

Separating Succinct Non-Interactive Arguments From All Falsifiable Assumptions

Craig Gentry*

Daniel Wichs†

June 2, 2011

Abstract

An *argument system* (computationally sound proof) for \mathcal{NP} is *succinct*, if its communication complexity is polylogarithmic the instance and witness sizes. The seminal works of Kilian '92 and Micali '94 show that such arguments can be constructed under standard cryptographic hardness assumptions with four rounds of interaction, and that they be made non-interactive in the random-oracle model. The latter construction also gives us some evidence that *succinct non-interactive arguments* (SNARGs) may exist in the standard model with a common reference string (CRS), by replacing the oracle with a sufficiently complicated hash function whose description goes in the CRS. However, we currently do not know of any construction of SNARGs with a proof of security under any simple cryptographic assumption.

In this work, we give a broad black-box separation result, showing that black-box reductions cannot be used to prove the security of *any* SNARG construction based on *any falsifiable cryptographic assumption*. This includes essentially all common assumptions used in cryptography (one-way functions, trapdoor permutations, DDH, RSA, LWE etc.). More generally, we say that an assumption is falsifiable if it can be modeled as an interactive game between an adversary and an efficient challenger that can efficiently decide if the adversary won the game. This is similar, in spirit, to the notion of falsifiability of Naor '03, and captures the fact that we can efficiently check if an adversarial strategy breaks the assumption.

Our separation result also extends to *designated verifier SNARGs*, where the verifier needs a trapdoor associated with the CRS to verify arguments, and *slightly succinct* SNARGs, whose size is only required to be sublinear in the statement and witness size.

*IBM T.J.Watson Research Center. cbgentry@us.ibm.com

†New York University. wichs@cs.nyu.edu. Work completed while visiting IBM T.J. Watson.

1 Introduction

The notion of a “proof” plays a fundamental role in theoretical computer science. For example, the class \mathcal{NP} can be defined as the class of languages that have efficiently verifiable proofs (or witnesses) of membership. Various aspects of such proofs have brought about many breakthroughs in the past decades, including the concepts of *zero-knowledge proofs* [GMR85], *interactive proofs* [GMR85, LFKN90, Sha90], *probabilistically checkable proofs (PCPs)* [BFLS91, FGL⁺91, ALM⁺98] and *computationally sound proofs* [Kil92, Mic94].

Succinct Arguments. One natural question that we can ask is how large do proofs of membership for \mathcal{NP} statements need to be? If we restrict ourselves to unconditionally convincing proofs, then it is easy to see that the proof length cannot be too small since we can decide membership by trying to all possible proofs and seeing if one verifies. In particular, \mathcal{NP} statements of size n cannot have sublinear proofs of size $o(n)$, unless $\mathcal{NP} \subseteq \mathbf{DTIME}(2^{o(n)})$, which we do not believe to be likely. Similar intuition can even be extended to showing that *interactive* proofs are unlikely to have sublinear communication complexity [GH98, GVW02, Wee05]. However, the seminal works of Kilian [Kil92] and Micali [Mic94] show that succinct proofs may be possible if we consider a relaxed notion of *computational soundness*, where only proofs generated by computationally bounded parties are convincing, and unbounded parties may be able to “prove” false statements. Such computationally sound proofs are also called *arguments*.

More concretely, we say that an argument system for \mathcal{NP} is *succinct* if, given a statement x with a witness w , the communication-complexity of the argument is bounded by $\text{poly}(n)\text{polylog}(|x| + |w|)$, where n is a *security parameter* and soundness holds for all $\text{poly}(n)$ -time bounded provers. Alternatively, one can view the above as saying that the size of the argument is bounded by some fixed polynomial in the security parameter n , and essentially independent of the exact (polynomial) size of the statement being proved and its witness. We will also consider *slightly succinct* arguments whose communication-complexity can be sublinear in the instance and witness size, and is bounded by $\text{poly}(n)o(|x| + |w|)$.

Succinct Non-Interactive Arguments. The work of [Kil92] showed how to construct succinct *interactive* arguments for \mathcal{NP} , requiring four rounds of interaction between the prover and the verifier, under the simple assumption that *collision resistant hash functions* (CRHFs) exist. The work of [Mic94] showed that such arguments can also be made fully non-interactive in the random-oracle model.¹ However, this leaves the question whether *succinct non-interactive arguments (SNARGs)* may exist in the standard model.

