciphers are secure against the same type of
 392 ASAs. Follow-up works improved the original paper in several aspects [18, 10], and explored
 393 the power of ASAs in other contexts, *e.g.* digital signatures [5], secret sharing [31], and
 394 message authentication codes [3].

395 **Backdoors**

396 Another form of subversion consists of all those attacks that surreptitiously generate public
 397 parameters (primes, curves, etc.) together with secret backdoors that allow to bypass security.
 398 The study of this type of subversion is motivated by the DUAL_EC_DRBG PRG incident.

399 A formal study of parameters subversion has been considered for several primitives, includ-
 400 ing pseudorandom generators [20, 19], hash functions [27], non-interactive zero knowledge [8],
 401 and public-key encryption [6].

402 **Cliptography**

403 Russell *et al.* [46] (see also [47, 4]) consider a different approach to the immunization of
 404 cryptosystems against complete subversion (*i.e.*, when all algorithms can be subverted by the
 405 attacker): offline/online black-box testing. This amounts to introducing an external entity,
 406 called the watchdog, whose goal is to test, either in an online or in an offline fashion, whether
 407 a given cryptographic implementation is compliant with its specification.

408 Hence, a cryptosystem is deemed secure against complete subversion if there exists a
 409 universal watchdog such that, for every attacker subverting all algorithms, either the watchdog
 410 detects subversion with high probability, or the cryptoscheme remains secure even when
 411 using its subverted implementation.

412 **Self-guarding**

413 Yet another approach towards thwarting subversion is that of self-guarding [28]. The idea
 414 here is to assume a trusted initialization phase in which the honest parties possess a genuine
 415 implementation of the cryptosystem, before subversion takes place. This phase is used
 416 in order to generate samples that will be exploited later, together with additional simple
 417 operations that need to be implemented from scratch, to prevent leakage in the face of
 418 subversion attacks.

419 **6 Conclusion**

420 We showed how to design cryptographic reverse firewalls allowing to preserve security of
 421 interactive proof systems in the face of subversion. Our firewalls apply to a large class
 422 of Sigma protocols meeting a natural malleability property, and can be extended to cover
 423 classical applications of Sigma protocols for designing zero-knowledge proofs and for proving
 424 compound statements.

425 We leave it as an intriguing open problem to design a reverse firewall for the OR
426 composition of Sigma protocols that are delayed input, as considered in [13, 14].

