

# Hardness vs. (Very Little) Structure in Cryptography: A Multi-Prover Interactive Proofs Perspective

Gil Segev\*

Ido Shahaf\*†

## Abstract

The hardness of highly-structured computational problems gives rise to a variety of public-key primitives. On one hand, the structure exhibited by such problems underlies the basic functionality of public-key primitives, but on the other hand it may endanger public-key cryptography in its entirety via potential algorithmic advances. This subtle interplay initiated a fundamental line of research on whether structure is inherently necessary for cryptography, starting with Rudich’s early work (PhD Thesis ’88) and recently leading to that of Bitansky, Degwekar and Vaikuntanathan (CRYPTO ’17).

Identifying the structure of computational problems with their corresponding complexity classes, Bitansky et al. proved that a variety of public-key primitives (e.g., public-key encryption, oblivious transfer and even functional encryption) cannot be used in a black-box manner to construct either any hard language that has NP-verifiers both for the language itself and for its complement, or any hard language (and even promise problem) that has a statistical zero-knowledge proof system – corresponding to hardness in the structured classes  $\text{NP} \cap \text{coNP}$  or SZK, respectively, from a black-box perspective.

In this work we prove that the same variety of public-key primitives do not inherently require *even very little structure* in a black-box manner: We prove that they do not imply any hard language that has *multi-prover interactive proof systems* both for the language and for its complement – corresponding to hardness in the class  $\text{MIP} \cap \text{coMIP}$  from a black-box perspective. Conceptually, given that  $\text{MIP} = \text{NEXP}$ , our result rules out languages with very little structure. Additionally, we prove a similar result for collision-resistant hash functions, and more generally for any cryptographic primitive that exists relative to a random oracle.

Already the cases of languages that have IP or AM proof systems both for the language itself and for its complement, which we rule out as immediate corollaries, lead to intriguing insights. For the case of IP, where our result can be circumvented using non-black-box techniques, we reveal a gap between black-box and non-black-box techniques. For the case of AM, where circumventing our result via non-black-box techniques would be a major development, we both strengthen and unify the proofs of Bitansky et al. for languages that have NP-verifiers both for the language itself and for its complement and for languages that have a statistical zero-knowledge proof system.

---

\*School of Computer Science and Engineering, Hebrew University of Jerusalem, Jerusalem 91904, Israel. Email: {segev,ido.shahaf}@cs.huji.ac.il. Supported by the European Union’s Horizon 2020 Framework Program (H2020) via an ERC Grant (Grant No. 714253).

†Supported by the Clore Israel Foundation via the Clore Scholars Programme.

**Contents**

- 1 Introduction** **1**
- 1.1 Our Contributions . . . . . 2
- 1.2 Overview of Our Approach . . . . . 3
- 1.3 Paper Organization . . . . . 8
  
- 2 Preliminaries** **8**
  
- 3 The Classes of Constructions** **10**
  
- 4 Impossibility Result for Constructions based on  $i\mathcal{O}$  and Injective OWEs** **12**
- 4.1 Our Generalized Decision Oracle . . . . . 13
- 4.2 Our Indistinguishability Obfuscation Oracle . . . . . 15
- 4.3 The Existence of an Injective One-Way Function . . . . . 16
- 4.4 The Existence of an Indistinguishability Obfuscator . . . . . 19
- 4.5 Putting it All Together . . . . . 24
  
- 5 Impossibility Result for Constructions based on RO-Implied Primitives** **25**
- 5.1 Deciding (MIP, coMIP) Protocol Pair Languages . . . . . 26
- 5.2 Putting it All Together . . . . . 31
  
- References** **31**

## 1 Introduction

Starting with the revolutionary invention of public-key cryptography [DH76, RSA78, GM84], the hardness of highly-structured computational problems (e.g., factoring, discrete log, or various lattice-based problems) has given rise to a variety of public-key primitives. On one hand, the structure exhibited by such problems underlies the basic functionality of nearly all such primitives, but on the other hand it may also danger their conjectured hardness. As noted by Barak [Bar13], this “*makes public-key cryptography somewhat of an endangered species that could be wiped out by a surprising algorithmic advance*”.

This subtle interplay has led to the long-studied question of whether structure is inherently necessary for certain cryptographic primitives, and most notably for public-key primitives. While there may be different approaches for measuring or quantifying “structure”, the main approach taken by the cryptography community over the years relies on computational complexity: Understanding which cryptographic primitives inherently require hardness in “structured” complexity classes such as  $\text{NP} \cap \text{coNP}$ ,  $\text{TFNP}$  and  $\text{SZK}$ .

There are only a few known examples of cryptographic primitives that require hardness in such classes. Most notably, one-way permutations imply hardness in  $\text{NP} \cap \text{coNP}$  [Bra79], homomorphic encryption and non-interactive computational private-information retrieval imply hardness in  $\text{SZK}$  [BL13, LV16], and indistinguishability obfuscation implies hardness in  $\text{PPAD} \subseteq \text{TFNP}$  unless  $\text{NP} \subseteq \text{ioBPP}$  [BPR15, GPS16, KMN<sup>+</sup>14].

Within the classic framework of black-box constructions, capturing “natural” cryptographic constructions [IR89, RTV04], Rudich [Rud88] showed (based on [BI87, HH87]) that a one-way function cannot be used in black-box manner to construct  $\text{NP}$ -verifiers for any hard language both for the language itself and for its complement – corresponding to hardness in  $\text{NP} \cap \text{coNP}$  from a black-box perspective (we note that the known examples stated above all follow in such a black-box manner).

For several decades no progress has been made in extending Rudich’s result to public-key primitives or to other complexity classes. This situation has recently changed dramatically with the work of Bitansky, Degwekar and Vaikuntanathan [BDV17] (see also the refinements in the more recent work of Bitansky and Degwekar [BD19]): They showed that even indistinguishability obfuscation cannot be used in a black-box manner to construct any hard language that has  $\text{NP}$  verifiers both for the language itself and for its complement, or any hard language (and even a promise problem) that has a statistical zero-knowledge proof system – corresponding to hardness in  $\text{NP} \cap \text{coNP}$  or  $\text{SZK}$ , respectively, from a black-box perspective. Proving their result within the framework of Asharov and Segev [AS15, AS16] capturing indistinguishability obfuscation for oracle-aided computations, Bitansky et al. in fact proved their result for all primitives that can be based on indistinguishability obfuscation for circuits that access an injective one-way function in a black-box manner. These include, in particular, a variety of public-key primitives including public-key encryption, oblivious transfer and even functional encryption.

Focusing on the classes  $\text{NP} \cap \text{coNP}$  and  $\text{SZK}$ , Bitansky et al. showed that, from a black-box perspective, public-key cryptography does not inherently require highly-structured hardness. However, going back to Barak’s concern [Bar13], even less stringent forms of structure may still endanger public-key cryptography in its entirety. This leads to the following fundamental question aiming at substantially refining our understanding of the interplay between hardness and structure:

Does public-key cryptography inherently require hardness in complexity classes that are “less structured” than  $\text{NP} \cap \text{coNP}$  or  $\text{SZK}$ ?

## 1.1 Our Contributions

In this work we show that a wide variety of public-key primitives do not inherently require even very little structure in a black-box manner. Specifically, we prove that such primitives do not naturally imply hard languages that have *multi-prover interactive proof systems* (MIP) [BGK<sup>+</sup>88] both for the language and for its complement.

Conceptually, given that  $\text{MIP} = \text{NEXP}$  [BFL91], our result considers languages with very little structure. Already the cases of languages that have  $\text{IP}$  or  $\text{AM}$  proof systems both for the language itself and for its complement, which we obtain as immediate corollaries, lead to intriguing insights. For the case of  $\text{IP}$ , where our result can be circumvented using non-black-box techniques, we reveal a gap between black-box and non-black-box techniques (as we discuss below). For the case of  $\text{AM}$ , where circumventing our result via non-black-box techniques would be a major development, we both strengthen and unify the proofs of Bitansky et al. for languages that have  $\text{NP}$ -verifiers both for the language itself and for its complement and for languages that have a statistical zero-knowledge proof system (since  $\text{NP} \subseteq \text{AM}$  by definition, and since  $\text{SZK} \subseteq \text{AM} \cap \text{coAM}$  in a black-box manner [For89, AH91]).<sup>1</sup>

The following theorem is an informal statement of our main result. We refer the reader to Section 1.2 for an overview of our result, and to Sections 3 and 4 for a formal definition of the class of constructions to which our result applies and for a formal theorem statement, respectively.

**Theorem 1.1** (Informal). *There is no fully black-box construction of a pair of multi-prover interactive proof systems,  $\Pi$  and  $\overline{\Pi}$ , corresponding to a worst-case hard language  $L$  and to its complement  $\overline{L}$ , respectively, from an injective one-way function  $f$  and an indistinguishability obfuscator for the class of all oracle-aided circuits  $C^f$ .*

Note that as our result rules out constructions of languages that are *worst-case hard*, then it rules out in particular constructions of languages that are *average-case hard*.

**Black-box vs. non-black-box constructions.** Our result might seem too strong and somewhat contradicting to the fact that any one-way function implies a hard (even on average) language in  $\text{NP} \subseteq \text{IP}$  in a *black-box manner*. Given that  $\text{IP}$  is closed under complement [LFK<sup>+</sup>92, Sha90], then

$$\text{NP} \subseteq \text{IP} \cap \text{coIP} \subseteq \text{MIP} \cap \text{coMIP}.$$

In particular, any one-way function implies a hard language that has  $\text{IP}$  proof systems both for the language itself and for its complement, which seemingly contradicts our result. However, this sequence of containments cannot be established via relativizing reductions, and thus there is in fact no contradiction (note that any black-box reduction relativizes [RTV04]), but rather a gap between black-box and non-black-box techniques. Specifically, Chang et al. [CCG<sup>+</sup>94] showed that there exists an oracle  $\Gamma$  relative to which  $\text{NP}^\Gamma \not\subseteq \text{coIP}^\Gamma$ , and in particular  $\text{IP}$  is *not* closed under complement with respect to relativizing reductions. Still, as mentioned above, our impossibility result already applies to  $\text{AM} \cap \text{coAM}$ , for which circumventing our result via non-black-box techniques would be a major development. We discuss this in much more detail in Section 1.2 in the context of black-box representations of complexity classes.

---

<sup>1</sup>We note that the result of Bitansky et al. for SZK holds not only for languages but in fact also for promise problems. This, however, cannot be covered by our result since already a hard promise problem that has  $\text{NP}$  verifiers both for its “YES” instances and for its “NO” instances can be constructed based on any one-way function in a black-box manner.

**Implications to public-key cryptography.** Similarly to Bitansky et al. [BDV17] we prove our result within the framework of Asharov and Segev [AS15, AS16], capturing indistinguishability obfuscation for oracle-aided circuits. Indistinguishability obfuscation for such circuits suffices for realizing a variety of public-key primitives (e.g., public-key encryption, oblivious transfer and even functional encryption) in a fully black-box manner [SW14, Wat15, AS15], and therefore as a corollary we obtain that there is no construction of the above form based on any of these primitives.

We strongly emphasize that our result is unconditional, and in particular does not depend on whether or not indistinguishability obfuscation actually exists. Even if it does not exist in the actual world, then within the framework of Asharov and Segev it does exist information theoretically, and it implies the above variety of public-key primitives to which our result applies (once again, in an unconditional manner).

**Collision-resistant hash functions.** As mentioned above, Theorem 1.1 rules out black-box constructions that are based on any cryptographic primitive which can be constructed based on injective one-way functions and indistinguishability obfuscation in a fully black-box manner. Theorem 1.1 does not rule out black-box constructions that are based on other primitives. One such example is domain-invariant one-way permutations, which are not implied by injective one-way functions and indistinguishability obfuscation in a fully black-box manner [AS16], but do imply multi-prover interactive proof systems corresponding to a hard language and to its complement.

An additional such example is collision-resistant hash functions, which are also not implied by injective one-way functions and indistinguishability obfuscation in a fully black-box manner [AS15]. Our second result shows that unlike one-way permutations, collision-resistant hash functions do not imply such proof systems. The following theorem is an informal statement of our result. We refer the reader to Sections 3 and 5 for a formal definition of the class of constructions to which our result applies and for a formal theorem statement, respectively.

**Theorem 1.2** (Informal). *There is no fully black-box construction of a pair of multi-prover interactive proof systems,  $\Pi$  and  $\bar{\Pi}$ , corresponding to a worst-case hard language  $L$  and to its complement  $\bar{L}$ , respectively, from a collision-resistant hash function  $f$ .*

Bitansky and Degwekar [BD19] already showed that it is impossible to construct a statistical zero-knowledge proof system corresponding to a worst-case hard problem from a collision-resistant hash function in a black-box way, so our result strengthens their result.<sup>2</sup> We note that, for simplicity, our result is stated with respect to collision-resistant hash functions. However, our techniques in fact apply to any cryptographic primitive that exists relative to a random oracle (e.g., multi-input correlation-intractable hash functions [CGH04]).

The techniques used for proving Theorem 1.2 rely on the approaches of Blum and Impagliazzo [BI87], Impagliazzo and Naor [IN88] and Nisan [Nis89, Nis91], and are significantly different from those used for proving Theorem 1.1 (which are overviewed in Section 1.2) although they share the same overall framework. Specifically, we demonstrate that the fundamental notions of block sensitivity and certificate complexity can be adapted to our setting for ruling out black-box constructions.

## 1.2 Overview of Our Approach

In this section we provide a high-level overview of the framework in which we prove our impossibility results, and then describe the main ideas and challenges underlying our proof of Theorem 1.1.

---

<sup>2</sup>However, their result holds not only for languages but in fact also for promise problems, while as mentioned above, our result cannot cover promise problems.

**Black-box constructions.** Our goal is to prove a statement along the lines of “a cryptographic primitive  $\mathcal{P}$  does not naturally imply a hard language in a complexity class  $\mathcal{C}$ ”. However, it is not clear how to prove such a statement in an unconditional manner, as it may be the case that the class  $\mathcal{C}$  (e.g.,  $\text{NP} \cap \text{coNP}$  as discussed by Bitansky et al. [BDV17]) does not contain hard languages. One possible approach is to prove a result that is conditioned on a specific assumption, but then it may be the case that the assumption itself already rules out the existence of hard languages in the class  $\mathcal{C}$ . Obtaining substantial insight using such an approach requires a deep understanding of the interplay between the primitive  $\mathcal{P}$ , the complexity class  $\mathcal{C}$  and the additional assumption – which is somewhat rare when considering cryptographic primitives and assumptions.

Faced with such difficulties, the cryptography community has relied over the years on the framework of black-box constructions [IR89, RTV04] for proving impossibility results for “natural” construction techniques. In our context, a fully black-box construction of a hard language  $L \in \mathcal{C}$  based on a cryptographic primitive  $\mathcal{P}$  consists of two ingredients. The first ingredient is a “construction” of a language  $L^{\mathcal{P}}$  that completely ignores the internal implementation of  $\mathcal{P}$  and only requires black-box access to any given implementation of  $\mathcal{P}$ . Here, the notion of a “construction” depends on the specific complexity class  $\mathcal{C}$ . For example, in a natural black-box interpretation of  $\text{NP} \cap \text{coNP}$ , Rudich [Rud88] and Bitansky et al. [BDV17] considered as a construction a pair of oracle-aided  $\text{NP}$ -verifiers,  $V$  and  $\bar{V}$ , for the language itself and for its complement, respectively, where the verifiers have black-box access to the primitive  $\mathcal{P}$ . That is, for any oracle realizing  $\mathcal{P}$ , the two verifiers must be valid in the sense that for any instance  $x \in \{0,1\}^*$  either there exists a “yes” witness for  $V^{\mathcal{P}}$  and there do not exist any “no” witnesses for  $\bar{V}^{\mathcal{P}}$  (i.e.,  $x \in L^{\mathcal{P}}$ ), or there exists a “no” witness for  $\bar{V}^{\mathcal{P}}$  and there do not exist any “yes” witnesses for  $V^{\mathcal{P}}$  (i.e.,  $x \notin L^{\mathcal{P}}$ ). The second ingredient, is a black-box proof of hardness, showing that for any implementation of the primitive  $\mathcal{P}$ , any algorithm that decides the language  $L^{\mathcal{P}}$  can be efficiently used in a black-box manner for breaking the security of the given implementation of  $\mathcal{P}$ .

At this point we would like to already emphasize that a “black-box representation” of a complexity class is in fact not unique, and that different representations are not always equivalent from a black-box perspective. For example, a natural black-box representation for the class  $\text{IP} \cap \text{coIP}$  relative to a given primitive  $\mathcal{P}$  is to consider all languages that have interactive proof systems both for the language itself and for its complement, where the two proof systems access  $\mathcal{P}$  in a black-box manner. However, since  $\text{IP}$  is closed under complement [LFK<sup>+</sup>92, Sha90] then  $\text{IP} \cap \text{coIP} = \text{IP}$  and therefore an additional representation is to consider all languages that have interactive proof systems for the language itself (without considering its complement) where the proof system accesses  $\mathcal{P}$  in a black-box manner. As discussed in Section 1.1, these two representations are not equivalent from a black-box perspective since  $\text{IP}$  is not closed under complement with respect to relativizing reductions.

**The structure of our proof of Theorem 1.1.** Following Bitansky et al. [BDV17] we prove our result within the framework of Asharov and Segev [AS15] for capturing black-box constructions based on indistinguishability obfuscation, utilizing the latter as a “central hub” for deriving impossibility results for a variety of public-key primitives. As observed by Asharov and Segev, although constructions that are based on indistinguishability obfuscation are almost always *non-black-box*, most of their non-black-box techniques have essentially the same flavor: The obfuscator itself is used in a black-box manner and applied to circuits that can be constructed in a fully black-box manner from a low-level primitive, such as a one-way function. Thus, even though the obfuscator requires concrete implementations of such circuits, by introducing the stronger primitive of an indistinguishability obfuscator for *oracle-aided* circuits (see Section 2), Asharov and Segev showed that such non-black-box techniques in fact directly translate into black-box ones. These include, in particular, non-black-box

techniques such as the punctured programming approach of Sahai and Waters [SW14] and Waters [Wat15] leading to the construction of a variety of public-key primitive. Relying on the transitivity of black-box reductions, this enables to rule out black-box constructions based on all of these primitives by focusing only on indistinguishability obfuscation for oracle-aided circuits and one-way functions.

In order to prove our impossibility result within this framework, we present a distribution over oracles  $\Gamma$  relative to which we prove the following two properties:

- Relative to a random instance of  $\Gamma$  there exist an injective one-way function  $f$  and an indistinguishability obfuscator  $i\mathcal{O}$  for the class of all oracle-aided circuits  $C^f$ .
- Relative to any instance of  $\Gamma$ , we can efficiently decide in the worst case any language that has multi-prover interactive proof systems,  $\Pi^{f,i\mathcal{O}}$  and  $\bar{\Pi}^{f,i\mathcal{O}}$ , for the language itself and for its complement, respectively.<sup>3</sup>

Our oracle  $\Gamma$  is a pair of the form  $(\Psi, \text{Decide}^\Psi)$ , where  $\Psi$  is based on the oracle of Asharov and Segev that realizes a one-way function and an indistinguishability obfuscator, and  $\text{Decide}^\Psi$  is a generalization of the “decision oracle” introduced by Bitansky et al. for deciding languages that rely on  $\Psi$  in a black-box manner (more specifically, whose black-box representation as discussed above relies on  $\Psi$ ). In the work of Bitansky et al. the decision oracle is defined in a manner that allows to easily decide any language  $L^\Psi$  that has NP-verifiers,  $V^\Psi$  and  $\bar{V}^\Psi$ , for the language itself and for its complement, and the main technical challenge underlying their work is proving that  $\Psi$  realizes a one-way function and an indistinguishability obfuscator relative to the decision oracle.