As pointed out in [BP04, Wee05], SNARGs that don’t require any initialization are unlikely to exist since we can always come up with a small adversarial prover that contains a single hard-coded false statement x along with some convincing proof π for it.² This is reminiscent of the reason why un-keyed CRHFs cannot exist, since there is always a small adversary with a hard-coded collision. Of course, the standard default notion of *keyed* CRHFs assumes that they are initialized with a short random public value (key), and we have many natural candidates for such constructions. Analogously, we will by default consider SNARGs that are initialized with a public value called a *common reference string (CRS)*. The central question considered in this work is whether such SNARGs exist and, if so, under what assumptions.

The construction of succinct argument in the random-oracle model from [Mic94] gives us some evidence that SNARGs may indeed exist. In particular, by replacing the random oracle with a sufficiently complicated hash function (whose description is placed in the CRS), we get a candidate SNARG construction in the standard model, whose security *seems plausible*. However, we currently do not have *any* construction of a SNARG with a formal proof of security under *any* simple cryptographic assumption.

¹In addition, [Mic94] also showed that computationally sound proofs have many other benefits such as very efficient verification procedures and the ability to prove all of $\mathcal{EXPTIME}$. Here we just focus on the succinctness and on the class \mathcal{NP} .

²At least if $\mathcal{NP} \not\subseteq \mathbf{BPTIME}(2^{o(n)})$ so that such false statements/proofs exist. The adversary is non-uniform.

1.1 Our Results

In this work we provide some explanation for the current lack of provably secure SNARG constructions. In particular, we give a broad black-box separation result showing that no construction of SNARGs can be proven secure via a *black-box reduction* from any *falsifiable cryptographic assumption* (unless that assumption is already false). The terms “black-box reduction” and “falsifiable assumption” are explained below, but on an informal level, this captures the types of assumptions and proof techniques that are used to prove the security of virtually every other primitive in cryptography.

Black-Box Reductions. A *black-box reduction* gets “oracle access” to an arbitrary successful attacker and must use it to break some underlying assumption. In other words, it only relies on the input/output behavior of the attacker, and does not use the description of the attacker otherwise. It should work for all attackers, including inefficient ones. Essentially all known proofs of cryptographic security in the literature are of this type, with one notable exception being Barak’s use of non-black-box techniques to construct zero-knowledge simulators and extractors in [Bar01]. However, it is currently unknown how to apply such techniques more broadly to other cryptographic tasks. We also note that, in contrast to most prior black-box separation results, we do not put any restrictions on the *construction* of our primitive, but only restrict its *proof of security* to using the adversary as a black-box.

Falsifiable Assumptions. A *falsifiable assumption* can be modeled as an interactive game between an *efficient* challenger and an adversary, at the conclusion of which the challenger can *efficiently* decide whether the adversary “won” the game. The assumption states that every efficient adversary has at most a negligible winning probability. Most standard cryptographic assumptions are falsifiable, including general notions (e.g. One Way Functions, Trapdoor Permutations, Oblivious Transfer, Identity Based Encryption, Fully Homomorphic Encryption etc.) and concrete assumptions (e.g. hardness of factoring, discrete logarithms, shortest vector problem, RSA, CDH, DDH, LWE etc.). The above notion of *falsifiability* is based, in spirit, on similar notions of [Nao03] and captures the fact that the challenger gives us an efficient process to test whether an adversarial strategy *falsifies* (i.e. breaks) the assumption. Intuitively, assumptions that are not falsifiable are harder to reason about, and therefore we have significantly less confidence in them.

Of course, it would be easy to prove the security of a SNARG construction, if we were to just make the *assumption* that the construction is secure. Indeed, this assumption turns out to *not* be falsifiable. In order for a challenger to decide whether an adversary produces a proof of a false statement, the challenger needs to decide whether an \mathcal{NP} statement is true or false, which may not be efficient. Our result therefore gives the first black-box separation between a meaningful non-falsifiable assumption and the class of *all* falsifiable ones.