427 — References

- 428 **1** Vulnerability summary for cve-2014-6271 (shellshock), September 2014. URL: <http://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2014-6271>.
429
- 430 **2** Juniper vulnerability, 2015. URL: <https://kb.juniper.net/InfoCenter/index?page=content&id=JSA10713>.
431
- 432 **3** Marcel Armour and Bertram Poettering. Substitution attacks against message authentication. *IACR Trans. Symmetric Cryptol.*, 2019(3):152–168, 2019.
433
- 434 **4** Giuseppe Ateniese, Danilo Francati, Bernardo Magri, and Daniele Venturi. Public immunization
435 against complete subversion without random oracles. In Robert H. Deng, Valérie Gauthier-
436 Umaña, Martín Ochoa, and Moti Yung, editors, *ACNS 19*, volume 11464 of *LNCS*, pages
437 465–485. Springer, Heidelberg, June 2019. doi:10.1007/978-3-030-21568-2_23.
- 438 **5** Giuseppe Ateniese, Bernardo Magri, and Daniele Venturi. Subversion-resilient signature
439 schemes. In Indrajit Ray, Ninghui Li, and Christopher Kruegel, editors, *ACM CCS 2015*,
440 pages 364–375. ACM Press, October 2015. doi:10.1145/2810103.2813635.
- 441 **6** Benedikt Auerbach, Mihir Bellare, and Eike Kiltz. Public-key encryption resistant to parameter
442 subversion and its realization from efficiently-embeddable groups. In Michel Abdalla and
443 Ricardo Dahab, editors, *PKC 2018, Part I*, volume 10769 of *LNCS*, pages 348–377. Springer,
444 Heidelberg, March 2018. doi:10.1007/978-3-319-76578-5_12.
- 445 **7** James Ball, Julian Borger, and Glenn Greenwald. Revealed: How US and UK spy agencies
446 defeat internet privacy and security. *Guardian Weekly*, September 2013.
- 447 **8** Mihir Bellare, Georg Fuchsbauer, and Alessandra Scafuro. NIZKs with an untrusted CRS:
448 Security in the face of parameter subversion. In Jung Hee Cheon and Tsuyoshi Takagi, editors,
449 *ASIACRYPT 2016, Part II*, volume 10032 of *LNCS*, pages 777–804. Springer, Heidelberg,
450 December 2016. doi:10.1007/978-3-662-53890-6_26.
- 451 **9** Mihir Bellare and Oded Goldreich. On defining proofs of knowledge. In Ernest F. Brickell,
452 editor, *CRYPTO'92*, volume 740 of *LNCS*, pages 390–420. Springer, Heidelberg, August 1993.
453 doi:10.1007/3-540-48071-4_28.
- 454 **10** Mihir Bellare, Joseph Jaeger, and Daniel Kane. Mass-surveillance without the state: Strongly
455 undetectable algorithm-substitution attacks. In Indrajit Ray, Ninghui Li, and Christopher
456 Kruegel, editors, *ACM CCS 2015*, pages 1431–1440. ACM Press, October 2015. doi:10.1145/
457 2810103.2813681.
- 458 **11** Mihir Bellare, Kenneth G. Paterson, and Phillip Rogaway. Security of symmetric encryption
459 against mass surveillance. In Juan A. Garay and Rosario Gennaro, editors, *CRYPTO 2014*,
460 *Part I*, volume 8616 of *LNCS*, pages 1–19. Springer, Heidelberg, August 2014. doi:10.1007/
461 978-3-662-44371-2_1.
- 462 **12** Rongmao Chen, Yi Mu, Guomin Yang, Willy Susilo, Fuchun Guo, and Mingwu Zhang.
463 Cryptographic reverse firewall via malleable smooth projective hash functions. In Jung Hee
464 Cheon and Tsuyoshi Takagi, editors, *ASIACRYPT 2016, Part I*, volume 10031 of *LNCS*, pages
465 844–876. Springer, Heidelberg, December 2016. doi:10.1007/978-3-662-53887-6_31.
- 466 **13** Michele Ciampi, Giuseppe Persiano, Alessandra Scafuro, Luisa Siniscalchi, and Ivan Visconti.
467 Improved OR-composition of sigma-protocols. In Eyal Kushilevitz and Tal Malkin, editors,
468 *TCC 2016-A, Part II*, volume 9563 of *LNCS*, pages 112–141. Springer, Heidelberg, January
469 2016. doi:10.1007/978-3-662-49099-0_5.
- 470 **14** Michele Ciampi, Giuseppe Persiano, Alessandra Scafuro, Luisa Siniscalchi, and Ivan Visconti.
471 Online/offline OR composition of sigma protocols. In Marc Fischlin and Jean-Sébastien Coron,
472 editors, *EUROCRYPT 2016, Part II*, volume 9666 of *LNCS*, pages 63–92. Springer, Heidelberg,
473 May 2016. doi:10.1007/978-3-662-49896-5_3.
- 474 **15** Ronald Cramer, Ivan Damgård, and Berry Schoenmakers. Proofs of partial knowledge and
475 simplified design of witness hiding protocols. In Yvo Desmedt, editor, *CRYPTO'94*, volume 839
476 of *LNCS*, pages 174–187. Springer, Heidelberg, August 1994. doi:10.1007/3-540-48658-5_19.