Our decision oracle is a natural generalization that allows to easily decide any language  $L^\Psi$  that has multi-prover proof systems,  $\Pi^\Psi$  and  $\bar{\Pi}^\Psi$ , for the language itself and for its complement. This decision oracle seems much more powerful than that of Bitansky et al. as it decides a significantly larger class of languages, and our technical effort is devoted to proving that the oracle  $\Psi$  still realizes a one-way function and an indistinguishability obfuscator even relative to our generalized decision oracle.

In what follows we describe the decision oracle of Bitansky et al. (to which we refer as the BDV decision oracle) and discuss its key property that underlies their approach. Then, we describe our generalized oracle, relative to which this key property no longer seems to hold, and then describe our the main ideas underlying our proof.

**The BDV decision oracle.** For any oracle  $\Psi$ , taken from an appropriate family  $\mathfrak{S}$  of oracles, the BDV decision oracle  $\text{Decide}_{\mathfrak{S}}^\Psi$  takes as input a triplet  $(V, \bar{V}, x)$ , where  $V$  and  $\bar{V}$  are oracle-aided circuits. The oracle first checks whether or not the pair  $(V, \bar{V})$  indeed consists of valid NP-verifiers for a language and for its complement in the standard black-box sense discussed above. That is, checks whether or not for any  $\Psi' \in \mathfrak{S}$  and  $x' \in \{0, 1\}^n$  exactly one out of the following two cases holds: (1) There exists a “yes” witness  $w'$  such that  $V^{\Psi'}(x', w') = 1$  and there do not exist any “no” witnesses  $w'$  such that  $\bar{V}^{\Psi'}(x', w') = 1$ ; (2) there exists a “no” witness  $w'$  such that  $\bar{V}^{\Psi'}(x', w') = 1$  and there do not exist any “yes” witnesses  $w'$  such that  $V^{\Psi'}(x', w') = 1$  (note that the witnesses are allowed to depend on  $\Psi'$ ). If  $(V, \bar{V})$  is not valid in this sense, then the oracle outputs  $\perp$ . If  $(V, \bar{V})$  is valid, then the oracle outputs 1 if  $x \in L^\Psi$  and 0 otherwise, where  $L^\Psi$  is the language defined by  $(V^\Psi, \bar{V}^\Psi)$ .

---

<sup>3</sup>In fact, as discussed below we allow the honest provers to depend on the one-way function and the obfuscator in an arbitrary non-black-box manner, and only require that the verifiers are constructed in a black-box manner (this makes our result stronger).

Then, any language that has oracle-aided NP-verifiers both for the language itself and for its complement with respect to any  $\Psi \in \mathfrak{S}$ , can be easily decided in the worst case by an algorithm that issues a single query to the BDV decision oracle. The main challenge in the work of Bitansky et al. was in showing that a random instance of  $\Psi$  that is sampled from the family  $\mathfrak{S}$  of oracles introduced by Asharov and Segev (or from any other appropriate family) realizes a one-way function and an indistinguishability obfuscator even relative to  $\text{Decide}_{\mathfrak{S}}^{\Psi}$ .

**The existence of small critical sets.** The key property underlying the proof of Bitansky et al. is the following observation on the existence of “small critical sets”. Fix an oracle  $\Psi \in \mathfrak{S}$  and let  $(V, \bar{V}, x)$  be a query to their decision oracle such that the pair  $(V, \bar{V})$  is valid in the above sense, and  $V$  and  $\bar{V}$  issue at most  $q$  oracle queries. Then, there exists a “critical set” of at most  $q$  queries, such that for any oracle  $\Psi' \in \mathfrak{S}$  that agrees with  $\Psi$  on the outputs of all queries from the critical set it holds that  $\text{Decide}_{\mathfrak{S}}^{\Psi}(V, \bar{V}, x) = \text{Decide}_{\mathfrak{S}}^{\Psi'}(V, \bar{V}, x)$ .

The existence of such a small critical set follows from the  $\text{NP} \cap \text{coNP}$  structure of the pair  $(V, \bar{V})$ . Specifically, assume without loss of generality that  $x \in L^{\Psi}$ , and let  $w$  be a witness such that  $V^{\Psi}(x, w) = 1$ . Define the set of critical queries as all  $\Psi$ -queries that are issued in the computation  $V^{\Psi}(x, w)$ , and let  $\Psi'$  be any oracle that agrees with  $\Psi$  on this set. Then clearly  $V^{\Psi'}(x, w) = V^{\Psi}(x, w) = 1$ , and the validity of the pair  $(V, \bar{V})$  guarantees that there is no witness  $\tilde{w}$  such that  $\bar{V}^{\Psi'}(x, \tilde{w}) = 1$ . Thus,  $\text{Decide}_{\mathfrak{S}}^{\Psi}(V, \bar{V}, x) = \text{Decide}_{\mathfrak{S}}^{\Psi'}(V, \bar{V}, x) = 1$ .

Relying on this key property, Bitansky et al. proved that  $\Psi$  realizes a one-way function and an indistinguishability obfuscator relative to their decision oracle via an elegant sequence of hybrids in each case. Specifically, in each sequence the first experiment is the actual security experiment of the one-way function or the indistinguishability obfuscator, the last experiment is one in which no algorithm can achieve any advantage, and the transition between each consecutive pair of experiment is enabled by this key property (or via standard arguments).

**Representing  $\text{MIP} \cap \text{coMIP}$  in a black-box manner.** In order to describe our approach, we first need to describe our black-box representation of languages in the complexity class  $\text{MIP} \cap \text{coMIP}$ . Naturally generalizing the approach of Rudich and Bitansky et al. for  $\text{NP} \cap \text{coNP}$ , we consider pairs of polynomial-time oracle-aided MIP-verifiers,  $V$  and  $\bar{V}$ , for the language itself and for its complement, respectively, subject to a similar validity requirement of their black-box flavor: For any oracle  $\Psi$  taken from an appropriate family  $\mathfrak{S}$  of oracles, there should exist a language  $L^{\Psi}$  such that the following two conditions are satisfied<sup>4</sup>:

- For every  $x \in L^{\Psi}$  there exist computationally-unbounded provers  $P_1, \dots, P_N$  such that<sup>5</sup>

$$\Pr_{r \leftarrow \{0,1\}^{\text{poly}(|x|)}} [\langle V^{\Psi}(x; r), P_1, \dots, P_N \rangle = 1] \geq 2/3,$$

and for every computationally-unbounded provers  $\bar{P}_1, \dots, \bar{P}_N$  it holds that

$$\Pr_{r \leftarrow \{0,1\}^{\text{poly}(|x|)}} [\langle \bar{V}^{\Psi}(x; r), \bar{P}_1, \dots, \bar{P}_N \rangle = 1] \leq 1/3.$$

---

<sup>4</sup>For an oracle  $\Psi$ , an instance  $x$ , a string  $r$ , a polynomial-time oracle-aided verifier  $V$ , and provers  $P_1, \dots, P_N$  we denote by  $\langle V^{\Psi}(x; r), P_1, \dots, P_N \rangle$  the output of  $V$  with oracle access to  $\Psi$  on input  $x$  and randomness  $r$  in the multi-prover execution with  $P_1, \dots, P_N$ . Note that whenever the provers are computationally unbounded we can assume that they are deterministic.

<sup>5</sup>It is usually assumed that the same provers are used for every  $x \in \{0,1\}^n$ , and that they obtain  $x$  as input. However, since the provers are computationally unbounded, our definition is clearly equivalent and easier to work with for our purposes.

- For every  $x \notin L^\Gamma$  there exist computationally-unbounded provers  $\bar{P}_1, \dots, \bar{P}_N$  such that

$$\Pr_{r \leftarrow \{0,1\}^{\text{poly}(|x|)}} \left[ \langle \bar{V}^\Psi(x; r), \bar{P}_1, \dots, \bar{P}_N \rangle = 1 \right] \geq 2/3,$$

and for every computationally-unbounded provers  $P_1, \dots, P_N$  it holds that

$$\Pr_{r \leftarrow \{0,1\}^{\text{poly}(|x|)}} \left[ \langle V^\Psi(x; r), P_1, \dots, P_N \rangle = 1 \right] \leq 1/3.$$

Note that instead of considering oracle-aided MIP proof systems we consider oracle-aided MIP verifiers, and allow the honest provers to depend on any given oracle in an arbitrary non-black-box manner (thus our result rules out, in particular, oracle-aided proof systems). We refer the reader to Section 3 where we formally describe the proof systems we consider and the class of constructions to which our result applies.

**Our generalized decision oracle.** For any oracle  $\Psi \in \mathfrak{S}$  our generalized decision oracle  $\text{Decide}_{\mathfrak{S}}^\Psi$  takes as input a triplet  $(V, \bar{V}, x)$ , where  $V$  and  $\bar{V}$  are oracle-aided MIP-verifiers and  $x \in \{0,1\}^n$ . The oracle first checks whether or not the pair  $(V, \bar{V})$  indeed consists of MIP-verifiers for a language and for its complement with respect to all oracles in  $\mathfrak{S}$  as discussed above. If  $(V, \bar{V})$  is not valid in this sense, then the oracle outputs  $\perp$ . If  $(V, \bar{V})$  is valid, then the oracle outputs 1 if  $x \in L^\Psi$  and 0 otherwise, where  $L^\Psi$  is the language defined by  $(V^\Psi, \bar{V}^\Psi)$ .

At this point, we would ideally like to follow the approach of Bitansky et al. in proving that  $\Psi$  realizes a one-way function and an indistinguishability obfuscator relative to our generalized decision oracle. Recall that their proof consists of a sequence of hybrid experiments, where the transition between each consecutive pair of experiments is enabled by the existence of a small set of critical queries. Specifically, in each transition they modify  $\Psi$  on some set of queries into an oracle  $\Psi'$ , and argue that unless these queries fall into the small critical set then the decision oracle behaves exactly the same.

**Are there small and useful critical query sets?** Fix an oracle  $\Psi \in \mathfrak{S}$ , and fix a query  $(V, \bar{V}, x)$  to our generalized decision oracle, where  $V$  and  $\bar{V}$  are valid MIP-verifiers in the above sense. Unlike the case of NP-verifiers, when considering MIP-verifiers then at a first glance there does not seem to be a small set of queries that completely determines whether or not  $x \in L^\Psi$ . Specifically, assuming for the current discussion that  $x \in L^\Psi$ , in the case of NP-verifiers this is completely determined by the polynomial number of queries to the oracle  $\Psi$  in the execution  $V^\Psi(x, w)$  where  $w$  is any specific witness (say, the lexicographically first such witness). However, in the case of MIP-verifiers, we are guaranteed that there exist provers  $P_1, \dots, P_N$  that lead the MIP-verifier  $V^\Psi(x; r)$  to accept with probability at least  $2/3$  over the randomness  $r \leftarrow \{0,1\}^{\text{poly}(|x|)}$  of the verifier – but this guarantee involves potentially exponentially-many executions and thus exponentially-many queries to the oracle  $\Psi$ . It may even be the case that any oracle  $\Psi'$  that agrees with  $\Psi$  on all of these queries, is in fact  $\Psi' = \Psi$ , and this is not very useful for the purpose of transitioning between two hybrid experiments.

Nevertheless, let us consider an oracle  $\Psi'$  that differs from  $\Psi$  on a *single* query  $z$ , and now suppose that suddenly  $x \notin L^{\Psi'}$  although we started with  $x \in L^\Psi$ . Thus, no provers can now lead  $V^{\Psi'}(x; r)$  to accept with probability larger than  $1/3$  over the randomness  $r \leftarrow \{0,1\}^{\text{poly}(|x|)}$ , and in particular this holds for the above provers  $P_1, \dots, P_N$  that led  $V^\Psi(x; r)$  to accept with probability at least  $2/3$ . The only way that  $V^{\Psi'}(x; r)$  can differ from  $V^\Psi(x; r)$  in an execution with the same  $P_1, \dots, P_N$  is by having  $V^\Psi(x; r)$  query  $\Psi$  on  $z$  – and we can deduce that with probability at least  $1/3$  over the choice of  $r \leftarrow \{0,1\}^{\text{poly}(|x|)}$  it holds that  $V^\Psi(x; r)$  queries  $\Psi$  on  $z$  when interacting with  $P_1, \dots, P_N$ .

Therefore, it is quite tempting to fix a distance parameter  $d \geq 1$ , and then for an oracle  $\Psi \in \mathfrak{S}$  and a query  $(V, \bar{V}, x)$  such that  $x \in L^\Psi$  to define the following “ $d$ -influential set” of queries: Let  $P_1, \dots, P_N$  be provers that lead  $V^\Psi(x; r)$  to accept with probability at least  $2/3$ , then the  $d$ -influential set consists of all queries that  $V^\Psi(x; r)$  issues to  $\Psi$  in at least a  $1/(3d)$ -fraction of these executions. Then, if  $V$  issues at most  $q$  queries in each execution, then this set consists of at most  $3qd$  queries. Moreover, for any oracle  $\Psi'$  that differs from  $\Psi$  on at most  $d$  queries, and these queries are not in the  $d$ -influential set, then it must hold that  $x \in L^{\Psi'}$  (the probability that  $V(x; r)$  accepts cannot drop from  $2/3$  to  $1/3$  when switching from  $\Psi$  to  $\Psi'$  since they differ on at most  $d$  queries and each of these queries affects less than a  $1/(3d)$ -fraction of the executions).

**From influential queries to influential labels.** Unfortunately, this observation is still insufficient for our purposes. In the proof of Bitansky et al. the number of differences between  $\Psi$  and  $\Psi'$  is irrelevant as long as these differences are not in the critical set. However, in our case more than  $d$  differences outside of the  $d$ -influential set may still cause the verifier’s acceptance probability to drop from  $2/3$  to below  $1/3$ .

Although our proof considers oracles  $\Psi$  and  $\Psi'$  that may differ on an exponential number of queries, we tailor the specific structure of our obfuscator in a way that enables us to “group together” related queries: We introduce labeling functions (depending on the specific structure of our oracles) that assign a label to each query to the oracle  $\Psi$ , where different queries may share the same label. We show that it now suffices to focus on the small number  $d \leq 3$  of labels that result from the potentially-exponential number of differences between the oracles  $\Psi$  and  $\Psi'$ .

Specifically, we prove that for any  $\Psi$  and for any query  $(V, \bar{V}, x)$  to our generalized decision oracle there exists a small set  $\mathbf{I}$  of “ $d$ -influential labels” such that any changes to  $\Psi$  involving at most  $d$  labels outside of  $\mathbf{I}$  do not change the answer to the query. That is, let  $\Psi' \in \mathfrak{S}$  be any oracle for which there exists a set  $\mathcal{D} \subseteq \mathcal{X} \setminus \mathbf{I}$  of at most  $d$  labels such that if  $\Psi'(\alpha) \neq \Psi(\alpha)$  then  $\text{lab}(\alpha) \in \mathcal{D}$ , where  $\mathcal{X}$  is the set of all possible labels and  $\text{lab}$  is a labeling function. Then, it holds that  $\text{Decide}_{\mathfrak{S}}^{\Psi'}(V, \bar{V}, x) = \text{Decide}_{\mathfrak{S}}^{\Psi}(V, \bar{V}, x)$ . This is a simplified description of the key property on which we rely in order to prove that a random instance of  $\Psi$  realizes a one-way function and an indistinguishability obfuscator relative to our generalized decision oracle, and we refer the reader to Section 4 for the proof of our impossibility result.

### 1.3 Paper Organization

The remainder of this paper is organized as follows. In Section 2 we introduce some standard notation as well as the cryptographic primitives under consideration in this paper. In Section 3 we define the classes of constructions to which our impossibility results apply. In Section 4 we formally state and prove Theorem 1.1, and in Section 5 we formally state and prove Theorem 1.2.

## 2 Preliminaries

In this section we present the notation and basic definitions that are used in this work. For a distribution  $X$  we denote by  $x \leftarrow X$  the process of sampling a value  $x$  from the distribution  $X$ . Similarly, for a set  $\mathcal{X}$  we denote by  $x \leftarrow \mathcal{X}$  the process of sampling a value  $x$  from the uniform distribution over  $\mathcal{X}$ . For an integer  $n \in \mathbb{N}$  we denote by  $[n]$  the set  $\{1, \dots, n\}$ . For every  $n \in \mathbb{N}$  and  $m \geq n$  we denote by  $\text{InjFunc}_n^m$  the set of all injective functions  $f : \{0, 1\}^n \rightarrow \{0, 1\}^m$ .

**Oracle-aided languages and complexity classes.** For a language  $L \subseteq \{0, 1\}^*$ , we let  $\chi_L : \{0, 1\}^* \rightarrow \{0, 1\}$  denote the characteristic function of  $L$ , that is,  $\chi_L(x) = 1$  if and only if  $x \in L$ . A

deterministic algorithm  $A$  decides a language  $L$  if for every  $x \in \{0, 1\}^*$  it holds that  $A(x) = \chi_L(x)$ .

We consider the standard notions of languages and complexity classes when naturally generalized to oracle-aided computations. In particular, an oracle-aided language  $L$  defines a set  $L^\Gamma \subseteq \{0, 1\}^*$  for any possible oracle  $\Gamma : \{0, 1\}^* \rightarrow \{0, 1\}^*$ . Our definitions throughout the paper follow the standard approach that was introduced in the classic complexity-theory literature for proving separations between complexity classes by considering type-2 languages and complexity classes (see, for example, [BCE<sup>+</sup>95, CIY97] and the references therein).

**Indistinguishability obfuscation for oracle-aided circuits.** We consider the standard notion of indistinguishability obfuscation [BGI<sup>+</sup>12, GGH<sup>+</sup>13] when naturally generalized to oracle-aided circuits (i.e., circuits that may contain oracle gates in addition to standard gates) [AS15, AS16]. We first define the notion of functional equivalence relative to a specific function (provided as an oracle), and then we define the notion of an indistinguishability obfuscation for a class of oracle-aided circuits. In what follows, when considering a class  $\mathcal{C} = \{C_n\}_{n \in \mathbb{N}}$  of oracle-aided circuits, we assume that each  $C_n$  consists of circuits of size at most  $n$ .

**Definition 2.1.** Let  $C_0$  and  $C_1$  be two oracle-aided circuits, and let  $f$  be a function. We say that  $C_0$  and  $C_1$  are functionally equivalent relative to  $f$ , denoted  $C_0^f \equiv C_1^f$ , if for any input  $x$  it holds that  $C_0^f(x) = C_1^f(x)$ .

**Definition 2.2.** A probabilistic polynomial-time oracle-aided algorithm  $i\mathcal{O}$  is an indistinguishability obfuscator relative to an oracle  $\Gamma$  for a class  $\mathcal{C} = \{C_n\}_{n \in \mathbb{N}}$  of oracle-aided circuits if the following conditions are satisfied:

- **Functionality.** For all  $n \in \mathbb{N}$  and for all  $C \in \mathcal{C}_n$  it holds that

$$\Pr \left[ C^\Gamma \equiv \widehat{C}^\Gamma : \widehat{C} \leftarrow i\mathcal{O}^\Gamma(1^n, C) \right] = 1.$$

- **Indistinguishability.** For any probabilistic polynomial-time oracle-aided distinguisher  $\mathcal{A} = (\mathcal{A}_1, \mathcal{A}_2)$  there exists a negligible function  $\nu(\cdot)$  such that

$$\text{Adv}_{\Gamma, i\mathcal{O}, \mathcal{A}, \mathcal{C}}^{i\mathcal{O}}(n) \stackrel{\text{def}}{=} \left| \Pr \left[ \text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{A}, \mathcal{C}}^{i\mathcal{O}}(n) = 1 \right] - \frac{1}{2} \right| \leq \nu(n)$$

for all sufficiently large  $n \in \mathbb{N}$ , where the random variable  $\text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{A}, \mathcal{C}}^{i\mathcal{O}}(n)$  is defined via the following experiment:

1.  $b \leftarrow \{0, 1\}$ .
2.  $(C_0, C_1, \text{state}) \leftarrow \mathcal{A}_1^\Gamma(1^n)$  where  $C_0, C_1 \in \mathcal{C}_n$  and  $C_0^\Gamma \equiv C_1^\Gamma$ .
3.  $\widehat{C} \leftarrow i\mathcal{O}^\Gamma(1^n, C_b)$ .
4.  $b' \leftarrow \mathcal{A}_2^\Gamma(\text{state}, \widehat{C})$ .
5. If  $b' = b$  then output 1, and otherwise output 0.