Extensions and Limitations. Our separation result also extends to the weaker notion of *designated-verifier* SNARGs. In such arguments, the verifier keeps a private verification key associated with the CRS and uses it to verify the arguments of the system. Security is only guaranteed if the private verification key is kept hidden from the adversary. It also extends to SNARGs that are only slightly succinct.

Syntactically, one may equate designated verifier SNARGs with interactive two-round (challenge-response) arguments by thinking of the CRS as the verifier’s challenge. However, the security properties we normally require from these two primitives are different. The natural security notion for SNARGs considers *adaptive soundness*, where the adversary can choose which false statement to prove *adaptively* depending on the CRS. For two-round interactive arguments, one traditionally only considers the weaker notion of *static soundness*, where the adversary chooses the false statement prior to seeing the verifier’s challenge. Therefore, succinct two-round arguments are a weaker primitive than designated verifier SNARGs. We do not know how to extend our separation to the former primitive, and it remains as an interesting open problem to do so (or to give a construction).

We also mention that our separations are not unconditional, since without assumptions it may be possible that all \mathcal{NP} statements x have witnesses of size $o(|x|)$ or even $\text{polylog}(|x|)$. This may even be possible if (say) one-way functions exist. Therefore, to get around this, we will just make the assumption that there exist some sub-exponentially hard subset-membership problems (e.g. distributions over an \mathcal{NP} language and its complement that are indistinguishable by sub-exponential attackers). In other words, our results essentially say that a black-box proof of security of some SNARG under some falsifiable assumptions (e.g. security of RSA) would imply that either (1) the falsifiable assumption is actually false (e.g. there is a poly-time attack on RSA), or (2) there are *no* sub-exponentially hard subset membership problems in \mathcal{NP} . Note that (1) and (2) may be incomparable. When extending our separation to *slightly* succinct SNARGs, we need to assume the existence of *exponentially* hard subset-membership problems.

1.2 Our Techniques

Separation via a Simulatable Adversary. Our main technique is to show that every SNARG for an \mathcal{NP} complete language L has a *simulatable adversary* $\overline{\mathcal{P}}$. This is an *inefficient* adversarial prover that, given a CRS, outputs a false statement $x \notin L$ and a verifying proof π for it. However, it also comes with an *efficient* simulator \mathcal{S} so that no efficient machine can tell whether it is interacting with $\overline{\mathcal{P}}$ or \mathcal{S} . Assuming a simulatable adversary, our black-box separation result (Theorem 5.1) follows almost immediately. In particular, a black-box reduction is an efficient oracle-access machine $\mathcal{R}^{(\cdot)}$ which, when given access to a successful adversary, breaks some falsifiable assumption. But if $\mathcal{R}^{\overline{\mathcal{P}}}$ breaks some falsifiable assumption then the *efficient* machine $\mathcal{R}^{\mathcal{S}}$ must break it as well since the *efficient* challenger of a falsifiable assumption cannot distinguish $\overline{\mathcal{P}}$ from \mathcal{S} . Therefore, we show that if there is a black-box reduction from some *falsifiable* assumption to the soundness of a SNARG, then the assumption must already be *false*.

Existence of a Simulatable Adversary. To show the existence of a simulatable adversary, we prove a basic lemma of independent interest about *indistinguishability with auxiliary information* (Lemma 3.1). Assume that two distributions, \mathcal{L} over the set L and $\overline{\mathcal{L}}$ over $\overline{L} = \{0, 1\}^* \setminus L$, are computationally indistinguishable. Then, for any short auxiliary information π that we can give about $x \leftarrow \mathcal{L}$, there exists some information $\overline{\pi}$ that we can give about $\overline{x} \leftarrow \overline{\mathcal{L}}$ so that (x, π) and $(\overline{x}, \overline{\pi})$ are also computationally indistinguishable, where the security degrades (exponentially) with the size of π . This holds even if the auxiliary information π is *not* efficiently computable from x (say, if it depends on a witness w for x), and $\overline{\pi}$ may not be efficiently computable from \overline{x} . Our proof relies on von Neumann’s min-max theorem [vN28].