- 477 16 Ronald Cramer, Rosario Gennaro, and Berry Schoenmakers. A secure and optimally efficient
478 multi-authority election scheme. In Walter Fumy, editor, *EUROCRYPT'97*, volume 1233 of
479 *LNCS*, pages 103–118. Springer, Heidelberg, May 1997. doi:10.1007/3-540-69053-0_9.
- 480 17 Ivan Damgård and Jens Groth. Non-interactive and reusable non-malleable commitment
481 schemes. In *35th ACM STOC*, pages 426–437. ACM Press, June 2003. doi:10.1145/780542.
482 780605.
- 483 18 Jean Paul Degabriele, Pooya Farshim, and Bertram Poettering. A more cautious approach to
484 security against mass surveillance. In Gregor Leander, editor, *FSE 2015*, volume 9054 of *LNCS*,
485 pages 579–598. Springer, Heidelberg, March 2015. doi:10.1007/978-3-662-48116-5_28.
- 486 19 Jean Paul Degabriele, Kenneth G. Paterson, Jacob C. N. Schuldt, and Joanne Woodage.
487 Backdoors in pseudorandom number generators: Possibility and impossibility results. In
488 Matthew Robshaw and Jonathan Katz, editors, *CRYPTO 2016, Part I*, volume 9814 of *LNCS*,
489 pages 403–432. Springer, Heidelberg, August 2016. doi:10.1007/978-3-662-53018-4_15.
- 490 20 Yevgeniy Dodis, Chaya Ganesh, Alexander Golovnev, Ari Juels, and Thomas Ristenpart. A
491 formal treatment of backdoored pseudorandom generators. In Elisabeth Oswald and Marc
492 Fischlin, editors, *EUROCRYPT 2015, Part I*, volume 9056 of *LNCS*, pages 101–126. Springer,
493 Heidelberg, April 2015. doi:10.1007/978-3-662-46800-5_5.
- 494 21 Yevgeniy Dodis, Ilya Mironov, and Noah Stephens-Davidowitz. Message transmission with
495 reverse firewalls—secure communication on corrupted machines. In Matthew Robshaw and
496 Jonathan Katz, editors, *CRYPTO 2016, Part I*, volume 9814 of *LNCS*, pages 341–372. Springer,
497 Heidelberg, August 2016. doi:10.1007/978-3-662-53018-4_13.
- 498 22 Sebastian Faust, Markulf Kohlweiss, Giorgia Azzurra Marson, and Daniele Venturi. On the non-
499 malleability of the Fiat-Shamir transform. In Steven D. Galbraith and Mridul Nandi, editors,
500 *INDOCRYPT 2012*, volume 7668 of *LNCS*, pages 60–79. Springer, Heidelberg, December 2012.
501 doi:10.1007/978-3-642-34931-7_5.
- 502 23 Uriel Feige, Dror Lapidot, and Adi Shamir. Multiple non-interactive zero knowledge proofs
503 based on a single random string (extended abstract). In *31st FOCS*, pages 308–317. IEEE
504 Computer Society Press, October 1990. doi:10.1109/FSCS.1990.89549.
- 505 24 Uriel Feige and Adi Shamir. Witness indistinguishable and witness hiding protocols. In *22nd*
506 *ACM STOC*, pages 416–426. ACM Press, May 1990. doi:10.1145/100216.100272.
- 507 25 Amos Fiat and Adi Shamir. How to prove yourself: Practical solutions to identification and
508 signature problems. In Andrew M. Odlyzko, editor, *CRYPTO'86*, volume 263 of *LNCS*, pages
509 186–194. Springer, Heidelberg, August 1987. doi:10.1007/3-540-47721-7_12.
- 510 26 Marc Fischlin. Communication-efficient non-interactive proofs of knowledge with online
511 extractors. In Victor Shoup, editor, *CRYPTO 2005*, volume 3621 of *LNCS*, pages 152–168.
512 Springer, Heidelberg, August 2005. doi:10.1007/11535218_10.
- 513 27 Marc Fischlin, Christian Janson, and Sogol Mazaheri. Backdoored hash functions: Immunizing
514 HMAC and HKDF. In *31st IEEE Computer Security Foundations Symposium, CSF 2018*,
515 *Oxford, United Kingdom, July 9-12, 2018*, pages 105–118, 2018.
- 516 28 Marc Fischlin and Sogol Mazaheri. Self-guarding cryptographic protocols against algorithm
517 substitution attacks. In *31st IEEE Computer Security Foundations Symposium, CSF 2018*,
518 *Oxford, United Kingdom, July 9-12, 2018*, pages 76–90, 2018.
- 519 29 Chaya Ganesh, Bernardo Magri, and Daniele Venturi. Cryptographic reverse firewalls for
520 interactive proof systems. *IACR Cryptology ePrint Archive*, 2020:204, 2020. URL: <https://eprint.iacr.org/2020/204>.
- 521
- 522 30 Juan A. Garay, Philip D. MacKenzie, and Ke Yang. Strengthening zero-knowledge pro-
523 tocols using signatures. *Journal of Cryptology*, 19(2):169–209, April 2006. doi:10.1007/
524 s00145-005-0307-3.
- 525 31 Irene Giacomelli, Ruxandra F. Olimid, and Samuel Ranellucci. Security of linear secret-sharing
526 schemes against mass surveillance. In Michael Reiter and David Naccache, editors, *CANS 15*,
527 *LNCS*, pages 43–58. Springer, Heidelberg, December 2015. doi:10.1007/978-3-319-26823-1_
528 4.