For simplicity, note that whenever the algorithm  $\mathcal{A}_1$  is deterministic there is in fact no need for  $\mathcal{A}_1$  to transfer any state information **state** to  $\mathcal{A}_2$  as the state can be reconstructed if needed by invoking  $\mathcal{A}_1$ . Looking ahead, in this paper we consider computationally-unbounded algorithms (i.e., we limit their query complexity but we do not limit their internal computation), and such algorithms can be assumed without loss of generality to be deterministic.

### 3 The Classes of Constructions

In this section we formalize the proof systems and the classes of constructions that we consider in this paper. The proof systems we consider in this paper can be formalized in a variety of seemingly equivalent manners, and here we choose a specific definition that we find to simplify the proof of our impossibility results:

**Definition 3.1.** For functions  $V, P : \{0, 1\}^* \rightarrow \{0, 1\}^*$ , an integer  $k \geq 0$  and a string  $s \in \{0, 1\}^*$ , we denote by  $\langle V(s), P \rangle_k$  the output of the following computation:

- Let  $m_0 = P(V(s, 0))$ .
- For  $1 \leq i < k$ , let  $m_i = P(V(s, i, m_0, \dots, m_{i-1}))$ .
- Output  $V(s, k, m_0, \dots, m_{k-1}) \in \{0, 1\}$ .

That is, we consider a sequential process that is executed by two parties, a verifier  $V$  that is given as input a string  $s$ , and a prover  $P$  that is not given any input. The process consists of  $k$  rounds, where in each round the verifier sends the prover a message that is computed as a function of its input  $s$ , the index  $i$  of the current round, and the prover's previous responses  $m_0, \dots, m_{i-1}$ . In turn, the prover replies with a response  $m_i$ , and following these  $k$  steps the verifier outputs a bit indicating acceptance or rejection.

A crucial property to notice is that the prover's response,  $m_i$ , in each step is a function of the verifier's  $i$ th message only, and not of the entire transcript which includes all of the verifier's previous messages as well (i.e., the prover is "memoryless"). A verifier may potentially include the entire transcript in each message, and then the definition would collapse to the class **IP** of languages that have an interactive proof system [GMR89].

In general, however, a verifier need not send the entire transcript in each message, and this enables us to capture the class **MIP** of languages that have a multi-prover interactive proof system [BGK<sup>+</sup>88]. Specifically, any such proof system  $\langle V, P_1, \dots, P_N \rangle$  in which each prover sends at most  $v$  messages can be transformed in a black-box manner into a proof system  $\langle V, P \rangle_k$  of the above form with  $k = v \cdot N$ . This can be done, for example, by defining  $P(i, \cdot) = P_i(\cdot)$  for every  $i \in [T]$ , and having the verifier include in each message the index of the prover to which this message is sent together with the entire transcript that this specific prover has seen so far. Although we have not yet defined the completeness and soundness properties for the above proof systems (these are defined as part of the following definition), we already note that this transformation naturally preserves them.

As discussed in Section 1.2, instead of considering oracle-aided **MIP** proof systems we consider oracle-aided **MIP** verifiers, and allow the honest provers to depend on any given oracle in an arbitrary non-black-box manner (thus our result rules out, in particular, oracle-aided proof systems). This is captured via the following definition:

**Definition 3.2.** A pair  $(V, \bar{V})$  of oracle-aided polynomial-time algorithms, together with polynomials  $\ell_r(\cdot)$  and  $k(\cdot)$ , define a **(MIP, coMIP) protocol pair relative to an oracle**  $\Psi : \{0, 1\}^* \rightarrow \{0, 1\}^*$  if there exists a language  $L^\Psi \subseteq \{0, 1\}^*$  and such that:

- For every  $x \in L^\Psi$  there exists a function  $P : \{0, 1\}^* \rightarrow \{0, 1\}$  such that

$$\Pr_{r \leftarrow \{0, 1\}^{\ell_r(|x|)}} \left[ \langle V^\Psi(x, r), P \rangle_{k(|x|)} = 1 \right] \geq 2/3,$$

and for every function  $\bar{P} : \{0, 1\}^* \rightarrow \{0, 1\}$  it holds that

$$\Pr_{r \leftarrow \{0, 1\}^{\ell_r(|x|)}} \left[ \langle \bar{V}^\Psi(x, r), \bar{P} \rangle_{k(|x|)} = 1 \right] \leq 1/3.$$

- For every  $x \notin L^\Psi$  there exists a function  $\bar{P} : \{0, 1\}^* \rightarrow \{0, 1\}$  such that

$$\Pr_{r \leftarrow \{0,1\}^{\ell_r(|x|)}} \left[ \langle \bar{V}^\Psi(x, r), \bar{P} \rangle_{k(|x|)} = 1 \right] \geq 2/3,$$

and for every function  $P : \{0, 1\}^* \rightarrow \{0, 1\}$  it holds that

$$\Pr_{r \leftarrow \{0,1\}^{\ell_r(|x|)}} \left[ \langle V^\Psi(x, r), P \rangle_{k(|x|)} = 1 \right] \leq 1/3.$$

Note that the above definition considers provers that output only a single bit in each step. This is just for syntactical reasons, making sure that the verifier runs in polynomial-time with respect to the length of the input  $x$ . For example, if the prover was allowed to be a length-doubling function, then after  $|x|$  rounds this would allow a polynomial-time verifier to run in time that is exponential in the length of  $|x|$ . There are naturally various ways in which this technical issue can be handled (e.g., providing the verifier with oracle access to the prover instead of direct communication), clearly without having any effect on our result.

The following two definitions are based on those of [AS15, AS16, BDV17] (which, in turn, are motivated by [Lub96, Gol00, RTV04]), and capture the classes of constructions that we consider in this paper. Definition 3.3 considers black-box constructions of (MIP, coMIP) protocols based on injective one-way functions and indistinguishability obfuscation, and will be used in Section 4. Definition 3.4 considers for black-box constructions of (MIP, coMIP) protocols based on collision-resistant hash functions, and will be used in Section 5. We remind the reader that two oracle-aided circuits,  $C_0$  and  $C_1$ , are functionally equivalent relative to a function  $f$ , denoted  $C_0^f \equiv C_1^f$ , if for any input  $x$  it holds that  $C_0^f(x) = C_1^f(x)$  (see Definition 2.1).

**Definition 3.3.** A fully black-box construction of a worst-case hard (MIP, coMIP) protocol pair from an injective one-way function  $f$  and an indistinguishability obfuscator for the class  $\mathcal{C}$  of all oracle-aided circuits  $C^f$ , consists of a pair of oracle-aided polynomial-time algorithms  $(V, \bar{V})$ , polynomials  $\ell_r(\cdot)$  and  $k(\cdot)$ , an oracle-aided polynomial-time algorithm  $M$ , and “security loss” functions  $\epsilon_{M,1}(\cdot)$  and  $\epsilon_{M,2}(\cdot)$ , such that the following conditions hold:

- **Correctness:** For every ensemble  $f = \{f_n : \{0, 1\}^n \rightarrow \{0, 1\}^{n+1}\}_{n \in \mathbb{N}}$  of injective functions, and for any function  $i\mathcal{O}$  such that  $i\mathcal{O}(C; r)^f \equiv C^f$  for any circuit  $C$  and  $r \in \{0, 1\}^*$ , the pair  $(V, \bar{V})$ , together with the polynomials  $\ell_r(\cdot)$  and  $k(\cdot)$ , define an (MIP, coMIP) protocol pair (with a corresponding language  $L^{f, i\mathcal{O}}$ ) relative to the oracle  $(f, i\mathcal{O})$ .
- **Black-box proof of hardness:** For every ensemble  $f = \{f_n : \{0, 1\}^n \rightarrow \{0, 1\}^{n+1}\}_{n \in \mathbb{N}}$  of injective functions, for any function  $i\mathcal{O}$  such that  $i\mathcal{O}(C; r)^f \equiv C^f$  for any circuit  $C$  and  $r \in \{0, 1\}^*$ , and for any oracle-aided algorithm  $\mathcal{A}$  that runs in time  $T_{\mathcal{A}}(\cdot)$ , if  $\mathcal{A}^{f, i\mathcal{O}}(x) = \chi_{L^{f, i\mathcal{O}}}(x)$  for every  $x \in \{0, 1\}^*$  then either

$$\Pr \left[ M^{f, i\mathcal{O}, \mathcal{A}}(f(x)) = x \right] \geq \epsilon_{M,1}(T_{\mathcal{A}}(n)) \cdot \epsilon_{M,2}(n)$$

for infinitely many values of  $n \in \mathbb{N}$ , where the probability is taken over the choice of  $x \leftarrow \{0, 1\}^n$  and over the internal randomness of  $M$ , or

$$\left| \Pr \left[ \text{Exp}_{(f, i\mathcal{O}), i\mathcal{O}, M^{\mathcal{A}, \mathcal{C}}}(n) = 1 \right] - \frac{1}{2} \right| \geq \epsilon_{M,1}(T_{\mathcal{A}}(n)) \cdot \epsilon_{M,2}(n)$$

for infinitely many values of  $n \in \mathbb{N}$ .

Intuitively, a black-box proof of hardness for  $L^{f,i\mathcal{O}}$ , as formalized by the above definition, means that any algorithm that decides  $L^{f,i\mathcal{O}}$  can be used to construct an adversary that breaks either the one-wayness of  $f$  or the indistinguishability property of  $i\mathcal{O}$  in a black-box way.

**Definition 3.4.** A fully black-box construction of a worst-case hard (MIP, coMIP) protocol pair from a collision-resistant hash function  $f$  consists of a pair of oracle-aided polynomial-time algorithms  $(V, \bar{V})$ , polynomials  $\ell_r(\cdot)$  and  $k(\cdot)$ , an oracle-aided polynomial-time algorithm  $M$ , and “security loss” functions  $\epsilon_{M,1}(\cdot)$  and  $\epsilon_{M,2}(\cdot)$ , such that the following conditions hold:

- **Correctness:** For every ensemble  $f = \{f_n : \{0, 1\}^{2n} \rightarrow \{0, 1\}^{n-1}\}_{n \in \mathbb{N}}$  of functions, the pair  $(V, \bar{V})$ , together with the polynomials  $\ell_r(\cdot)$  and  $k(\cdot)$ , define an (MIP, coMIP) protocol pair (with a corresponding language  $L^f$ ) relative to the oracle  $f$ .
- **Black-box proof of hardness:** For every ensemble  $f = \{f_n : \{0, 1\}^{2n} \rightarrow \{0, 1\}^{n-1}\}_{n \in \mathbb{N}}$  of functions, and for any oracle-aided algorithm  $\mathcal{A}$  that runs in time  $T_{\mathcal{A}}(\cdot)$ , if  $\mathcal{A}^{f,i\mathcal{O}}(x) = \chi_{L^{f,i\mathcal{O}}}(x)$  for every  $x \in \{0, 1\}^*$ , then for infinitely many values of  $n \in \mathbb{N}$  it holds that

$$\Pr \left[ \left( x_1 \neq x_2 \right) \wedge \left( f_n(s, x_1) = f_n(s, x_2) \right) \right] \geq \epsilon_{M,1}(T_{\mathcal{A}}(n)) \cdot \epsilon_{M,2}(n),$$

where  $s \leftarrow \{0, 1\}^n$  and  $(x_1, x_2) \leftarrow M^{f,\mathcal{A}}(s)$ .

Intuitively, a black-box proof of hardness for  $L^f$ , as formalized in the above definition, means that any algorithm that decides  $L^f$  can be used to construct an adversary that breaks the collision-resistance of  $f$  in a black-box way.

Note that in both definitions, restricting  $\mathcal{A}$  to be deterministic and to decide the language in the worst case (i.e., on all inputs) only makes our impossibility result stronger. Furthermore, in Definition 3.4 we consider *public-coin* collision-resistant hash functions, and this once again only makes our impossibility result stronger. Also note that, following Asharov and Segev [AS15, AS16], we split the security loss in the above definition to an adversary-dependent security loss (the function  $\epsilon_{M,1}(\cdot)$ ) and an adversary-independent security loss (the function  $\epsilon_{M,2}(\cdot)$ ), as this allows us to also rule out constructions in which one of these losses is super-polynomial while the other is polynomial.

## 4 Impossibility Result for Constructions based on $i\mathcal{O}$ and Injective OWFs

Equipped with a formal definition of a black-box construction of a hard (MIP, coMIP) protocol pair from an injective one-way function and an indistinguishability obfuscator (recall Definition 3.3), in this section we prove the following theorem:

**Theorem 4.1.** *Let  $((V, \bar{V}), \ell_r, k, M, T_M, \epsilon_{M,1}, \epsilon_{M,2})$  be a fully black-box construction of a worst-case hard (MIP, coMIP) protocol pair from an injective one-way function  $f$  and an indistinguishability obfuscator for all oracle-aided circuits  $C^f$ . Then, it holds that*

$$\epsilon_{M,1}(n) \cdot \epsilon_{M,2}(n) \leq 2^{-\Omega(n)}.$$

*That is, at least one out of the adversary-dependent security loss  $\epsilon_{M,1}(\cdot)$  and the adversary-independent security loss  $\epsilon_{M,2}(\cdot)$  is exponential.*

Theorem 4.1 rules out, in particular, standard “polynomial-time polynomial-loss” reductions. More generally, the theorem implies that if the adversary-dependent security loss  $\epsilon_{M,1}(\cdot)$  is polynomial (as is typically the case in cryptographic reductions), then the adversary-independent security

loss  $\epsilon_{M,2}(\cdot)$  must be exponential. Thus, this also rules out constructions that are based on indistinguishability obfuscation with sub-exponential security (e.g., [BPR15, BPW16]).

In what follows we first introduce our generalized decision oracle, and capture its main property on which we rely in our proof, as discussed in Section 1.2. Then, in Section 4.2 we introduce the additional oracles on which we rely, and in Sections 4.3 and 4.4 we prove that relative to these oracles and to our decision oracle there exist an injective one-way function and an indistinguishability obfuscator, respectively. Finally, in Section 4.5 we derive the proof of Theorem 4.1.

#### 4.1 Our Generalized Decision Oracle

For a family of oracles  $\mathfrak{G}$  and for any specific oracle  $\Psi \in \mathfrak{G}$ , we define the oracle  $\text{Decide}_{\mathfrak{G}}^{\Psi}$  as the following function: Given as input tuple  $(C_0, C_1, 1^{\ell_r}, 1^k)$ , where  $C_0$  and  $C_1$  are oracle-aided circuits, and  $\ell_r$  and  $k$  are non-negative integers, for every  $\Phi \in \mathfrak{G}$  the oracle checks if one of the following two cases holds:

- There exists a function  $P_1 : \{0, 1\}^* \rightarrow \{0, 1\}$  such that

$$\Pr_{r \leftarrow \{0,1\}^{\ell_r}} [\langle C_1^{\Phi}(r), P_1 \rangle_k = 1] \geq 2/3 ,$$

and for every function  $P_0 : \{0, 1\}^* \rightarrow \{0, 1\}$  it holds that

$$\Pr_{r \leftarrow \{0,1\}^{\ell_r}} [\langle C_0^{\Phi}(r), P_0 \rangle_k = 1] \leq 1/3 .$$

In this case, we say that  $(C_0^{\Phi}, C_1^{\Phi}, 1^{\ell_r}, 1^k)$  is a *yes-instance*.

- There exists a function  $P_0 : \{0, 1\}^* \rightarrow \{0, 1\}$  such that

$$\Pr_{r \leftarrow \{0,1\}^{\ell_r}} [\langle C_0^{\Phi}(r), P_0 \rangle_k = 1] \geq 2/3 ,$$

and for every function  $P_1 : \{0, 1\}^* \rightarrow \{0, 1\}$  it holds that

$$\Pr_{r \leftarrow \{0,1\}^{\ell_r}} [\langle C_1^{\Phi}(r), P_1 \rangle_k = 1] \leq 1/3 .$$

In this case, we say that  $(C_0^{\Phi}, C_1^{\Phi}, 1^{\ell_r}, 1^k)$  is a *no-instance*.

If there exists an oracle  $\Phi \in \mathfrak{G}$  such that none of the above cases hold, then we say that the input  $(C_0, C_1, 1^{\ell_r}, 1^k)$  is *invalid* and set  $\text{Decide}_{\mathfrak{G}}^{\Psi}$  to output  $\perp$ . Otherwise,  $\text{Decide}_{\mathfrak{G}}^{\Psi}$  outputs 1 or 0 according to whether  $(C_0^{\Psi}, C_1^{\Psi}, 1^{\ell_r}, 1^k)$  is a yes-instance or a no-instance.<sup>6</sup>

The following simple lemma shows that the oracle  $\text{Decide}_{\mathfrak{G}}^{\Psi}$  can be easily used in order to decide any language that is defined via a (MIP, coMIP) protocol pair:

**Lemma 4.2.** *Let  $\mathfrak{G}$  be a family of oracles, and let  $(V, \bar{V})$  be a pair of oracle-aided polynomial-time algorithms that is an (MIP, coMIP) protocol pair, with respect to polynomials  $\ell_r(\cdot)$  and  $k(\cdot)$ , relative to every oracle  $\Psi \in \mathfrak{G}$ . Then, there exists a polynomial-time single-query algorithm  $\mathcal{A}$  such that for every  $\Psi \in \mathfrak{G}$ , given oracle access to  $\text{Decide}_{\mathfrak{G}}^{\Psi}$  the algorithm  $\mathcal{A}$  decides the language  $L^{\Psi} \subseteq \{0, 1\}^*$  defined by  $(V, \bar{V}, \ell_r, k)$  relative to  $\Psi$ . That is, for every  $\Psi \in \mathfrak{G}$  and  $x \in \{0, 1\}^*$  the algorithm  $\mathcal{A}^{\text{Decide}_{\mathfrak{G}}^{\Psi}}(x)$  outputs 1 if and only if  $x \in L^{\Psi}$ .*

<sup>6</sup>Note that for an input  $(C_0, C_1, 1^{\ell_r}, 1^k)$ , either it is invalid and then  $\text{Decide}_{\mathfrak{G}}^{\Psi}$  outputs  $\perp$  for every  $\Psi \in \mathfrak{G}$ , or it is valid and then  $\text{Decide}_{\mathfrak{G}}^{\Psi}$  outputs 0 or 1 depending on  $\Psi$ .