Given the above, we can show the existence of a simulatable adversary $\overline{\mathcal{P}}$ and its corresponding simulator \mathcal{S} (Lemma 4.1). Assuming the existence of a sub-exponentially hard subset-membership problem, there is an \mathcal{NP} language L along with distributions \mathcal{L} and $\overline{\mathcal{L}}$ as above, that are computationally indistinguishable. On a high level, the simulator \mathcal{S} efficiently samples $x \leftarrow \mathcal{L}$ along with a witness w and efficiently computes an honest proof π for x . The unbounded simulatable adversary $\overline{\mathcal{P}}$ samples $\overline{x} \leftarrow \overline{\mathcal{L}}$ along with some inefficiently samplable auxiliary information $\overline{\pi}$, as defied by our lemma. Because the proofs π of a SNARG are sufficiently short, the distributions $(\overline{x}, \overline{\pi})$ produced by $\overline{\mathcal{P}}$ and (x, π) produced by \mathcal{S} are computationally indistinguishable by efficient parties. In particular, that means that $\overline{\mathcal{P}}$ produces valid proofs for false statements, and hence is a successful adversary, but it can also be simulated by the efficient simulator \mathcal{S} .

Our actual proof is somewhat more involved and also deals with the fact that a reduction may call the oracle $\overline{\mathcal{P}}$ many times and may “lie” about the value of the security parameter it gives to $\overline{\mathcal{P}}$. In the latter case, $\overline{\mathcal{P}}$ may output very short false statements \overline{x} that *can* be efficiently distinguished from true statements x . Therefore, our simulator \mathcal{S} must sometimes also output false statements in a careful manner.

1.3 Related Work

Succinct Arguments. The works of [Kil92, Mic94] introduced the concept of succinct arguments. Their constructions both rely on the PCP theorem of [ALM⁺98], and the construction of [Mic94] can be seen as

an application of the Fiat-Shamir heuristic [FS86] to convert the interactive four-round protocol of [Kil92] into a non-interactive argument in the random-oracle model. The works [GH98, GVW02, Wee05] show that only languages which are “easy” have unconditional succinct proofs (even interactive ones). Therefore, computational assumptions are necessary to construct succinct argument for \mathcal{NP} .

The work of Aiello et al. [ABOR00] suggested a concrete approach for constructing succinct two-round arguments for \mathcal{NP} using private information retrieval (PIR). However, Dwork et al. [DLN⁺04] showed that such an approach is fundamentally flawed and is unlikely to follow from PIR security alone. Therefore, the question of constructing succinct arguments with fewer than four rounds of interaction under standard assumptions remains open. Recently, the works of [Mie08, CL08, Gro10] construct (designated-verifier and publicly verifiable) SNARGs under various non-falsifiable “knowledge” assumptions.

The work of Rothblum and Vadhan [RV10] also considers succinct arguments with black-box reductions under falsifiable assumptions. It shows that such arguments (even interactive ones) can be efficiently converted into PCP systems. Therefore, the heavy machinery of PCPs is “inherent” in the constructions of such arguments, explaining why all known works rely on it. However, since PCPs exist unconditionally, this result does not help *separate* such arguments from computational assumptions. Our techniques for showing the existence for a simulatable adversary seem to significantly differ from those of [RV10].

Black-Box Separations. Black-Box separation results in cryptography go back to Impagliazzo and Rudich [IR89], who showed that key-agreement (KA) cannot be *constructed* from one-way permutations (OWP) if the construction uses the OWP in a black-box manner. Since then, we have many other results in this vein (e.g. [Sim98, GKM⁺00, GMR01, RTV04, BPR⁺08]) showing that the *construction* of one primitive cannot just use another (simpler) primitive in a black-box manner. A natural criticism of such separations is that they do not address natural cryptographic constructions that use an underlying primitive in a non-black-box way, for example by using its description to run zero-knowledge proofs. In contrast, our separation *only* places limitations on the proof of security, requiring that it uses the adversary as a black-box, but does not place any restrictions on the construction. Separations of this type also appear in [DOP05, AF07, HH09]. However, in all of these previous works, some additional technical restrictions are placed on the construction and/or reduction, beyond just requiring that the proof of security is black-box in the adversary. To the best of our knowledge, ours is the first black-box separation result that does not place any other restrictions on the reduction, other than that it treats the adversary as a black box. Concurrently to our work, Pass [Pas11] provides a similarly broad black-box separations for the Schnorr identification scheme and several related primitives.