- 529 **32** Oded Goldreich. *Foundations of Cryptography: Basic Tools*, volume 1. Cambridge University
530 Press, Cambridge, UK, 2001.
- 531 **33** Oded Goldreich and Ariel Kahan. How to construct constant-round zero-knowledge proof
532 systems for NP. *Journal of Cryptology*, 9(3):167–190, June 1996.
- 533 **34** Oded Goldreich, Silvio Micali, and Avi Wigderson. Proofs that yield nothing but their validity
534 or all languages in NP have zero-knowledge proof systems. *Journal of the ACM*, 38(3):691–729,
535 1991.
- 536 **35** Shafi Goldwasser and Silvio Micali. Probabilistic encryption and how to play mental poker
537 keeping secret all partial information. In *14th ACM STOC*, pages 365–377. ACM Press, May
538 1982. doi:10.1145/800070.802212.
- 539 **36** Shafi Goldwasser, Silvio Micali, and Charles Rackoff. The knowledge complexity of interactive
540 proof systems. *SIAM Journal on Computing*, 18(1):186–208, 1989.
- 541 **37** Louis C. Guillou and Jean-Jacques Quisquater. A practical zero-knowledge protocol fitted to
542 security microprocessor minimizing both transmission and memory. In C. G. Günther, editor,
543 *EUROCRYPT’88*, volume 330 of *LNCS*, pages 123–128. Springer, Heidelberg, May 1988.
544 doi:10.1007/3-540-45961-8_11.
- 545 **38** Arjen K. Lenstra, James P. Hughes, Maxime Augier, Joppe W. Bos, Thorsten Kleinjung,
546 and Christophe Wachter. Public keys. In Reihaneh Safavi-Naini and Ran Canetti, editors,
547 *CRYPTO 2012*, volume 7417 of *LNCS*, pages 626–642. Springer, Heidelberg, August 2012.
548 doi:10.1007/978-3-642-32009-5_37.
- 549 **39** Ueli M. Maurer. Unifying zero-knowledge proofs of knowledge. In Bart Preneel, editor,
550 *AFRICACRYPT 09*, volume 5580 of *LNCS*, pages 272–286. Springer, Heidelberg, June 2009.
- 551 **40** Ilya Mironov and Noah Stephens-Davidowitz. Cryptographic reverse firewalls. In Elisabeth
552 Oswald and Marc Fischlin, editors, *EUROCRYPT 2015, Part II*, volume 9057 of *LNCS*, pages
553 657–686. Springer, Heidelberg, April 2015. doi:10.1007/978-3-662-46803-6_22.
- 554 **41** Tatsuaki Okamoto. Provably secure and practical identification schemes and corresponding
555 signature schemes. In Ernest F. Brickell, editor, *CRYPTO’92*, volume 740 of *LNCS*, pages
556 31–53. Springer, Heidelberg, August 1993. doi:10.1007/3-540-48071-4_3.
- 557 **42** Tatsuaki Okamoto and Kazuo Ohta. Divertible zero knowledge interactive proofs and com-
558 mutative random self-reducibility. In Jean-Jacques Quisquater and Joos Vandewalle, editors,
559 *EUROCRYPT’89*, volume 434 of *LNCS*, pages 134–148. Springer, Heidelberg, April 1990.
560 doi:10.1007/3-540-46885-4_16.
- 561 **43** Rafail Ostrovsky, Vanishree Rao, and Ivan Visconti. On selective-opening attacks against en-
562 cryption schemes. In Michel Abdalla and Roberto De Prisco, editors, *SCN 14*, volume
563 8642 of *LNCS*, pages 578–597. Springer, Heidelberg, September 2014. doi:10.1007/
564 978-3-319-10879-7_33.
- 565 **44** Torben P. Pedersen. Non-interactive and information-theoretic secure verifiable secret sharing.
566 In Joan Feigenbaum, editor, *CRYPTO’91*, volume 576 of *LNCS*, pages 129–140. Springer,
567 Heidelberg, August 1992. doi:10.1007/3-540-46766-1_9.
- 568 **45** Nicole Perlroth, Jeff Larson, and Scott Shane. N.S.A. able to foil basic safeguards of privacy
569 on web. *The New York Times*, September 2013.
- 570 **46** Alexander Russell, Qiang Tang, Moti Yung, and Hong-Sheng Zhou. Cliptography: Clipping
571 the power of kleptographic attacks. In Jung Hee Cheon and Tsuyoshi Takagi, editors,
572 *ASIACRYPT 2016, Part II*, volume 10032 of *LNCS*, pages 34–64. Springer, Heidelberg,
573 December 2016. doi:10.1007/978-3-662-53890-6_2.
- 574 **47** Alexander Russell, Qiang Tang, Moti Yung, and Hong-Sheng Zhou. Generic semantic security
575 against a kleptographic adversary. In Bhavani M. Thuraisingham, David Evans, Tal Malkin,
576 and Dongyan Xu, editors, *ACM CCS 2017*, pages 907–922. ACM Press, October / November
577 2017. doi:10.1145/3133956.3133993.
- 578 **48** Alessandra Scafuro and Ivan Visconti. On round-optimal zero knowledge in the bare public-key
579 model. In David Pointcheval and Thomas Johansson, editors, *EUROCRYPT 2012*, volume 7237