**Proof.** Since  $V$  and  $\bar{V}$  are polynomial time, there exists a polynomial  $p(n)$  such that on input of size  $n$  their output is of size at most  $p(n)$ . Given  $x \in \{0, 1\}^*$  as input and oracle access to  $\text{Decide}_{\mathfrak{G}}^{\Psi}$ , the algorithm  $\mathcal{A}$  queries  $\text{Decide}_{\mathfrak{G}}^{\Psi}$  on  $(C_0, C_1, 1^{\ell_r(|x|)}, 1^{k(|x|)})$ , where  $C_0$  and  $C_1$  are the hardwired oracle-aided circuits  $\bar{V}(x, \cdot)$  and  $V(x, \cdot)$  respectively, the input size of both circuits is  $\lceil \log(k(|x|) + 1) \rceil + k(|x|)$  (where  $\lceil \log(k(|x|) + 1) \rceil$  bits are for the index of the communication round and  $k(|x|)$  bits are for the messages of the prover) and the output size is  $p(|x| + \lceil \log(k(|x|) + 1) \rceil + k(|x|))$ . Finally, the algorithm  $\mathcal{A}$  outputs 1 if and only if the oracle's response to the query is 1.  $\blacksquare$

The following lemma captures the key property of our oracle, as discussed in Section 1.2:

**Lemma 4.3.** *Let  $\mathfrak{S}$  be a family of oracles, let  $\mathcal{Q}$  be the set of all possible queries for every oracle in the family, let  $\text{lab} : \mathcal{Q} \rightarrow \mathcal{X}$  be a “labeling” of the possible queries, and let  $d \in \mathbb{N}$  be a parameter.*

*For any  $\Psi \in \mathfrak{S}$  and for any  $\text{Decide}_{\mathfrak{G}}^{\Psi}$ -query  $(C_0, C_1, 1^{\ell_r}, 1^k)$  such that each of the circuits  $C_0$  and  $C_1$  contains at most  $q$  oracle gates, there exists a set of labels  $\mathbf{I} = \mathbf{I}(\mathfrak{S}, \Psi, C_0, C_1, \ell_r, k, \text{lab}, d) \subseteq \mathcal{X}$ , which we call the influential labels, satisfying the following two properties:*

1. *The set is small:  $|\mathbf{I}| \leq 3 \cdot q \cdot k \cdot d$ .*
2. *Any changes to  $\Psi$  involving at most  $d$  labels outside of  $\mathbf{I}$  do not change the answer of the query: Let  $\Phi \in \mathfrak{S}$  be another oracle, such that there exists a set  $\mathcal{D} \subseteq \mathcal{X} \setminus \mathbf{I}$  of labels with cardinality at most  $d$  such that if  $\Phi(q) \neq \Psi(q)$  then  $\text{lab}(q) \in \mathcal{D}$ . Then, it holds that*

$$\text{Decide}_{\mathfrak{G}}^{\Phi}(C_0, C_1, 1^{\ell_r}, 1^k) = \text{Decide}_{\mathfrak{G}}^{\Psi}(C_0, C_1, 1^{\ell_r}, 1^k).$$

**Proof.** If  $\text{Decide}_{\mathfrak{G}}^{\Psi}(C_0, C_1, 1^{\ell_r}, 1^k) = \perp$  this means that the input  $(C_0, C_1, 1^{\ell_r}, 1^k)$  is invalid, and then  $\text{Decide}_{\mathfrak{G}}^{\Phi}(C_0, C_1, 1^{\ell_r}, 1^k) = \perp$  holds for every  $\Phi \in \mathfrak{S}$  and the claim follows for  $\mathbf{I} = \emptyset$ . Otherwise, suppose without loss of generality that  $\text{Decide}_{\mathfrak{G}}^{\Psi}(C_0, C_1, 1^{\ell_r}, 1^k) = 1$  and let  $P_1 : \{0, 1\}^* \rightarrow \{0, 1\}$  such that

$$\Pr_{r \leftarrow \{0, 1\}^{\ell_r}} [\langle C_1^{\Psi}(r), P_1 \rangle_k = 1] \geq 2/3.$$

Roughly speaking, we define  $\mathbf{I} \subseteq \mathcal{X}$  to be the set of all labels for which a query with that label is performed during the execution of the protocol  $\langle C_1^{\Psi}(\cdot), P_1 \rangle_k$  with high probability over the choice of  $r$ . More formally, we define

$$\mathbf{I} = \left\{ \text{label} \in \mathcal{X} \left| \Pr_{r \leftarrow \{0, 1\}^{\ell_r}} \left[ \begin{array}{c} \text{A query } q \in \mathcal{Q} \text{ such that } \text{lab}(q) = \text{label} \text{ is performed} \\ \text{during the execution of } \langle C_1^{\Psi}(r), P_1 \rangle_k \end{array} \right] \geq \frac{1}{3 \cdot d} \right. \right\}.$$

First, for every  $r \in \{0, 1\}^{\ell_r}$  at most  $q \cdot k$  queries are performed during the execution of  $\langle C_1^{\Psi}(r), P_1 \rangle_k$ . Therefore, for any  $0 < \epsilon \leq 1$  there are at most  $q \cdot k / \epsilon$  labels such that

$$\Pr_{r \leftarrow \{0, 1\}^{\ell_r}} \left[ \begin{array}{c} \text{A query } q \in \mathcal{Q} \text{ such that } \text{lab}(q) = \text{label} \text{ is performed} \\ \text{during the execution of } \langle C_1^{\Psi}(r), P_1 \rangle_k \end{array} \right] \geq \epsilon.$$

In our case, this means that  $|\mathbf{I}| \leq q \cdot k \cdot 3 \cdot d$  as claimed.

Next, let  $\Phi \in \mathfrak{S}$  such that there exists a set  $\mathcal{D} \subseteq \mathcal{X} \setminus \mathbf{I}$  of labels with cardinality at most  $d$  such that if  $\Phi(q) \neq \Psi(q)$  then  $\text{lab}(q) \in \mathcal{D}$ . By a union bound it holds that

$$\Pr_{r \leftarrow \{0, 1\}^{\ell_r}} \left[ \begin{array}{c} \text{A query } q \in \mathcal{Q} \text{ such that } \text{lab}(q) \in \mathcal{D} \text{ is performed} \\ \text{during the execution of } \langle C_1^{\Psi}(r), P_1 \rangle_k \end{array} \right] < \frac{|\mathcal{D}|}{3 \cdot d} \leq \frac{1}{3}.$$

If the above event does not occur then  $\langle C_1^\Phi(r), P_1 \rangle_k = \langle C_1^\Psi(r), P_1 \rangle_k$ . Hence,

$$\begin{aligned} & \Pr_{r \leftarrow \{0,1\}^{\ell_r}} [\langle C_1^\Phi(r), P_1 \rangle_k = 1] \\ & \geq \Pr_{r \leftarrow \{0,1\}^{\ell_r}} [\langle C_1^\Psi(r), P_1 \rangle_k = 1] - \Pr_{r \leftarrow \{0,1\}^{\ell_r}} [\langle C_1^\Phi(r), P_1 \rangle_k \neq \langle C_1^\Psi(r), P_1 \rangle_k] \\ & > \frac{2}{3} - \frac{1}{3} = \frac{1}{3}, \end{aligned}$$

so  $(C_0^\Phi, C_1^\Phi, 1^{\ell_r}, 1^k)$  is not a no-instance. Since  $\text{Decide}_{\mathfrak{S}}^\Psi(C_0, C_1, 1^{\ell_r}, 1^k) \neq \perp$ ,  $(C_0^\Phi, C_1^\Phi, 1^{\ell_r}, 1^k)$  must be a yes-instance and therefore  $\text{Decide}_{\mathfrak{S}}^\Phi(C_0, C_1, 1^{\ell_r}, 1^k) = 1 = \text{Decide}_{\mathfrak{S}}^\Psi(C_0, C_1, 1^{\ell_r}, 1^k)$  as claimed. ■

## 4.2 Our Indistinguishability Obfuscation Oracle

In what follows we define the family  $\mathfrak{S}$  of oracles that realize injective functions and strongly-unambiguous obfuscations relative to our decision oracle, and define a distribution  $\mathcal{D}(\mathfrak{S})$  over that family. The family  $\mathfrak{S}$  consists of all triplets  $(f, \mathcal{O}, E) = (\{f_n\}_{n \in \mathbb{N}}, \{\mathcal{O}_n\}_{n \in \mathbb{N}}, \{E_n\}_{n \in \mathbb{N}})$ , satisfying the following three conditions for every  $n \in \mathbb{N}$ :

1. The function  $f_n : \{0, 1\}^n \rightarrow \{0, 1\}^{n+1}$  is injective. Looking ahead,  $f$  will serve as an injective one-way function.
2. The function  $\mathcal{O}_n : \{0, 1\}^{2n} \rightarrow \{0, 1\}^{10n}$  is injective. Looking ahead, for an oracle-aided circuit  $C \in \{0, 1\}^n$  with  $f$ -gates and randomness  $r \in \{0, 1\}^n$ , the output  $\mathcal{O}_n(C, r)$  will serve as an obfuscation of  $C$ , and the restriction that  $\mathcal{O}_n$  is injective means that the obfuscation is *strongly-unambiguous* in the sense that any obfuscation  $\widehat{C} \in \text{Image}(\mathcal{O}_n)$  only comes from a single circuit with a single randomness string.
3. The function  $E_n : \{0, 1\}^{11n} \rightarrow \{0, 1\}^n$  satisfies the following condition: For every oracle-aided circuit  $C \in \{0, 1\}^n$  with  $f$ -gates, every randomness  $r \in \{0, 1\}^n$  and every input  $\alpha \in \{0, 1\}^n$ , it holds that  $E_n(\mathcal{O}_n(C, r), \alpha) = C^f(x)$ . Namely, given an obfuscation  $\widehat{C} = \mathcal{O}_n(C, r)$  and an input  $\alpha$ , the function  $E_n$  evaluates  $C$  on input  $\alpha$  with respect to the oracle  $f$ .

We emphasize that for any  $\widehat{C} \in \{0, 1\}^{10n} \setminus \text{Image}(\mathcal{O}_n)$ , there is no restriction on  $E_n(\widehat{C}, \cdot)$ , so there is no clear way to verify whether some  $\widehat{C} \in \{0, 1\}^{10n}$  is a valid obfuscation. As noted by Bitansky et al. [BDV17], it is necessary for the obfuscation to not be verifiable since an unambiguous and verifiable indistinguishability obfuscator does imply hardness in  $\text{NP} \cap \text{coNP}$ .

Now, we define a distribution  $\mathcal{D}(\mathfrak{S})$  over  $\mathfrak{S}$ , relative to which we prove that an oracle  $\Psi \leftarrow \mathcal{D}(\mathfrak{S})$  realizes an injective one-way function and an indistinguishability obfuscator. The distribution  $\mathcal{D}(\mathfrak{S})$  is obtained by sampling a triplet  $(f, \mathcal{O}, \text{Eval}^{f, \mathcal{O}})$  from  $\mathfrak{S}$  as follows:

1. For every  $n \in \mathbb{N}$  the function  $f_n$  is uniformly chosen from the set  $\text{InjFunc}_n^{n+1}$  of all injective functions  $f_n : \{0, 1\}^n \rightarrow \{0, 1\}^{n+1}$ .
2. For every  $n \in \mathbb{N}$  the function  $\mathcal{O}_n : \{0, 1\}^{2n} \rightarrow \{0, 1\}^{10n}$  is sampled as follows: A function  $h$  is uniformly chosen from the set  $\text{InjFunc}_n^{5n}$ , and for every  $r \in \{0, 1\}^n$  a function  $g_r$  is uniformly chosen from the set  $\text{InjFunc}_n^{5n}$ . Then, for every input  $(C, r) \in \{0, 1\}^n \times \{0, 1\}^n$  we define  $\mathcal{O}_n(C, r) = (h(r), g_r(C))$ . Note that  $\mathcal{O}_n$  is injective as required, and that this distribution of the function  $\mathcal{O}$  differs from that of Asharov and Segev [AS15] and Bitansky et al. [BDV17], where  $\mathcal{O}_n$  was a uniformly chosen injective function.

3. For every  $n \in \mathbb{N}$ , the function  $\text{Eval}^{f, \mathcal{O}}$  on input  $(\widehat{C}, \alpha) \in \{0, 1\}^{10n} \times \{0, 1\}^n$  is defined as follows: If there exists a pair  $(C, r) \in \{0, 1\}^n \times \{0, 1\}^n$  such that  $\widehat{C} = \mathcal{O}_n(C, r)$  then it outputs  $C^f(\alpha)$ , and otherwise it outputs  $\perp$ . Note that  $\text{Eval}^{f, \mathcal{O}}$  satisfies the above third condition for membership in  $\mathfrak{S}$ .

### 4.3 The Existence of an Injective One-Way Function

In this section we prove that the injective function  $f$  is one way relative to  $(\Psi, \text{Decide}_{\mathfrak{S}}^{\Psi})$ , where  $\Psi = (f, \mathcal{O}, \text{Eval}^{f, \mathcal{O}})$  is sampled from the distribution  $\mathcal{D}(\mathfrak{S})$  over  $\mathfrak{S}$  (see Section 4.2 for the description of this distribution). Our proof follows the structure of that of Bitansky, Degwekar and Vaikuntanathan [BDV17], while strengthened to deal with our generalized decision oracle as explained in Section 1.2.

In what follows we call an oracle-aided algorithm  $\mathcal{A}$  a  $q$ -query algorithm, for a function  $q = q(n)$ , if when given any input  $x \in \{0, 1\}^n$  it issues at most  $q(n)$  queries to the oracle  $\Gamma$ , each of its queries to  $\text{Eval}$  and  $\text{Decide}$  consists of circuits with at most  $q(n)$  oracle gates, and the number of communication rounds in the proof systems corresponding to each of its queries to  $\text{Decide}$  is at most  $q(n)$ .

**Theorem 4.4.** *For any oracle-aided  $2^{n/12}$ -query algorithm  $\mathcal{A}$  it holds that*

$$\Pr_{\substack{\Psi \leftarrow \mathcal{D}(\mathfrak{S}) \\ x \leftarrow \{0, 1\}^n}} \left[ \mathcal{A}^{\Psi, \text{Decide}_{\mathfrak{S}}^{\Psi}}(f(x)) = x \right] \leq O(2^{-n/2})$$

for all sufficiently large  $n \in \mathbb{N}$ .

In what follows, we let  $\mathfrak{F}$  denote the family of ensembles  $f = \{f_n\}_{n \in \mathbb{N}}$  where  $f_n \in \text{InjFunc}_n^{n+1}$  for all  $n \in \mathbb{N}$ . As our first step, we prove that  $f \leftarrow \mathfrak{F}$  is one way relative to the oracle  $(f, \text{Decide}_{\mathfrak{F}}^f)$ .

**Lemma 4.5.** *For any oracle-aided  $2^{n/6}$ -query algorithm  $\mathcal{A}$  it holds that*

$$\Pr_{\substack{f \leftarrow \mathfrak{F} \\ x \leftarrow \{0, 1\}^n}} \left[ \mathcal{A}^{f, \text{Decide}_{\mathfrak{F}}^f}(f(x)) = x \right] \leq O(2^{-n/2}).$$

**Proof.** We prove that the lemma holds when even fixing the oracles  $f_{-n} = \{f_k\}_{k \neq n}$  and only sampling  $f_n$ . We introduce a sequence of three hybrid experiments such that the first hybrid experiment  $\mathcal{H}_1$  is the real one-wayness experiment and the last hybrid experiment  $\mathcal{H}_3$  is an experiment in which the probability of the adversary is of winning is  $1/2^n$ . Then, by upper bounding the difference in the winning probability between each pair of consecutive hybrid experiments we deduce our claim.

**The hybrid  $\mathcal{H}_1$ .** This is the real experiment in which we sample  $x \leftarrow \{0, 1\}^n$ , give  $f_n(x) \in \{0, 1\}^{n+1}$  to  $\mathcal{A}$  as input, and give  $\mathcal{A}$  oracle access to  $\Gamma = (f, \text{Decide}_{\mathfrak{F}}^f)$ .

**The hybrid  $\mathcal{H}_2$ .** In this experiment, we sample  $y \leftarrow \{0, 1\}^{n+1} \setminus \text{Image}(f_n)$ , give  $y$  to  $\mathcal{A}$  as input, and give  $\mathcal{A}$  oracle access to  $\Gamma' = (f_{x \mapsto y}, \text{Decide}_{\mathfrak{F}}^{f_{x \mapsto y}})$ , where  $f_{x \mapsto y}$  is defined as

$$f_{x \mapsto y}(z) = \begin{cases} y & \text{if } z = x \\ f(z) & \text{otherwise} \end{cases}.$$

That is, we “plant”  $y$  as the challenge and as the image of  $x$ .

**The hybrid  $\mathcal{H}_3$ .** This experiment is obtained from  $\mathcal{H}_2$  by giving  $\mathcal{A}$  oracle access to the original oracle  $\Gamma$  instead of the oracle  $\Gamma'$  with the planted  $y$ , while still giving  $\mathcal{A}$  the planted  $y$  as input.

The following table summarizes our hybrid experiments:

Hybrid	$\mathcal{H}_1$	$\mathcal{H}_2$	$\mathcal{H}_3$
<b>Challenger Randomness</b>	$x \leftarrow \{0, 1\}^n$		
<b>Injective Function</b>	$f_n \leftarrow \text{InjFunc}_n^{n+1}$		
<b>Challenge</b>	$f_n(x)$	$y \leftarrow \{0, 1\}^{n+1} \setminus \text{Image}(f_n)$	
<b>Oracle</b>	$\Gamma = (f, \text{Decide}_{\mathfrak{F}}^f)$	$\Gamma' = (f_{x \rightarrow y}, \text{Decide}_{\mathfrak{F}}^{f_{x \rightarrow y}})$	$\Gamma = (f, \text{Decide}_{\mathfrak{F}}^f)$
<b>Winning Condition</b>	$\mathcal{A}$ outputs $x$		

**Claim 4.6.**  $\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_1] = \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_2]$ .

**Proof.** We couple the experiments  $\mathcal{H}_1$  and  $\mathcal{H}_2$  as follows.<sup>7</sup> First, we sample the same  $x \leftarrow \{0, 1\}^n$  for both experiments. Then, we uniformly sample a random injective function  $\widehat{f} : \{0, 1\}^n \setminus \{x\} \rightarrow \{0, 1\}^{n+1}$ . Next, we sample two distinct  $y, y' \leftarrow \{0, 1\} \setminus \text{Image}(\widehat{f})$ . Now, in  $\mathcal{H}_1$  we let the injective function be

$$f_n(z) = \begin{cases} y & \text{if } z = x \\ \widehat{f}(z) & \text{otherwise} \end{cases},$$

whereas in  $\mathcal{H}_2$  we let the injective function be

$$f'_n(z) = \begin{cases} y' & \text{if } z = x \\ \widehat{f}(z) & \text{otherwise} \end{cases},$$

and let  $y$  be the planted challenge. It is easy to see that the marginal distribution in both experiments is correct, and that both experiments are identical. That is,  $\mathcal{A}$  gets the same challenge as input and gets access to the same oracle, thus the claim follows.  $\blacksquare$

**Claim 4.7.**  $|\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_2] - \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_3]| \leq 3 \cdot q^3 / 2^n$ .

**Proof.** We observe that the view of  $\mathcal{A}$  in  $\mathcal{H}_3$  is independent of the choice of  $x$ . Therefore, if a query to  $f_n$  is made, then the probability of it to be  $x$  is at most  $1/2^n$ . In any other case, the answer to this query is the same in  $\mathcal{H}_2$  and  $\mathcal{H}_3$ , and both executions proceed the same way.

Now, if a query  $(C_0, C_1, 1^{\ell_r}, 1^k)$  to  $\text{Decide}_{\mathfrak{F}}^f$  is made, then we apply Lemma 4.3. We take the label function  $\text{lab} : \mathcal{Q} \rightarrow \mathcal{X}$  to be the identity function. The set  $\mathbf{I} = \mathbf{I}(\mathfrak{F}, f, C_0, C_1, \ell_r, k, \text{lab}, 1) \subseteq \mathcal{X}$  of influential labels is independent of the choice of  $x$ . Therefore, the probability that  $\mathbf{I}$  contains  $x$  is at most  $|\mathbf{I}|/2^n \leq 3q^2/2^n$ . In any other case, the oracle  $f_{x \rightarrow y}$  of  $\mathcal{H}_2$  is obtained from  $f$  of  $\mathcal{H}_3$  by changes involving one label outside of  $\mathbf{I}$ , and therefore by Lemma 4.3 it holds that

$$\text{Decide}_{\mathfrak{F}}^{f_{x \rightarrow y}}(C_0, C_1, 1^{\ell_k}, 1^k) = \text{Decide}_{\mathfrak{F}}^f(C_0, C_1, 1^{\ell_k}, 1^k),$$

<sup>7</sup>To couple two probability distributions means to define a joint distribution whose marginals are exactly those two distributions.

and both executions proceed the same way. Applying a union bound we deduce that

$$\begin{aligned} & |\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_2] - \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_3]| \\ & \leq \Pr \left[ \text{A query was answered differently in } \mathcal{A}^\Gamma(y) \text{ and } \mathcal{A}^{\Gamma'}(y) \right] \\ & \leq \frac{3q^3}{2^n}, \end{aligned}$$

and the claim follows. ■

**Claim 4.8.**  $\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_3] = 1/2^n$ .

**Proof.** In this experiment the view of  $\mathcal{A}$  is independent of  $x$ . ■

Now we turn back to proving Lemma 4.5. It holds that

$$\begin{aligned} \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_1] & \leq |\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_1] - \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_2]| \\ & \quad + |\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_2] - \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_3]| + \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_3] \\ & \leq 0 + \frac{3q^3}{2^n} + \frac{1}{2^n} = O\left(\frac{q^3}{2^n}\right), \end{aligned}$$

and by plugging  $q = 2^{n/6}$  we obtain Lemma 4.5. ■

Now, as our second and final step, we show how to deduce Theorem 4.4 from Lemma 4.5.

**Proof of Theorem 4.4.** We prove that the theorem holds when even fixing the oracle  $\mathcal{O}$  and only sampling  $f$ . Similar to Bitansky, Degwekar and Vaikuntanathan [BDV17], we show how to convert a  $q$ -query adversary  $\mathcal{A}$  that inverts  $f_n$  when given access to the oracle  $(\Psi, \text{Decide}_{\mathfrak{G}}^\Psi)$ , where  $\Psi = (f, \mathcal{O}, \text{Eval}^{f, \mathcal{O}})$ , into a  $q^2$ -query adversary  $\mathcal{B}$  with the same winning probability but that only has access to the oracle  $(f, \text{Decide}_{\mathfrak{F}}^f)$ . The algorithm  $\mathcal{B}$  simulates the algorithm  $\mathcal{A}$ , and upon each query to  $(\Psi, \text{Decide}_{\mathfrak{G}}^\Psi)$ ,  $\mathcal{B}$  acts as follows:

- A query to  $f$  is answered by forwarding it to the oracle  $f$  (to which  $\mathcal{B}$  has access) and providing  $\mathcal{A}$  with its response.
- A query to  $\mathcal{O}$  is answered according to the fixed  $\mathcal{O}$  without any queries.
- A query  $(\tilde{C}, \alpha) \in \{0, 1\}^{10m} \times \{0, 1\}^m$  to  $\text{Eval}^{f, \mathcal{O}}$  is answered as follows: If there exists an oracle-aided circuit  $C \in \{0, 1\}^m$  and  $r \in \{0, 1\}^m$  for which  $\mathcal{O}(C, r) = \tilde{C}$  according to the fixed  $\mathcal{O}$ , then  $\mathcal{B}$  computes  $C^f(\alpha)$  and return the answer to  $\mathcal{A}$ . Otherwise,  $\mathcal{B}$  answers with  $\perp$ . This step requires at most  $q$  queries to  $f$  per query.
- A query  $(C_0, C_1, 1^{\ell_r}, 1^k)$  to  $\text{Decide}_{\mathfrak{G}}^\Psi$  is answered as follows: If the input  $(C_0, C_1, 1^{\ell_r}, 1^k)$  is invalid, then  $\mathcal{B}$  returns the answer  $\perp$  to  $\mathcal{A}$ . Otherwise,  $\mathcal{B}$  constructs  $q^2$ -query circuits  $C'_0$  and  $C'_1$  with only  $f$  gates from  $C_0$  and  $C_1$  by hard-wiring the entire function  $\mathcal{O}$ , replacing each  $\mathcal{O}$ -gate with a direct computation from the hardwired  $\mathcal{O}$ , and replacing each  $\text{Eval}^{f, \mathcal{O}}$ -gate with the same computation done above for queries to  $\text{Eval}^{f, \mathcal{O}}$ , which requires at most  $q$  additional  $f$ -gates per  $\text{Eval}^{f, \mathcal{O}}$ -gate. Thus, the number of  $f$ -gates in each of the circuits  $C'_0$  and  $C'_1$  is at most  $q^2$ . Then,  $\mathcal{B}$  queries  $\text{Decide}_{\mathfrak{F}}^f$  with  $(C'_0, C'_1, 1^{\ell_r}, 1^k)$  and returns the answer to  $\mathcal{A}$ .

Overall,  $\mathcal{B}$  perfectly simulates  $\mathcal{A}$ , requires at most  $q^2$  queries, and the circuits in each query to  $\text{Decide}_{\mathfrak{G}}^\Psi$  consist of at most  $q^2$  queries. Thus, given a  $2^{n/12}$ -query algorithm  $\mathcal{A}$ , the resulting  $\mathcal{B}$  is a

$2^{n/6}$ -query algorithm, and we deduce Theorem 4.4 (or rather, the strengthened version with fixed  $\mathcal{O}$ ) by applying Lemma 4.5

$$\Pr_{\substack{f \leftarrow \mathfrak{F} \\ \Psi = (f, \mathcal{O}, \text{Eval}^{f, \mathcal{O}}) \\ x \leftarrow \{0,1\}^n}} \left[ \mathcal{A}^{\Psi, \text{Decide}_{\mathfrak{S}}^{\Psi}}(f(x)) = x \right] = \Pr_{\substack{f \leftarrow \mathfrak{F} \\ x \leftarrow \{0,1\}^n}} \left[ \mathcal{B}^{f, \text{Decide}_{\mathfrak{S}}^f}(f(x)) = x \right] \leq O(2^{-n/2}).$$

■

#### 4.4 The Existence of an Indistinguishability Obfuscator

In this section we prove that relative to the oracle  $\Gamma = (\Psi, \text{Decide}_{\mathfrak{S}}^{\Psi})$ , where  $\Psi = (f, \mathcal{O}, \text{Eval}^{f, \mathcal{O}})$  is sampled from the distribution  $\mathcal{D}(\mathfrak{S})$  defined in Section 4.2, there exists an indistinguishability obfuscator  $i\mathcal{O}$  for all circuits with  $f$ -gates.

Our obfuscator is based on those of Asharov and Segev [AS15] and Bitansky et al. [BDV17] but has a somewhat different structure. Similarly to their obfuscator, for every  $n \in \mathbb{N}$ , given an oracle-aided circuit  $C \in \{0,1\}^n$ , the obfuscator  $i\mathcal{O}$  samples  $r \leftarrow \{0,1\}^n$  and outputs the obfuscated circuit  $\widehat{C} = \mathcal{O}_n(C, r) \in \{0,1\}^{10n}$ . In turn, the oracle  $\text{Eval}^{f, \mathcal{O}}$  can be used for evaluating such an obfuscated circuit at any given point  $\alpha$ : If there exists a pair  $(C, r) \in \{0,1\}^n \times \{0,1\}^n$  such that  $\widehat{C} = \mathcal{O}_n(C, r)$  then  $\text{Eval}^{f, \mathcal{O}}$  outputs  $C^f(\alpha)$  and otherwise it outputs  $\perp$ .

However, unlike their obfuscator of Asharov and Segev [AS15] and Bitansky et al. [BDV17], which was sampled uniformly at random among all injective functions (of the appropriate input and output lengths), recall that according to our definition of the distribution  $\mathcal{D}(\mathfrak{S})$  it holds that  $\mathcal{O}_n(C, r) = (h(r), g_r(C))$ , where the function  $h$  is uniformly-chosen from the set  $\text{InjFunc}_n^{5n}$ , and for every  $r \in \{0,1\}^n$  a function  $g_r$  is uniformly-chosen from the set  $\text{InjFunc}_n^{5n}$ .

Recall that we call an oracle-aided algorithm  $\mathcal{A}$  a  $q$ -query algorithm, for a function  $q = q(n)$ , if when given any input  $x \in \{0,1\}^n$  it issues at most  $q(n)$  queries to the oracle  $\Gamma$ , each of its queries to  $\text{Eval}$  and  $\text{Decide}$  consists of circuits with at most  $q(n)$  oracle gates, and the number of communication rounds in the proof systems corresponding to each of its queries to  $\text{Decide}$  is at most  $q(n)$ . Letting  $\mathcal{C}$  denote the class of all oracle-aided circuit with  $f$ -gates, we prove the following theorem:

**Theorem 4.9.** *For any oracle-aided  $2^{n/6}$ -query algorithm  $\mathcal{A}$  it holds that*

$$\mathbb{E}_{\Gamma} \left| \Pr \left[ \text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{A}, \mathcal{C}}^{\text{iO}}(n) = 1 \right] - \frac{1}{2} \right| \leq O(2^{-n/4})$$

where the expectation is taken over the choice of  $\Gamma = (\Psi, \text{Decide}_{\mathfrak{S}}^{\Psi})$  where  $\Psi \leftarrow \mathcal{D}(\mathfrak{S})$ , and the inner probability is taken over the randomness of the experiment  $\text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{A}, \mathcal{C}}^{\text{iO}}(n)$ .

Toward proving Theorem 4.9, we first prove the following lemma.

**Lemma 4.10.** *For any oracle-aided  $4 \cdot 2^{n/6}$ -query algorithm  $\mathcal{A}$  it holds that*

$$\left| \Pr \left[ \text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{A}, \mathcal{C}}^{\text{iO}}(n) = 1 \right] - \frac{1}{2} \right| \leq O(2^{-n/2})$$

where the probability is taken both over the choice of  $\Gamma = (\Psi, \text{Decide}_{\mathfrak{S}}^{\Psi})$  where  $\Psi \leftarrow \mathcal{D}(\mathfrak{S})$ , and over the randomness of the experiment  $\text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{A}, \mathcal{C}}^{\text{iO}}(n)$ .

**Proof.** We prove that the lemma holds when even fixing the oracle  $f$  and  $\mathcal{O}_{-n} = \{\mathcal{O}_k\}_{k \neq n}$ , and only sampling  $\mathcal{O}_n$ . We introduce a sequence of 5 hybrid experiments such that the first hybrid

experiment  $\mathcal{H}_1$  is the real indistinguishability-obfuscation experiment  $\text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{A}, \mathcal{C}}^{\text{IO}}(n)$  and the last hybrid experiment  $\mathcal{H}_5$  is an experiment in which the advantage of the adversary is 0. Then, by upper bounding the difference in the advantage between each pair of consecutive hybrid experiments we deduce our lemma.

In what follows we first describe the hybrid experiments (see also the table below for a summary – where we omit the function  $f$  since it has been fixed), and then present a sequence of claims for bounding the differences in the advantages.

**The hybrid  $\mathcal{H}_1$ .** This is the real experiment in which we sample  $\mathcal{O}_n$  by sampling  $h \leftarrow \text{InjFunc}_n^{5n}$ , sampling  $g_r \leftarrow \text{InjFunc}_n^{5n}$  for every  $r \in \{0, 1\}^n$ , and setting  $\mathcal{O}_n(C, r) = (h(r), g_r(C))$ .

**The hybrid  $\mathcal{H}_2$ .** In this experiment, instead of giving the pre-challenge adversary  $\mathcal{A}_0$  access to the oracle  $\Gamma = (\Psi, \text{Decide}_{\mathfrak{S}}^{\Psi})$  where  $\Psi = (f, \mathcal{O}, \text{Eval}^{\mathcal{O}})$ , we sample a string  $\hat{h} \leftarrow \{0, 1\}^{5n} \setminus \text{Image}(h)$  and a function  $\hat{g} \leftarrow \text{InjFunc}_n^{5n}$ , then we give  $\mathcal{A}_0$  access to the oracle  $\Gamma' = (\Psi', \text{Decide}_{\mathfrak{S}}^{\Psi'})$  where  $\Psi' = (f, \mathcal{O}_{(\cdot, r^*) \rightarrow (\hat{h}, \hat{g}(\cdot))}, \text{Eval}^{\mathcal{O}})$  and for every  $C, r \in \{0, 1\}^n$  we define

$$\mathcal{O}_{(\cdot, r^*) \rightarrow (\hat{h}, \hat{g}(\cdot))}(C, r) = \begin{cases} (\hat{h}, \hat{g}(C)) & \text{if } r = r^* \\ \mathcal{O}(C, r) & \text{otherwise} \end{cases}.$$

That is, for the challenge randomness  $r^*$ , instead of obfuscating using  $h(r^*)$  and  $g_{r^*}(\cdot)$  we use our “planted obfuscation”  $\hat{h}$  and  $\hat{g}(\cdot)$ . The rest of the experiment proceeds as before.

**The hybrid  $\mathcal{H}_3$ .** In this experiment, we return to giving the pre-challenge adversary  $\mathcal{A}_0$  access to the real oracle  $\Gamma$ . However, we now give the post-challenge adversary  $\mathcal{A}_1$  a “planted challenge”  $(\hat{h}, \hat{g}(C_b))$ , and we give  $\mathcal{A}_1$  access to the oracle  $\Gamma' = (\Psi', \text{Decide}_{\mathfrak{S}}^{\Psi'})$  where  $\Psi' = (f, \mathcal{O}_{(\cdot, r^*) \rightarrow (\hat{h}, \hat{g}(\cdot))}, \text{Eval}^{\mathcal{O}})$ .

**The hybrid  $\mathcal{H}_4$ .** For an obfuscator function of the form  $\mathcal{O}(C, r) = (h(r), g_r(C))$ ,  $\hat{h} \in \{0, 1\}^{5n} \setminus \text{Image}(h)$  and  $\hat{g} \in \text{InjFunc}_n^{5n}$ , we define the planted evaluation function  $\text{PEval}_{(\hat{h}, \hat{g})}^{\mathcal{O}}$  as

$$\text{PEval}_{(\hat{h}, \hat{g})}^{\mathcal{O}}(\tilde{C}, \alpha) = \begin{cases} C^f(\alpha) & \text{if } \tilde{C} = (\hat{h}, \hat{g}(C)) \text{ for a circuit } C \in \{0, 1\}^n \\ \text{Eval}^{\mathcal{O}}(\tilde{C}, \alpha) & \text{otherwise} \end{cases}.$$

Note that since  $\hat{h} \notin \text{Image}(h)$ , it holds that  $\text{PEval}_{(\hat{h}, \hat{g})}^{\mathcal{O}}$  is a valid evaluation function and therefore  $(f, \mathcal{O}, \text{PEval}_{(\hat{h}, \hat{g})}^{\mathcal{O}}) \in \mathfrak{S}$ . Now, the experiment  $\mathcal{H}_4$  is obtained from  $\mathcal{H}_3$  by replacing the post-challenge oracle  $\Gamma'$  with the oracle  $\Gamma'' = (\Psi'', \text{Decide}_{\mathfrak{S}}^{\Psi''})$  where  $\Psi'' = (f, \mathcal{O}, \text{PEval}_{(\hat{h}, \hat{g})}^{\mathcal{O}})$ . Note that in this experiment, the randomness  $r^*$  has no role.

**The hybrid  $\mathcal{H}_5$ .** This experiment is obtained from  $\mathcal{H}_4$  by replacing the challenge obfuscation  $(\hat{h}, \hat{g}(C_b))$  with  $(\hat{h}, \hat{g}(C_0))$ . Note that in this experiment, the bit  $b$  has no role except for the winning condition, namely,  $\mathcal{A}$  wins if  $\mathcal{A}_1$  outputs  $b$ .

Hybrid	$\mathcal{H}_1$	$\mathcal{H}_2$	$\mathcal{H}_3$	$\mathcal{H}_4$	$\mathcal{H}_5$
<b>Challenger Randomness</b>	$b \leftarrow \{0, 1\}, r^* \leftarrow \{0, 1\}^n$			$b \leftarrow \{0, 1\}$	
<b>Obfuscator Fuction</b>	$\mathcal{O}_n(C, r) = (h(r), g_r(C))$ , where $h \leftarrow \text{InjFunc}_n^{5n}$ and $g_r \leftarrow \text{InjFunc}_n^{5n}$ for every $r \in \{0, 1\}^n$				
<b>Planted Obfuscation</b>	N/A	$\hat{h} \leftarrow \{0, 1\}^{5n} \setminus \text{Image}(h), \hat{g} \leftarrow \text{InjFunc}_n^{5n}$			
<b>Pre-challenge Oracle</b>	$\Psi = (\mathcal{O}, \text{Eval}^\mathcal{O})$ $\text{Decide}_{\mathfrak{S}}^\Psi$	$\Psi' = (\mathcal{O}' = \mathcal{O}_{(\cdot, r^*) \rightarrow (\hat{h}, \hat{g}(\cdot))}, \text{Eval}^{\mathcal{O}'})$ , $\text{Decide}_{\mathfrak{S}}^{\Psi'}$	$\Psi = (\mathcal{O}, \text{Eval}^\mathcal{O})$ $\text{Decide}_{\mathfrak{S}}^\Psi$		
<b>Challenge Obfuscation</b>	$\mathcal{O}(C_b, r^*) = (h(r^*), g_{r^*}(C_b))$		$(\hat{h}, \hat{g}(C_b))$	$(\hat{h}, \hat{g}(C_0))$	
<b>Post-challenge Oracle</b>	$\Psi = (\mathcal{O}, \text{Eval}^\mathcal{O})$ $\text{Decide}_{\mathfrak{S}}^\Psi$		$\Psi' = (\mathcal{O}' = \mathcal{O}_{(\cdot, r^*) \rightarrow (\hat{h}, \hat{g}(\cdot))}, \text{Eval}^{\mathcal{O}'})$ , $\text{Decide}_{\mathfrak{S}}^{\Psi'}$	$\Psi'' = (\mathcal{O}, \text{PEval}_{(\hat{h}, \hat{g})}^\mathcal{O})$ $\text{Decide}_{\mathfrak{S}}^{\Psi''}$	
<b>Winning Condition</b>	$\mathcal{A}_1$ outputs $b$				

**Claim 4.11.**  $|\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_1] - \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_2]| \leq 27q^3/2^n$ .

**Proof.** We observe that if  $\mathcal{A}_0$  has the same output both when given access to  $\Gamma$  or  $\Gamma'$ , then the rest of the experiment proceeds the same way and  $\mathcal{A}$  wins in  $\mathcal{H}_1$  if and only if he wins in  $\mathcal{H}_2$ . Hence,

$$|\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_1] - \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_2]| \leq \Pr[\mathcal{A}_0^\Gamma(1^n) \neq \mathcal{A}_0^{\Gamma'}(1^n)].$$

We define a label function  $\text{lab} : \mathcal{Q} \rightarrow \mathcal{X}$ , where the label of a query  $(C, r) \in \{0, 1\}^n \times \{0, 1\}^n$  to  $\mathcal{O}_n$  is  $r$ , the label of a query  $((\hat{h}, \hat{g}), \alpha) \in (\{0, 1\}^{5n} \times \{0, 1\}^{5n}) \times \{0, 1\}^n$  to  $E_n$  is  $h$ , and the label of any other query is  $\perp$ .