- 580 of *LNCS*, pages 153–171. Springer, Heidelberg, April 2012. doi:[10.1007/978-3-642-29011-4_](https://doi.org/10.1007/978-3-642-29011-4_11)
581 [11](https://doi.org/10.1007/978-3-642-29011-4_11).
- 582 **49** Claus-Peter Schnorr. Efficient identification and signatures for smart cards. In Gilles Brassard,
583 editor, *CRYPTO'89*, volume 435 of *LNCS*, pages 239–252. Springer, Heidelberg, August 1990.
584 doi:[10.1007/0-387-34805-0_22](https://doi.org/10.1007/0-387-34805-0_22).
- 585 **50** Gustavus J. Simmons. The prisoners' problem and the subliminal channel. In David Chaum,
586 editor, *CRYPTO'83*, pages 51–67. Plenum Press, New York, USA, 1983.
- 587 **51** Dominique Unruh. Quantum proofs of knowledge. In David Pointcheval and Thomas Johansson,
588 editors, *EUROCRYPT 2012*, volume 7237 of *LNCS*, pages 135–152. Springer, Heidelberg,
589 April 2012. doi:[10.1007/978-3-642-29011-4_10](https://doi.org/10.1007/978-3-642-29011-4_10).
- 590 **52** Adam Young and Moti Yung. Kleptography: Using cryptography against cryptography.
591 In Walter Fumy, editor, *EUROCRYPT'97*, volume 1233 of *LNCS*, pages 62–74. Springer,
592 Heidelberg, May 1997. doi:[10.1007/3-540-69053-0_6](https://doi.org/10.1007/3-540-69053-0_6).