We observe that the view of  $\mathcal{A}_0$  in  $\mathcal{H}_1$  is independent of the choice of  $r^*$ . Also, the view of  $\mathcal{A}_0$  does not provide any information about  $\hat{h}$  apart from being outside the image of  $h$ . Therefore, if a query to  $\mathcal{O}$  is made, then the probability of its label to be  $r^*$  is at most  $1/2^n$ . In any other case, the answer to this query is the same in  $\mathcal{H}_1$  and  $\mathcal{H}_2$ , and both executions proceed the same way. Also, if a query to  $\text{Eval}^\mathcal{O}$  is made, then the probability of its label to be  $h(r^*)$  or  $\hat{h}$  is at most  $1/2^n + 1/(5^n - 2^n)$ . In any other case, the answer to this query is the same in  $\mathcal{H}_1$  and  $\mathcal{H}_2$ , and both executions proceed the same way. Finally, if a query  $(C_0, C_1, 1^{\ell_r}, 1^k)$  to  $\text{Decide}_{\mathfrak{S}}^\Psi$  is made, then we apply Lemma 4.3. The set  $\mathbf{I} = \mathbf{I}(\mathfrak{S}, \Psi, C_0, C_1, \ell_r, k, \text{lab}, 3) \subseteq \mathcal{X}$  of influential labels is independent of  $r^*$  and  $\hat{h}$  (subject to  $\hat{h}$  being outside the image of  $h$ ). Therefore, the probability of  $\mathbf{I}$  containing  $r^*$ ,  $h(r^*)$  or  $\hat{h}$  is at most  $3|I|/2^n \leq 27q^2/2^n$ . In any other case, the oracle  $\Psi'$  is obtained from  $\Psi$  by changes involving 3 labels outside of  $\mathbf{I}$ , and therefore by Lemma 4.3 it holds that

$$\text{Decide}_{\mathfrak{S}}^{\Psi'}(C_0, C_1, 1^{\ell_k}, 1^k) = \text{Decide}_{\mathfrak{S}}^\Psi(C_0, C_1, 1^{\ell_k}, 1^k),$$

and both executions proceed the same way. Applying a union bound we deduce that

$$\begin{aligned} \Pr[\mathcal{A}_0^\Gamma(1^n) \neq \mathcal{A}_0^{\Gamma'}(1^n)] &\leq \Pr[\text{A query was answered differently in } \mathcal{A}_0^\Gamma(1^n) \text{ and } \mathcal{A}_0^{\Gamma'}(1^n)] \\ &\leq \frac{27q^3}{2^n}, \end{aligned}$$

and the claim follows. ■

**Claim 4.12.**  $\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_2] = \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_3]$ .

**Proof.** We couple the experiments  $\mathcal{H}_2$  and  $\mathcal{H}_3$  as follows. First, we sample the same  $r^* \leftarrow \{0, 1\}^n$  and  $b \leftarrow \{0, 1\}$  for both experiments. Then, we uniformly sample a random injective function  $h : \{0, 1\}^n \setminus \{r^*\} \rightarrow \{0, 1\}^{5n}$  and sample  $g_r \leftarrow \text{InjFunc}_n^{5n}$  for every  $r \in \{0, 1\}^n \setminus \{r^*\}$ . Next, we sample distinct  $\widehat{h}, \widehat{h}' \leftarrow \{0, 1\} \setminus \text{Image}(h)$  and sample two injective functions  $\widehat{g}, \widehat{g}' \leftarrow \text{InjFunc}_n^{5n}$ . Now, in  $\mathcal{H}_2$  we let the obfuscator function be

$$\mathcal{O}(C, r) = \begin{cases} (\widehat{h}, \widehat{g}(C)) & \text{if } r = r^* \\ (h(r), g_r(C)) & \text{otherwise} \end{cases},$$

and let  $\widehat{h}', \widehat{g}'$  be the planted obfuscation, whereas in  $\mathcal{H}_3$  we let the obfuscator function be

$$\mathcal{O}(C, r) = \begin{cases} (\widehat{h}', \widehat{g}'(C)) & \text{if } r = r^* \\ (h(r), g_r(C)) & \text{otherwise} \end{cases},$$

and let  $\widehat{h}, \widehat{g}$  be the planted obfuscation. It is easy to see that the marginal distribution in both experiments is correct, and that both experiments are identical, thus the claim follows.  $\blacksquare$

**Claim 4.13.**  $|\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_3] - \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_4]| \leq 12q^3/2^n$ .

**Proof.** Here we use an argument similar to that of Claim 4.11. If the queries of  $\mathcal{A}_1$  are answered in the same way both in  $\mathcal{H}_3$  and  $\mathcal{H}_4$ , then  $\mathcal{A}_1$  outputs the same guess and wins with the same probability. Therefore, it suffices to bound the probability that some query was answered differently.

We define the same label function  $\text{lab} : \mathcal{Q} \rightarrow \mathcal{X}$  as in Claim 4.11. That is, the label of a query  $(C, r) \in \{0, 1\}^n \times \{0, 1\}^n$  to  $\mathcal{O}_n$  is  $r$ , the label of a query  $((\widehat{h}, \widehat{g}), \alpha) \in (\{0, 1\}^{5n} \times \{0, 1\}^{5n}) \times \{0, 1\}^n$  to  $E_n$  is  $\widehat{h}$ , and the label of any other query is  $\perp$ .

We observe that the view of  $\mathcal{A}_1$  in  $\mathcal{H}_4$  is independent of the choice of  $r^*$ . In fact,  $r^*$  does not have any role in that experiment. Therefore, if a query to  $\mathcal{O}$  is made, then the probability of its label to be  $r^*$  is at most  $1/2^n$ . In any other case, the answer to this query is the same in  $\mathcal{H}_3$  and  $\mathcal{H}_4$ , and both executions proceed the same way. Also, if a query to  $\text{PEval}_{(\widehat{h}, \widehat{g})}^{\mathcal{O}}$  is made (which corresponds to a query to  $\text{Eval}^{\mathcal{O}}$  in  $\mathcal{H}_3$ ), then the probability of its label to be  $h(r^*)$  is at most  $1/2^n$ . In any other case, the answer to this query is the same in  $\mathcal{H}_3$  and  $\mathcal{H}_4$ , and both executions proceed the same way. Finally, if a query  $(C_0, C_1, 1^{\ell_r}, 1^k)$  to  $\text{Decide}_{\mathfrak{S}}^{\Psi''}$  is made, then we apply Lemma 4.3. The set  $\mathbf{I} = \mathbf{I}(\mathfrak{S}, \Psi'', C_0, C_1, \ell_r, k, \text{lab}, 2) \subseteq \mathcal{X}$  of influential labels is independent of  $r^*$ . Therefore, the probability that  $\mathbf{I}$  contains  $r^*$  or  $h(r^*)$  is at most  $2|I|/2^n \leq 12q^2/2^n$ . In any other case, the oracle  $\Psi'$  is obtained from  $\Psi''$  by changes involving two labels outside of  $\mathbf{I}$ , and therefore by Lemma 4.3 it holds that

$$\text{Decide}_{\mathfrak{S}}^{\Psi'}(C_0, C_1, 1^{\ell_k}, 1^k) = \text{Decide}_{\mathfrak{S}}^{\Psi''}(C_0, C_1, 1^{\ell_k}, 1^k),$$

and both executions proceed the same way. Applying a union bound we deduce that

$$\Pr \left[ \text{A query was answered differently in } \mathcal{A}_1^{\Gamma''}(\widehat{h}, \widehat{g}(C_b)) \text{ and } \mathcal{A}_1^{\Gamma'}(\widehat{h}, \widehat{g}(C_b)) \right] \leq \frac{12q^3}{2^n},$$

and the claim follows.  $\blacksquare$

**Claim 4.14.**  $\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_4] = \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_5]$ .

**Proof.** Let  $\widehat{g} \leftarrow \text{InjFunc}_n^{5n}$ ,  $b \leftarrow \{0, 1\}$ , and let

$$\widehat{g}'(C) = \begin{cases} \widehat{g}(C_{1-\sigma}) & b = 1 \text{ and } C = C_\sigma \text{ for some } \sigma \in \{0, 1\} \\ \widehat{g}(C) & \text{otherwise} \end{cases}.$$

Then,  $\widehat{g}'$  is also uniformly-distributed in  $\text{InjFunc}_n^{5n}$ . Also, since  $C_0^f \equiv C_1^f$  it holds that

$$\text{PEval}_{(\widehat{h}, \widehat{g})}^{\mathcal{O}} \equiv \text{PEval}_{(\widehat{h}, \widehat{g}')}^{\mathcal{O}},$$

and thus if we replace  $\widehat{g}$  with  $\widehat{g}'$  the entire oracle  $\Psi''$  stays the same. So if we couple the experiments  $\mathcal{H}_4$  and  $\mathcal{H}_5$  by using the same randomness expect for replacing  $\widehat{g}$  in  $\mathcal{H}_5$  with  $\widehat{g}'$ , then we obtain the same challenge in both experiments  $(\widehat{h}, \widehat{g}(C_b)) = (\widehat{h}, \widehat{g}'(C_0))$  and the same oracle  $\Psi''$ . As a result, both experiments proceed the same way, thus the claim follows. ■

**Claim 4.15.**  $\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_5] = 1/2$ .

**Proof.** In this experiment the view of  $\mathcal{A}$  is independent of  $b$ . ■

Now we turn back to proving Lemma 4.10. It holds that

$$\begin{aligned} |\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_1] - \frac{1}{2}| &= |\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_1] - \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_5]| \\ &\leq \sum_{i=1}^4 |\Pr[\mathcal{A} \text{ wins in } \mathcal{H}_i] - \Pr[\mathcal{A} \text{ wins in } \mathcal{H}_{i+1}]| \leq \frac{40q^3}{2^n}, \end{aligned}$$

and by plugging  $q = 4 \cdot 2^{n/6}$  we obtain Lemma 4.10. ■

Lastly, we show how to deduce Theorem 4.9 from Lemma 4.10. In what follows, for an event  $\mathcal{E}$  we denote by  $\Pr_{\Gamma}[\mathcal{E}]$  its probability over the choice of  $\Gamma \leftarrow \mathcal{D}(\mathfrak{S})$ , and by  $\Pr_{\text{Exp}}[\mathcal{E}]$  its probability over the randomness of the indistinguishability-obfuscation experiment. Thus, if we write  $\Pr_{\text{Exp}}[\dots]$  then  $\Gamma$  is already fixed in the experiment, and if we write  $\Pr_{\Gamma, \text{Exp}}[\dots]$  then  $\Gamma$  is sampled for the experiment.

**Proof of Theorem 4.9.** Let  $\mathcal{A} = (\mathcal{A}_0, \mathcal{A}_1)$  be a  $q(n)$ -query algorithm, where  $q(n) \leq 2^{n/6}$ . Our goal is to prove that

$$\mathbb{E}_{\Gamma} \left| \Pr_{\text{Exp}} \left[ \text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{A}, \mathcal{C}}^{\text{IO}}(n) = 1 \right] - \frac{1}{2} \right| \leq O(2^{-n/4})$$

where the expectation is taken over the choice of  $\Gamma = (\Psi, \text{Decide}_{\mathfrak{S}}^{\Psi})$  where  $\Psi \leftarrow \mathcal{D}(\mathfrak{S})$ , and the inner probability is taken over the randomness of the experiment  $\text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{A}, \mathcal{C}}^{\text{IO}}(n)$ . Our proof transforms  $\mathcal{A}$  into a  $4q(n)$ -query algorithm  $\mathcal{B} = (\mathcal{B}_0, \mathcal{B}_1)$ , where for every oracle  $\Gamma$  it holds that  $\Pr_{\text{Exp}}[\text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{B}, \mathcal{C}}^{\text{IO}}(n) = 1] \geq 1/2$  (and the advantage itself is polynomially related to that of  $\mathcal{A}$ ). We note that ideas and transformations along these lines are quite common, and we refer the reader to the work of Brakerski and Goldreich for similar results in a more general context [BG11].

First, since  $\mathcal{A}$  is computationally unbounded, we may assume without loss of generality that it is deterministic (e.g., by fixing the randomness that maximizes its expected advantage). As discussed in Section 2, we can further assume that  $\mathcal{A}$  is stateless (e.g., by letting  $\mathcal{A}_1$  recompute  $\mathcal{A}_0^{\Gamma}(1^n)$  at the cost of at most  $q$  additional queries).

Now, we let  $\mathcal{B}_0 = \mathcal{A}_0$ , and given an obfuscation  $\widehat{C}$  as input and oracle access to  $\Gamma$ , the algorithm  $\mathcal{B}_1$  is defined as follows:

1. Compute  $\sigma = \mathcal{A}_1^{\Gamma}(\widehat{C}) \in \{0, 1\}$ .
2. Sample  $\tilde{b} \leftarrow \{0, 1\}$  and  $\tilde{r}^* \leftarrow \{0, 1\}^n$ , and then compute  $\tilde{\sigma} = \mathcal{A}_1^{\Gamma}(\mathcal{O}(C_{\tilde{b}}, \tilde{r}^*)) \in \{0, 1\}$ .
3. Output  $\sigma \oplus \tilde{b} \oplus \tilde{\sigma}$ .

Let  $\Pr_{\text{Exp}}[\text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{A}, c}^{\text{iO}}(n) = 1] = 1/2 + \epsilon$ , where  $\epsilon \in [-1/2, 1/2]$ . Then,

$$\begin{aligned} \Pr_{\text{Exp}}[\text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{B}, c}^{\text{iO}}(n) = 1] &= \Pr_{\text{Exp}}[\sigma \oplus \tilde{b} \oplus \tilde{\sigma} = b] \\ &= \Pr_{\text{Exp}}[\sigma = b \wedge \tilde{\sigma} = \tilde{b}] + \Pr_{\text{Exp}}[\sigma \neq b \wedge \tilde{\sigma} \neq \tilde{b}] \\ &= \left(\frac{1}{2} + \epsilon\right)^2 + \left(\frac{1}{2} - \epsilon\right)^2 = \frac{1}{2} + \epsilon^2. \end{aligned}$$

So we conclude that for any oracle  $\Gamma$  it holds that  $\Pr_{\text{Exp}}[\text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{B}, c}^{\text{iO}}(n) = 1] \geq 1/2$  and that

$$\left| \Pr_{\text{Exp}}[\text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{A}, c}^{\text{iO}}(n) = 1] - \frac{1}{2} \right| = \left( \Pr_{\text{Exp}}[\text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{B}, c}^{\text{iO}}(n) = 1] - \frac{1}{2} \right)^{1/2}.$$

Noting that  $\mathcal{B}$  is a  $4q(n)$ -query algorithm where  $q(n) = 2^{n/6}$ , by applying Lemma 4.10 to  $\mathcal{B}$  we obtain

$$\Pr_{\Gamma, \text{Exp}} \left[ \text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{B}, c}^{\text{iO}}(n) = 1 \right] - \frac{1}{2} = \left| \Pr_{\Gamma, \text{Exp}} \left[ \text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{B}, c}^{\text{iO}}(n) = 1 \right] - \frac{1}{2} \right| \leq O(2^{n/2}).$$

Finally, Jensen's inequality settles the proof of Theorem 4.9 as follows

$$\begin{aligned} &\mathbb{E}_{\Gamma} \left| \Pr_{\text{Exp}} \left[ \text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{A}, c}^{\text{iO}}(n) = 1 \right] - \frac{1}{2} \right| \\ &= \mathbb{E}_{\Gamma} \left( \Pr_{\text{Exp}} \left[ \text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{B}, c}^{\text{iO}}(n) = 1 \right] - \frac{1}{2} \right)^{1/2} \\ &\leq \left( \mathbb{E}_{\Gamma} \left( \Pr_{\text{Exp}} \left[ \text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{B}, c}^{\text{iO}}(n) = 1 \right] - \frac{1}{2} \right) \right)^{1/2} \\ &= \left( \Pr_{\Gamma, \text{Exp}} \left[ \text{Exp}_{\Gamma, i\mathcal{O}, \mathcal{B}, c}^{\text{iO}}(n) = 1 \right] - \frac{1}{2} \right)^{1/2} \leq O(2^{n/4}). \end{aligned}$$

■

## 4.5 Putting it All Together

Given Theorems 4.4 and 4.9 we can now derive Theorem 4.1.

**Proof of Thm. 4.1.** Let  $((V, \bar{V}), \ell_r, k, M, T_M, \epsilon_{M,1}, \epsilon_{M,2})$  be a fully black-box construction of a worst-case hard (MIP, coMIP) protocol pair from an injective one-way function  $f$  and an indistinguishability obfuscator for all oracle-aided circuits  $C^f$ . Lemma 4.2 guarantees the existence of a polynomial-time single-query algorithm  $\mathcal{A}$  such that for every  $\Psi = (f, \mathcal{O}, \text{Eval}^{f, \mathcal{O}})$  in the support of the distribution  $\mathcal{D}(\mathfrak{S})$ , the algorithm  $\mathcal{A}$  with oracle access to  $\Gamma = (\Psi, \text{Decide}_{\mathfrak{S}}^{\Psi})$  decides the language  $L^{\Psi} \subseteq \{0, 1\}^*$  defined by  $(V, \bar{V}, \ell_r, k)$  relative to  $\Psi$ . That is, for every  $x \in \{0, 1\}^*$  it holds that  $\mathcal{A}^{\Gamma}(x) = \chi_{L^{\Psi}}(x)$ . For every  $n \in \mathbb{N}$ , denote by  $T_{\mathcal{A}}(n)$  the polynomial running time of  $\mathcal{A}$  on inputs of length  $n$ .

Definition 3.3 then states that there are two possible cases to consider:  $\mathcal{A}$  can be used either for inverting the injective one-way function  $f$ , or for breaking the security of the indistinguishability obfuscator  $i\mathcal{O}$ . Specifically, in the first case we obtain from Definition 3.3 that for every  $\Psi = (f, \mathcal{O}, \text{Eval}^{f, \mathcal{O}})$  in the support of the distribution  $\mathcal{D}(\mathfrak{S})$  it holds that

$$\Pr [M^{\Gamma, \mathcal{A}}(f(x)) = x] \geq \epsilon_{M,1}(T_{\mathcal{A}}(n)) \cdot \epsilon_{M,2}(n)$$

for infinitely many values of  $n \in \mathbb{N}$ , where  $\Gamma = (\Psi, \text{Decide}_{\mathfrak{G}}^{\Psi})$  and the probability is taken over the choice of  $x \leftarrow \{0, 1\}^n$  and over the internal randomness of  $M$ . The algorithm  $M$  may invoke  $\mathcal{A}$  on various input lengths (i.e., in general  $M$  is not restricted to invoking  $\mathcal{A}$  only on input length  $n$ ), and we denote by  $\ell(n)$  the maximal input length on which  $M$  invokes  $\mathcal{A}$  (when  $M$  itself is invoked on input  $f(x)$  for  $x \in \{0, 1\}^n$ ). Thus, viewing  $M^{\mathcal{A}}$  as a single oracle-aided algorithm that has access to  $\Gamma$ , its running time  $T_{M^{\mathcal{A}}}(n)$  satisfies  $T_{M^{\mathcal{A}}}(n) \leq T_M(n) \cdot T_{\mathcal{A}}(\ell(n))$  (this follows since  $M$  may invoke  $\mathcal{A}$  at most  $T_M(n)$  times, and the running time of  $\mathcal{A}$  on each such invocation is at most  $T_{\mathcal{A}}(\ell(n))$ ). In particular, viewing  $M' = M^{\mathcal{A}}$  as a single oracle-aided algorithm that has oracle access to  $\Gamma$ , implies that  $M'$  is a  $q$ -query algorithm where  $q(n) = T_{M^{\mathcal{A}}}(n)$ . This holds for any  $\Psi$  in the support of the distribution  $\mathcal{D}(\mathfrak{G})$ , and given that  $q(n)$  is polynomial in  $n$  then Theorem 4.4 guarantees that  $\epsilon_{M,1}(T_{\mathcal{A}}(n)) \cdot \epsilon_{M,2}(n) \leq O(2^{-n/2})$ .

In the second case we obtain from Definition 3.3 that for every  $\Psi = (f, \mathcal{O}, \text{Eval}^{f, \mathcal{O}})$  in the support of the distribution  $\mathcal{D}(\mathfrak{G})$  it holds that

$$\left| \Pr \left[ \text{Exp}_{\Gamma, i, \mathcal{O}, M^{\mathcal{A}}, \mathcal{C}}^{\text{iO}}(n) = 1 \right] - \frac{1}{2} \right| \geq \epsilon_{M,1}(T_{\mathcal{A}}(n)) \cdot \epsilon_{M,2}(n)$$

for infinitely many values of  $n \in \mathbb{N}$ , where the probability is taken over the randomness of the experiment  $\text{Exp}_{\Gamma, i, \mathcal{O}, M^{\mathcal{A}}, \mathcal{C}}^{\text{iO}}(n)$ . The same reasoning applied to the first case, together with Theorem 4.9 guarantee that  $\epsilon_{M,1}(T_{\mathcal{A}}(n)) \cdot \epsilon_{M,2}(n) \leq O(2^{-n/4})$ .

We conclude the proof noting that the algorithm  $\mathcal{A}$  provided by Lemma 4.2 runs in fact in linear time. That is,  $T_{\mathcal{A}}(n) = O(n)$ , and thus from the above two cases we obtain  $\epsilon_{M,1}(n) \cdot \epsilon_{M,2}(n) \leq 2^{-\Omega(n)}$ .

■

## 5 Impossibility Result for Constructions based on RO-Implied Primitives

Equipped with a formal definition of a black-box construction of a hard (MIP, coMIP) protocol pair from a collision-resistant hash function (recall Definition 3.4), in this section we prove the following theorem:

**Theorem 5.1.** *Let  $((V, \bar{V}), \ell_r, k, M, T_M, \epsilon_{M,1}, \epsilon_{M,2})$  be a fully black-box construction of a worst-case hard (MIP, coMIP) protocol pair from a collision-resistant hash function. Then, there exists a constant  $c \geq 1$  such that*

$$\epsilon_{M,1}(n^c) \cdot \epsilon_{M,2}(n) \leq 2^{-\Omega(n)}.$$

*That is, at least one out of the adversary-dependent security loss  $\epsilon_{M,1}(\cdot)$  and the adversary-independent security loss  $\epsilon_{M,2}(\cdot)$  is exponential.*

We make use of the following standard lemma, which states that a uniformly-random ensemble  $f = \{f_n : \{0, 1\}^{2n} \rightarrow \{0, 1\}^{n-1}\}_{n \in \mathbb{N}}$  is an exponentially-hard collision-resistant hash function, for which we provide a proof for completeness.

**Lemma 5.2.** *For every  $q(n)$ -query algorithm  $M$  it holds that*

$$\Pr \left[ \left( x_1 \neq x_2 \right) \wedge \left( f_n(s, x_1) = f_n(s, x_2) \right) \right] \leq \frac{(q(n))^2 + 2}{2^n},$$

*where  $s \leftarrow \{0, 1\}^n$ ,  $f = \{f_n : \{0, 1\}^{2n} \rightarrow \{0, 1\}^{n-1}\}_{n \in \mathbb{N}}$  is uniformly sampled, and  $(x_1, x_2) \leftarrow M^f(s)$ .*

**Proof.** Fix any  $n \in \mathbb{N}$  and let  $q = q(n)$ . For any  $i \in [q]$  denote by  $\alpha_i$  the random variable corresponding to the  $i$ th distinct query in the computation  $(x_1, x_2) \leftarrow M^f(s)$ . Then,

$$\begin{aligned} & \Pr \left[ (x_1 \neq x_2) \wedge (f_n(s, x_1) = f_n(s, x_2)) \right] \\ & \leq \Pr \left[ (x_1 \neq x_2) \wedge (f_n(s, x_1) = f_n(s, x_2)) \wedge (x_1 \notin \{\alpha_1, \dots, \alpha_q\} \vee x_2 \notin \{\alpha_1, \dots, \alpha_q\}) \right] \end{aligned} \quad (5.1)$$

$$+ \Pr \left[ (x_1 \neq x_2) \wedge (f_n(s, x_1) = f_n(s, x_2)) \wedge (x_1, x_2 \in \{\alpha_1, \dots, \alpha_q\}) \right]. \quad (5.2)$$

For bounding the term in Eq. (5.1), note that since either  $x_1$  or  $x_2$  were not queried during the computation, then the view of  $M$  is completely independent of either  $f_n(s, x_1)$  or  $f_n(s, x_2)$ . Therefore,

$$\Pr \left[ (x_1 \neq x_2) \wedge (f_n(s, x_1) = f_n(s, x_2)) \wedge (x_1 \notin \{\alpha_1, \dots, \alpha_q\} \vee x_2 \notin \{\alpha_1, \dots, \alpha_q\}) \right] \leq \frac{1}{2^{n-1}}.$$

For bounding the term in Eq. (5.2), note that since  $f_n(s, x_1) = f_n(s, x_2)$  and since both  $x_1$  and  $x_2$  were queried, then there must exist  $i \in [q-1]$  for which  $f_n(s, \alpha_{i+1}) \in \{f_n(s, \alpha_1), \dots, f_n(s, \alpha_i)\}$ . The view of  $M$  up to the  $(i+1)$ st query is completely independent of  $f_n(s, \alpha_{i+1})$ , so we can bound

$$\begin{aligned} & \Pr \left[ (x_1 \neq x_2) \wedge (f_n(s, x_1) = f_n(s, x_2)) \wedge (x_1, x_2 \in \{\alpha_1, \dots, \alpha_q\}) \right] \\ & \leq \sum_{i=1}^{q-1} \Pr [f_n(s, \alpha_{i+1}) \in \{f_n(s, \alpha_1), \dots, f_n(s, \alpha_i)\}] \\ & \leq \sum_{i=1}^{q-1} \frac{i}{2^{n-1}} \leq \frac{q^2}{2^n}. \end{aligned}$$

■

## 5.1 Deciding (MIP, coMIP) Protocol Pair Languages

In this section we show how to efficiently decide any language that is defined via an oracle-aided (MIP, coMIP) protocol pair. We prove the following lemma:

**Lemma 5.3.** *Let  $(V, \bar{V})$  be a pair of oracle-aided polynomial-time algorithms that is an (MIP, coMIP) protocol pair, with respect to polynomials  $\ell_r(\cdot)$  and  $k(\cdot)$ , relative to every oracle  $f = \{f_n : \{0, 1\}^{2^n} \rightarrow \{0, 1\}^{n-1}\}_{n \in \mathbb{N}}$ . Then, there exists an algorithm  $\mathcal{A}$  that issues a polynomial number of queries, such that for every ensemble  $f = \{f_n : \{0, 1\}^{2^n} \rightarrow \{0, 1\}^{n-1}\}_{n \in \mathbb{N}}$ , given oracle access to  $f$  the algorithm  $\mathcal{A}$  decides the language  $L^f \subseteq \{0, 1\}^*$  defined by  $(V, \bar{V}, \ell_r, k)$  relative to  $f$ . That is, for every ensemble  $f$  and for every input  $x \in \{0, 1\}^*$ , the algorithm  $\mathcal{A}^f(x)$  outputs 1 if and only if  $x \in L^f$ . Moreover, the algorithm  $\mathcal{A}$  can be implemented to run in polynomial time given access to an RE-complete oracle.*

In the above lemma, RE-completeness is defined with respect to polynomial-time reductions (e.g., the halting problem of Turing machines). We note that in fact the claim holds for weaker oracles (e.g., an EXPSPACE-complete oracle), but we chose an RE-complete oracle to simplify the proof.

For convenience, similarly to the proof of Lemma 4.2 we convert the algorithms  $V$  and  $\bar{V}$  as considered in Lemma 5.3 to circuits as follows. Since  $V$  and  $\bar{V}$  run in polynomial time, there exists a polynomial  $p(n)$  such that on input of size  $n$  their output is of size at most  $p(n)$ . Then, similar to the proof of Lemma 4.2, given  $x \in \{0, 1\}^*$  as input, we consider the tuple  $(C_0, C_1, 1^{\ell_r(|x|)}, 1^{k(|x|)})$ , where  $C_0$  and  $C_1$  are the hardwired oracle-aided circuits  $\bar{V}(x, \cdot)$  and  $V(x, \cdot)$  respectively, the input size of both circuits is  $\lceil \log(k(|x|) + 1) \rceil + k(|x|)$  (where  $\lceil \log(k(|x|) + 1) \rceil$  bits are for the index of

the communication round and  $k(|x|)$  bits are for the messages of the prover) and the output size is  $p(|x| + \lceil \log(k(|x|) + 1) \rceil + k(|x|))$ .

Similar to Section 4.1 but adapted to our setting, the above tuple  $t = (C_0, C_1, 1^{\ell_r}, 1^k)$  is a *valid tuple* in the sense that for every ensemble  $g = \{g_n : \{0, 1\}^{2n} \rightarrow \{0, 1\}^{n-1}\}_{n \in \mathbb{N}}$  exactly one of the following two cases holds:

- There exists a function  $P_1 : \{0, 1\}^* \rightarrow \{0, 1\}$  such that

$$\Pr_{r \leftarrow \{0, 1\}^{\ell_r}} [\langle C_1^g(r), P_1 \rangle_k = 1] \geq 2/3 ,$$

and for every function  $P_0 : \{0, 1\}^* \rightarrow \{0, 1\}$  it holds that

$$\Pr_{r \leftarrow \{0, 1\}^{\ell_r}} [\langle C_0^g(r), P_0 \rangle_k = 1] \leq 1/3 .$$

In this case, we say that  $t^g = (C_0^g, C_1^g, 1^{\ell_r}, 1^k)$  is a *yes-instance*.

- There exists a function  $P_0 : \{0, 1\}^* \rightarrow \{0, 1\}$  such that

$$\Pr_{r \leftarrow \{0, 1\}^{\ell_r}} [\langle C_0^g(r), P_0 \rangle_k = 1] \geq 2/3 ,$$

and for every function  $P_1 : \{0, 1\}^* \rightarrow \{0, 1\}$  it holds that

$$\Pr_{r \leftarrow \{0, 1\}^{\ell_r}} [\langle C_1^g(r), P_1 \rangle_k = 1] \leq 1/3 .$$

In this case, we say that  $t^g = (C_0^g, C_1^g, 1^{\ell_r}, 1^k)$  is a *no-instance*.

Towards showing how to decide whether  $t^f$  is a yes or no instance, we borrow definitions and techniques from Blum and Impagliazzo [BI87], Impagliazzo and Naor [IN88] and Nisan [Nis89, Nis91], and adapt them to our setting. We define the block sensitivity and certificate complexity of  $t$  and prove upper bounds for them.

**Definition 5.4.** For ensembles  $f$  and  $g$ , we define their *disagreement set* as

$$\text{dis}(f, g) = \{x \in \{0, 1\}^* \mid f(x) \neq g(x)\} .$$

**Definition 5.5.** For a valid tuple  $t$  and a function  $f$  such that  $t^f$  is a yes-instance (respectively, no-instance), we define the *block sensitivity*  $\text{bs}(t, f)$  of  $t$  on  $f$  as the maximal number  $b$  such that there exist ensembles  $g_1, \dots, g_b$  such that the sets  $\text{dis}(f, g_1), \dots, \text{dis}(f, g_b)$  are pairwise disjoint and  $t^{g_1}, \dots, t^{g_b}$  are no-instances (respectively, yes-instances). We define the block sensitivity of  $t$  as  $\text{bs}(t) = \max_f \text{bs}(t, f)$ .

**Definition 5.6.** For a valid tuple  $t$ , an ensemble  $f$  such that  $t^f$  is a yes-instance (respectively, no-instance), and a set  $X$  of inputs, we say that  $X$  is a *certificate* for  $f$  if for every ensemble  $g$  such that  $\text{dis}(f, g) \cap X = \emptyset$  (i.e.,  $g$  agrees with  $f$  on all inputs from  $X$ ) it holds that  $t^g$  is a yes-instance (respectively, no-instance). We define the certificate complexity  $\mathbf{N}(t, f)$  of  $t$  on  $f$  to be the minimal number  $b$  such that there exists a certificate for  $f$  of cardinality  $b$ . We define the certificate complexity of  $t$  as  $\mathbf{N}(t) = \max_f \mathbf{N}(t, f)$ .

The following claim upper bounds the block sensitivity of a valid tuple.

**Claim 5.7.** For a valid tuple  $t = (C_0, C_1, 1^{\ell_r}, 1^k)$  it holds that

$$\text{bs}(t) \leq 3 \cdot k \cdot q ,$$

where  $q$  is a bound on the number of oracle gates in  $C_0$  and  $C_1$ .

**Proof.** Let  $f$  be an ensemble. We assume w.l.o.g. that  $t^f$  is a yes-instance, and let  $P_1$  such that

$$\Pr_{r \leftarrow \{0,1\}^{\ell_r}} \left[ \langle C_1^f(r), P_1 \rangle_k = 1 \right] \geq 2/3 ,$$

For  $b = \text{bs}(t, f)$ , let  $g_1, \dots, g_b$  such that the sets  $\text{dis}(f, g_1), \dots, \text{dis}(f, g_b)$  are pairwise disjoint and  $t^{g_1}, \dots, t^{g_b}$  are no-instances. For every  $i \in [b]$  it holds that

$$\Pr_{r \leftarrow \{0,1\}^{\ell_r}} \left[ \langle C_1^{g_i}(r), P_1 \rangle_k = 1 \right] \leq 1/3 ,$$

and since  $\langle C_1^f(r), P_1 \rangle_k = \langle C_1^{g_i}(r), P_1 \rangle_k$  unless  $C_1$  queries an input from  $\text{dis}(f, g_i)$ , it must hold that

$$\Pr_{r \leftarrow \{0,1\}^{\ell_r}} \left[ \langle C_1^f(r), P_1 \rangle_k \text{ queries } f \text{ on an input from } \text{dis}(f, g_i) \right] \geq 1/3 . \quad (5.3)$$

Note that fixing  $r$ ,  $\langle C_1^f(r), P_1 \rangle_k$  queries  $f$  at most  $k \cdot q$  times. Therefore,

$$\begin{aligned} & \sum_{i=1}^b \Pr_{r \leftarrow \{0,1\}^{\ell_r}} \left[ \langle C_1^f(r), P_1 \rangle_k \text{ queries } f \text{ on an input from } \text{dis}(f, g_i) \right] \\ & \leq \sum_{j=1}^{k \cdot q} \sum_{i=1}^b \Pr_{r \leftarrow \{0,1\}^{\ell_r}} \left[ \langle C_1^f(r), P_1 \rangle_k \text{ queries } f \text{ on an input from } \text{dis}(f, g_i) \text{ in its } j\text{th query} \right] \\ & = \sum_{j=1}^{k \cdot q} \Pr_{r \leftarrow \{0,1\}^{\ell_r}} \left[ \langle C_1^f(r), P_1 \rangle_k \text{ queries } f \text{ on an input from } \bigcup_{i=1}^b \text{dis}(f, g_i) \text{ in its } j\text{th query} \right] \\ & \leq k \cdot q \end{aligned}$$

Summing Eq. (5.3) over all  $i \in [b]$ , we obtain  $b/3 \leq k \cdot q$ , and thus  $\text{bs}(t) \leq 3 \cdot k \cdot q$ .  $\blacksquare$

The following claim upper bounds the certificate complexity of a valid tuple in terms of its block sensitivity.

**Claim 5.8.** For a valid tuple  $t$  it holds that

$$\mathbf{N}(t) \leq (\text{bs}(t))^2 .$$

**Proof.** Let  $f$  be an ensemble and assume without loss of generality that  $t^f$  is a yes-instance. For  $b = \text{bs}(t, f)$ , let  $g_1, \dots, g_b$  such that the sets  $\text{dis}(f, g_1), \dots, \text{dis}(f, g_b)$  are pairwise disjoint and  $t^{g_1}, \dots, t^{g_b}$  are no-instances. We may further assume that for every  $i \in [b]$  the set  $\text{dis}(f, g_i)$  is minimal with respect to containment, i.e., for every  $h$  such that  $\text{dis}(f, h) \subsetneq \text{dis}(f, g_i)$ ,  $t^h$  is a yes-instance. This can be done by replacing each  $g_i$  with a  $g'_i$  such that  $\text{dis}(f, g'_i) \subseteq \text{dis}(f, g_i)$  and  $g'_i$  is minimal in the above sense. For every  $i \in [b]$  we claim that  $|\text{dis}(f, g_i)| \leq \text{bs}(t, g_i)$ . It holds that  $\text{bs}(t, g_i) \geq 1$ , so if  $|\text{dis}(f, g_i)| = 1$  then the claim is trivial. Otherwise, let  $\text{dis}(f, g_i) = \{x_1, \dots, x_s\}$  where  $s \geq 2$  and define  $h_1, \dots, h_s$  by

$$h_j(x) = \begin{cases} f(x) & \text{if } x = x_j \\ g_i(x) & \text{otherwise} \end{cases} .$$

Then,  $\text{dis}(f, h_j) = \{x_j\} \subsetneq \text{dis}(f, g_i)$  and therefore by minimality  $t^{h_j}$  is a yes-instance. Also, since  $\text{dis}(g_i, h_j) = \{x_j\}$ , the sets  $\text{dis}(g_i, h_1), \dots, \text{dis}(g_i, h_s)$  are pairwise disjoint. Hence,  $s \leq \text{bs}(t, g_i)$  as claimed.

Finally, we show that  $X = \bigcup_{i \in [b]} \text{dis}(f, g_i)$  is a certificate for  $f$ . Assume towards contradiction that there exists an ensemble  $h$  such that  $t^h$  is a no-instance but  $\text{dis}(f, h) \cap X = \emptyset$ . Then  $g_1, \dots, g_b, h$  are  $b + 1$  ensembles such that the sets  $\text{dis}(f, g_1), \dots, \text{dis}(f, g_b), \text{dis}(f, h)$  are pairwise disjoint and  $t^{g_1}, \dots, t^{g_b}, t^h$  are no-instances in contradiction to the definition of  $b = \text{bs}(t, f)$ . Note that  $|X| \leq \sum_{i=1}^b \text{bs}(t, g_i) \leq (\text{bs}(t))^2$ . For every  $f$ , we showed the existence of a certificate for  $f$  with cardinality at most  $(\text{bs}(t))^2$  and therefore the claim follows.  $\blacksquare$

From Claim 5.7 and Claim 5.8 we can easily derive the following upper bound on the certificate complexity of a valid tuple.

**Corollary 5.9.** *For a valid tuple  $t = (C_0, C_1, 1^{\ell_r}, 1^k)$  it holds that*

$$\mathbf{N}(t) \leq 9 \cdot k^2 \cdot q^2,$$

where  $q$  is a bound on the number of oracle gates in  $C_0$  and  $C_1$ .

Now that we have shown a bound on the certificate complexity of  $t$ , we can show an algorithm that decides whether  $t^f$  is a yes-instance or no-instance.

**Claim 5.10.** *There exists an algorithm  $\mathcal{A}$  that given a valid tuple  $t = (C_0, C_1, 1^{\ell_r}, 1^k)$  as input and oracle access to an ensemble  $f$  determines whether  $t^f$  is a yes-instance or no-instance by issuing at most  $\text{poly}(q, k)$  queries to  $f$ , where  $q$  is the number of oracle gates in  $C_0$  and  $C_1$ . Moreover, given access to an RE-complete oracle,  $\mathcal{A}$  can be implemented to run in time  $\text{poly}(|t|)$ .*

**Proof.** In what follows, we interchangeably refer to a yes-instance as a 1-instance and to a no-instance as a 0-instance. Given a valid tuple  $t = (C_0, C_1, 1^{\ell_r}, 1^k)$  as input and oracle access to  $f$ , the algorithm  $\mathcal{A}$  is defined as follows:

1. Initialize an empty set  $I$  of query inputs and an empty set  $Q$  of query-response pairs.
2. Do for  $r = 9 \cdot k^2 \cdot q^2 + 1$  iterations:
  - (a) Do for every  $\sigma \in \{0, 1\}$ :

- i. Find an oracle  $g$  that is consistent with  $Q$  such that  $t^g$  is a  $\sigma$ -instance.
- ii. If no such  $g$  exists, then output  $1 - \sigma$  and terminate.
- iii. Find  $P$  such that

$$\Pr_{r \leftarrow \{0, 1\}^{\ell_r}} [\langle C_\sigma^g(r), P \rangle_k = 1] \geq 2/3.$$

- iv. For every  $x \in \{0, 1\}^* \setminus I$  such that

$$\Pr_{r \leftarrow \{0, 1\}^{\ell_r}} [\langle C_\sigma^g(r), P \rangle_k \text{ queries } g \text{ on input } x] \geq \frac{1}{27 \cdot k^2 \cdot q^2},$$

query  $f$  on  $x$  and, add  $x$  to  $I$  and add the tuple  $(x, f(x))$  to  $Q$ .

3. Output  $\perp$ .

First, we upper bound the number of queries that  $\mathcal{A}$  performs. Fixing an iteration and looking at step 2(a)iv, we note that fixing  $r$ ,  $\langle C_\sigma^g(r), P \rangle_k$  queries  $g$  at most  $q \cdot k$  times. Therefore, there are at most  $27 \cdot k^3 \cdot q^3$  inputs satisfying the condition in step 2(a)iv. Summing up over all iterations, we get that  $\mathcal{A}$  queries  $f$  at most  $486 \cdot k^5 \cdot q^5$  times.

Moreover, given oracle access a RE-complete oracle, the algorithm  $\mathcal{A}$  can be implemented to run in time  $\text{poly}(|t|)$ . To see this, we observe that if there exists an ensemble  $g$  that satisfies the requirements in step 2(a)i, then there exists such ensemble that can be described in space  $\text{poly}(|t|)$ . This is because given a  $g$  that is consistent with  $Q$  such that  $t^g$  is a  $\sigma$ -instance, and given a minimal certificate  $X$  for  $g$ , we can define an ensemble  $g'$  that agrees with  $g$  on  $I \cup X$  and is 0 anywhere else. The ensemble  $g'$  is consistent with  $Q$ ,  $t^{g'}$  is a  $\sigma$ -instance, and  $g'$  can be described in space  $\text{poly}(|t|, \mathbf{N}(t))$ , where we recall that by Corollary 5.9 it holds that  $\mathbf{N}(t) = \text{poly}(|t|)$ . The description of such  $g$  can be efficiently computed using an oracle that decides the following RE language:

$$\left\{ (t, Q, \sigma, i, b) \mid \begin{array}{l} \text{The } i\text{th bit of the description of } g \text{ is } b, \text{ where } g \text{ is an ensemble such that} \\ g \text{ is consistent with } Q, t^g \text{ is a } \sigma\text{-instance, and the description of } g \text{ is the} \\ \text{lexicographically first description of an ensemble satisfying these conditions} \end{array} \right\}.$$

This shows how to implement step 2(a)i. Step 2(a)iv can be implemented with respect to the lexicographically first  $P$  satisfying the condition in step 2(a)iii using the following RE language:

$$\left\{ (t, g, i, b) \mid \begin{array}{l} \text{The } i\text{th bit of } (x_1, \dots, x_s) \text{ is } b, \text{ where } x_1, \dots, x_s \text{ are the } x\text{'s satisfying} \\ \text{the condition in step 2(a)iv with respect to the tuple } t, \text{ the ensemble } g, \\ \text{and the lexicographically first } P \text{ satisfying the condition in step 2(a)iii} \end{array} \right\}.$$

Finally, we prove the correctness of the algorithm. It is clear that if  $\mathcal{A}$  terminates at step 2(a)ii then the output is correct, since for every function  $h$  consistent with  $Q$  it holds that  $t^h$  is a  $(1 - \sigma)$ -instance, and in particular it holds for  $f$ . We now assume towards contradiction that  $\mathcal{A}$  reaches step 3, and we further assume w.l.o.g. that  $t^f$  is a yes-instance. Let  $X$  be a certificate of  $t$  for  $f$  of size  $\mathbf{N}(t, f)$ . For the  $i$ th iteration of the algorithm, consider the oracle  $g_i$  chosen when  $\sigma = 0$ , and the corresponding  $P_i$  such that

$$\Pr_{r \leftarrow \{0,1\}^{\ell_r}} [\langle C_0^{g_i}(r), P_i \rangle_k = 1] \geq 2/3. \quad (5.4)$$

Let  $I_i$  be the state of  $I$  in the beginning of the  $i$ th iteration, and let  $X_i = I_{i+1} \setminus I_i$ , that is,  $X_i$  are the queries issued to  $f$  in the  $i$ th iteration. We claim that  $X_i \cap X \neq \emptyset$ , and since  $X_1, \dots, X_r$  are pairwise disjoint, it follows that  $|X| \geq r = 9 \cdot k^2 \cdot q^2 + 1$  in contradiction to Corollary 5.9. Indeed, we define

$$h_i(x) = \begin{cases} f(x) & x \in X \\ g_i(x) & \text{otherwise} \end{cases}.$$

Then  $h_i$  is a yes-instance because it agrees with  $f$  on the certificate  $X$ . Therefore, it holds that

$$\Pr_{r \leftarrow \{0,1\}^{\ell_r}} [\langle C_0^{h_i}(r), P_i \rangle_k = 1] \leq 1/3. \quad (5.5)$$

From Eq. (5.4) and Eq. (5.5) it follows that

$$\Pr_{r \leftarrow \{0,1\}^{\ell_r}} [\langle C_0^{g_i}(r), P_i \rangle_k \text{ queries } g_i \text{ on input from } \text{dis}(g_i, h_i)] \geq 1/3.$$

Note that  $\text{dis}(g_i, h_i) \subseteq X \setminus I_i$ , so there exists  $x \in X \setminus I_i$  such that

$$\Pr_{r \leftarrow \{0,1\}^{\ell_r}} [\langle C_0^{g_i}(r), P_i \rangle_k \text{ queries } g_i \text{ on } x] \geq \frac{1}{3 \cdot |X \setminus I_i|} \geq \frac{1}{27 \cdot k^2 \cdot q^2}.$$

Therefore, the algorithm adds  $x$  to  $I$  and we get that  $x \in X_i$ , so  $X_i \cap X \neq \emptyset$  as claimed.  $\blacksquare$

## 5.2 Putting it All Together

Given Lemma 5.2 and Lemma 5.3 we can now derive Theorem 5.1.

**Proof of Thm. 5.1.** Let  $((V, \bar{V}), \ell_r, k, M, T_M, \epsilon_{M,1}, \epsilon_{M,2})$  be a fully black-box construction of a worst-case hard (MIP, coMIP) protocol pair from a collision-resistant hash function  $f$ . Lemma 5.3 guarantees the existence of a polynomial-time algorithm  $\mathcal{A}$  such that for every ensemble  $f = \{f_n : \{0, 1\}^{2n} \rightarrow \{0, 1\}^{n-1}\}_{n \in \mathbb{N}}$ , the algorithm  $\mathcal{A}$  with oracle access to  $f$  and an RE-complete oracle decides the language  $L^f \subseteq \{0, 1\}^*$  defined by  $(V, \bar{V}, \ell_r, k)$  relative to  $f$ . That is, for every  $x \in \{0, 1\}^*$  it holds that  $\mathcal{A}^{f, \text{RE}}(x) = \chi_{L^\Psi}(x)$ . For every  $n \in \mathbb{N}$ , denote by  $T_{\mathcal{A}}(n)$  the polynomial running time of  $\mathcal{A}$  on inputs of length  $n$ .

Definition 3.4 then states that  $\mathcal{A}$  can be used for breaking the collision resistance of  $f$ . Specifically, for every ensemble  $f$  it holds that

$$\Pr \left[ x_1 \neq x_2 \wedge f_n(s, x_1) = f_n(s, x_2) \mid (x_1, x_2) \leftarrow M^{f, \mathcal{A}^{\text{RE}}}(s) \right] \geq \epsilon_{M,1}(T_{\mathcal{A}}(n)) \cdot \epsilon_{M,2}(n)$$

for infinitely many values of  $n \in \mathbb{N}$ , where the probability is taken over the choice of  $s \leftarrow \{0, 1\}^n$ , and over the internal randomness of  $M$ . The algorithm  $M$  may invoke  $\mathcal{A}$  on various input lengths (i.e., in general  $M$  is not restricted to invoking  $\mathcal{A}$  only on input length  $n$ ), and we denote by  $\ell(n)$  the maximal input length on which  $M$  invokes  $\mathcal{A}$  (when  $M$  itself is invoked on input  $s$  for  $s \in \{0, 1\}^n$ ). Thus, viewing  $M^{\mathcal{A}, \text{RE}}$  as a single oracle-aided algorithm that has access to  $f$ , its running time  $T_{M^{\mathcal{A}, \text{RE}}}(n)$  satisfies  $T_{M^{\mathcal{A}, \text{RE}}}(n) \leq T_M(n) \cdot T_{\mathcal{A}}(\ell(n))$  (this follows since  $M$  may invoke  $\mathcal{A}$  at most  $T_M(n)$  times, and the running time of  $\mathcal{A}$  on each such invocation is at most  $T_{\mathcal{A}}(\ell(n))$ ). In particular, viewing  $M' = M^{\mathcal{A}, \text{RE}}$  as a single oracle-aided algorithm that has oracle access to  $f$ , implies that  $M'$  is a  $q$ -query algorithm where  $q(n) = T_{M^{\mathcal{A}, \text{RE}}}(n)$ . This holds for any ensemble  $f$ , and given that  $q(n)$  is polynomial in  $n$  then Lemma 5.2 guarantees that  $\epsilon_{M,1}(T_{\mathcal{A}}(n)) \cdot \epsilon_{M,2}(n) \leq O(2^{-n/2})$ .

We conclude the proof noting that the algorithm  $\mathcal{A}$  provided by Lemma 5.3 runs in polynomial time. That is,  $T_{\mathcal{A}}(n) = n^c$  for a constant  $c$ , and thus we obtain  $\epsilon_{M,1}(n^c) \cdot \epsilon_{M,2}(n) \leq 2^{-\Omega(n)}$ . ■

## References

- [AH91] W. Aiello and J. Håstad. Statistical zero-knowledge languages can be recognized in two rounds. *Journal of Computer and System Sciences*, 42(3):327–345, 1991.
- [AS15] G. Asharov and G. Segev. Limits on the power of indistinguishability obfuscation and functional encryption. In *Proceedings of the 56th Annual IEEE Symposium on Foundations of Computer Science*, pages 191–209, 2015.
- [AS16] G. Asharov and G. Segev. On constructing one-way permutations from indistinguishability obfuscation. In *Proceedings of the 13th Theory of Cryptography Conference*, pages 512–541, 2016.
- [Bar13] B. Barak. Structure vs. combinatorics in computational complexity. Windows on Theory (available at <https://windowsontheory.org/2013/10/07/structure-vs-combinatorics-in-computational-complexity/>), 2013.
- [BCE<sup>+</sup>95] P. Beame, S. A. Cook, J. Edmonds, R. Impagliazzo, and T. Pitassi. The relative complexity of NP search problems. In *Proceedings of the 27th Annual ACM Symposium on Theory of Computing*, pages 303–314, 1995.

- [BD19] N. Bitansky and A. Degwekar. On the complexity of collision resistant hash functions: New and old black-box separations. In *Proceedings of the 17th Theory of Cryptography Conference*, pages 422–450, 2019.
- [BDV17] N. Bitansky, A. Degwekar, and V. Vaikuntanathan. Structure vs. hardness through the obfuscation lens. In *Advances in Cryptology – CRYPTO ’17*, pages 696–723, 2017.
- [BFL91] L. Babai, L. Fortnow, and C. Lund. Non-deterministic exponential time has two-prover interactive protocols. *Computational Complexity*, 1:3–40, 1991.
- [BG11] Z. Brakerski and O. Goldreich. From absolute distinguishability to positive distinguishability. In *Studies in Complexity and Cryptography*, pages 141–155, 2011.
- [BGI<sup>+</sup>12] B. Barak, O. Goldreich, R. Impagliazzo, S. Rudich, A. Sahai, S. P. Vadhan, and K. Yang. On the (im)possibility of obfuscating programs. *Journal of the ACM*, 59(2):6, 2012.
- [BGK<sup>+</sup>88] M. Ben-Or, S. Goldwasser, J. Kilian, and A. Wigderson. Multi-prover interactive proofs: How to remove intractability assumptions. In *Proceedings of the 20th Annual ACM Symposium on Theory of Computing*, pages 113–131, 1988.
- [BI87] M. Blum and R. Impagliazzo. Generic oracles and oracle classes. In *Proceedings of the 28th Annual IEEE Symposium on Foundations of Computer Science*, pages 118–126, 1987.
- [BL13] A. Bogdanov and C. H. Lee. Limits of provable security for homomorphic encryption. In *Advances in Cryptology – CRYPTO ’13*, pages 111–128, 2013.
- [BPR15] N. Bitansky, O. Paneth, and A. Rosen. On the cryptographic hardness of finding a Nash equilibrium. In *Proceedings of the 56th Annual IEEE Symposium on Foundations of Computer Science*, pages 1480–1498, 2015.
- [BPW16] N. Bitansky, O. Paneth, and D. Wichs. Perfect structure on the edge of chaos – trapdoor permutations from indistinguishability obfuscation. In *Proceedings of the 13th Theory of Cryptography Conference*, pages 474–502, 2016.
- [Bra79] G. Brassard. Relativized cryptography. In *Proceedings of the 20th Annual IEEE Symposium on Foundations of Computer Science*, pages 383–391, 1979.
- [CCG<sup>+</sup>94] R. Chang, B. Chor, O. Goldreich, J. Hartmanis, J. Håstad, D. Ranjan, and P. Rohatgi. The random oracle hypothesis is false. *Journal of Computer and System Sciences*, 49(1):24–39, 1994.
- [CGH04] R. Canetti, O. Goldreich, and S. Halevi. The random oracle methodology, revisited. *J. ACM*, 51(4):557–594, 2004.
- [CIY97] S. A. Cook, R. Impagliazzo, and T. Yamakami. A tight relationship between generic oracles and type-2 complexity theory. *Information and Computation*, 137(2):159–170, 1997.
- [DH76] W. Diffie and M. E. Hellman. New directions in cryptography. *IEEE Transactions on Information Theory*, 22(6):644–654, 1976.
- [For89] L. J. Fortnow. Complexity-theoretic aspects of interactive proof systems. Ph.D. Thesis, MIT, 1989.

- [GGH<sup>+</sup>13] S. Garg, C. Gentry, S. Halevi, M. Raykova, A. Sahai, and B. Waters. Candidate indistinguishability obfuscation and functional encryption for all circuits. In *Proceedings of the 54th Annual IEEE Symposium on Foundations of Computer Science*, pages 40–49, 2013.
- [GM84] S. Goldwasser and S. Micali. Probabilistic encryption. *Journal of Computer and System Sciences*, 28(2):270–299, 1984.
- [GMR89] S. Goldwasser, S. Micali, and C. Rackoff. The knowledge complexity of interactive proof systems. *SIAM Journal on Computing*, 18(1):186–208, 1989.
- [Gol00] O. Goldreich. On security preserving reductions – revised terminology. Cryptology ePrint Archive, Report 2000/001, 2000.
- [GPS16] S. Garg, O. Pandey, and A. Srinivasan. Revisiting the cryptographic hardness of finding a Nash equilibrium. In *Advances in Cryptology – CRYPTO ’16*, pages 579–604, 2016.
- [HH87] J. Hartmanis and L. A. Hemachandra. One-way functions, robustness, and the non-isomorphism of NP-complete sets. In *Proceedings of the 2nd Annual Conference on Structure in Complexity Theory*, 1987.
- [IN88] R. Impagliazzo and M. Naor. Decision trees and downward closures. In *Proceedings of the 3rd Annual Structure in Complexity Theory Conference*, pages 29–38, 1988.
- [IR89] R. Impagliazzo and S. Rudich. Limits on the provable consequences of one-way permutations. In *Proceedings of the 21st Annual ACM Symposium on Theory of Computing*, pages 44–61, 1989.
- [KMN<sup>+</sup>14] I. Komargodski, T. Moran, M. Naor, R. Pass, A. Rosen, and E. Yogev. One-way functions and (im)perfect obfuscation. In *Proceedings of the 55th Annual IEEE Symposium on Foundations of Computer Science*, pages 374–383, 2014.
- [LFK<sup>+</sup>92] C. Lund, L. Fortnow, H. J. Karloff, and N. Nisan. Algebraic methods for interactive proof systems. *Journal of the ACM*, 39(4):859–868, 1992.
- [Lub96] M. Luby. Pseudorandomness and Cryptographic Applications. Princeton University Press, 1996.
- [LV16] T. Liu and V. Vaikuntanathan. On basing private information retrieval on NP-hardness. In *Proceedings on the 13th Theory of Cryptography Conference*, pages 372–386, 2016.
- [Nis89] N. Nisan. CREW prams and decision trees. In *Proceedings of the 21st Annual ACM Symposium on Theory of Computing, May 14-17, 1989, Seattle, Washington, USA*, pages 327–335. ACM, 1989.
- [Nis91] N. Nisan. CREW prams and decision trees. *SIAM J. Comput.*, 20(6):999–1007, 1991.
- [RSA78] R. L. Rivest, A. Shamir, and L. M. Adleman. A method for obtaining digital signatures and public-key cryptosystems. *Communication of the ACM*, 21(2):120–126, 1978.
- [RTV04] O. Reingold, L. Trevisan, and S. P. Vadhan. Notions of reducibility between cryptographic primitives. In *Proceedings of the 1st Theory of Cryptography Conference, TCC 2004*, pages 1–20, 2004.

- [Rud88] S. Rudich. Limits on the Provable Consequences of One-way Functions. PhD thesis, EECS Department, University of California, Berkeley, 1988.
- [Sha90] A. Shamir.  $IP=PSPACE$ . In *Proceedings of the 31st Annual IEEE Symposium on Foundations of Computer Science*, pages 11–15, 1990.
- [SW14] A. Sahai and B. Waters. How to use indistinguishability obfuscation: Deniable encryption, and more. In *Proceedings of the 46th Annual ACM Symposium on Theory of Computing*, pages 475–484, 2014.
- [Wat15] B. Waters. A punctured programming approach to adaptively secure functional encryption. In *Advances in Cryptology – CRYPTO ’15*, pages 678–697, 2015.