2 Preliminaries and Definitions

All of the results in this work hold with respect to a *non-uniform* model of computation. Given a (sometimes implicit) security parameter n , we identify *efficient* algorithms with $\text{poly}(n)$ -sized *randomized circuits* or, equivalently, probabilistic $\text{poly}(n)$ -time Turing Machines with $\text{poly}(n)$ -sized advice. A function $\epsilon(n)$ is called *negligible* if $\epsilon(n) = \frac{1}{n^{\omega(1)}}$ and we write $\epsilon(n) = \text{negl}(n)$ for short. We say that two distributions X_1, X_2 are $(s(n), \epsilon(n))$ -*indistinguishable* if for every circuit \mathcal{D} of size $s(n)$, we have $|\Pr[\mathcal{D}(X_1) = 1] - \Pr[\mathcal{D}(X_2) = 1]| \leq \epsilon(n)$. We say that they are (plain) *computationally indistinguishable* if for every $s(n) = \text{poly}(n)$ there is some $\epsilon(n) = \text{negl}(n)$ such that the distributions are $(s(n), \epsilon(n))$ -indistinguishable.

2.1 Succinct Non-Interactive Arguments (SNARGs)

We first formally define the properties that we expect for a SNARG. Since our focus will be on negative results, we will give a weaker definition than what one might expect. In particular, our default notions will be a *designated verifier* SNARG that requires some secrete associated with the CRS for verification. Also, we will give a slightly relaxed notion of succinctness.

A SNARG system Π consists of three efficient machines $\Pi = (\mathcal{G}, \mathcal{P}, \mathcal{V})$. The generation algorithm $(\text{crs}, \text{priv}) \leftarrow \mathcal{G}(1^n)$ produces a common reference string crs along with some private verification state priv . The prover algorithm $\pi \leftarrow \mathcal{P}(\text{crs}, x, w)$ produces a proof π for a statement x using a witness w . The verification algorithm $\mathcal{V}(\text{priv}, x, \pi)$ decides if π is a valid proof for x , using the private state priv .

Definition 2.1. We say that $\Pi = (\mathcal{G}, \mathcal{P}, \mathcal{V})$ is a succinct non-interactive argument (SNARG) for an \mathcal{NP} language L with a corresponding \mathcal{NP} relation R , if it satisfies the following three properties:

Completeness: For all $(x, w) \in R$, $\Pr \left[\mathcal{V}(\text{priv}, x, \pi) = 0 \mid \begin{array}{l} (\text{crs}, \text{priv}) \leftarrow \mathcal{G}(1^n) \\ \pi \leftarrow \mathcal{P}(\text{crs}, x, w) \end{array} \right] = \text{negl}(n)$.

Soundness: For all efficient $\bar{\mathcal{P}}$, $\Pr \left[\begin{array}{l} \mathcal{V}(\text{priv}, x, \pi) = 1 \\ \wedge x \notin L \end{array} \mid \begin{array}{l} (\text{crs}, \text{priv}) \leftarrow \mathcal{G}(1^n) \\ (x, \pi) \leftarrow \bar{\mathcal{P}}(1^n, \text{crs}) \end{array} \right] = \text{negl}(n)$.

Succinctness: The length of a proof is given by $|\pi| = \text{poly}(n)(|x| + |w|)^{o(1)}$.

The above will be our default notion. However, we will also consider a weaker notion of the succinctness property given below, and we call arguments that satisfy only this weaker notion slightly succinct.

Slightly Succinct: The length of a proof is given by $|\pi| = \text{poly}(n)o(|x| + |w|)$.

One can view the succinctness property as saying that there is a fixed polynomial bound on the size of the proof π , no matter how large polynomial-sized statement/witness pair is used to create the proof.

Public vs. Designated Verifiability. We say that a SNARG is *publicly verifiable* if the private verification state is just $\text{priv} = \text{crs}$. In that case, proofs can be verified by all parties. Otherwise, we call it a *designated-verifier* SNARG, in which case only the party that knows priv can verify proofs and soundness only holds if priv is kept private. The latter weaker definition is our default notion and all our negative results hold *even* for the case of designated-verifier SNARGs. In the case of designated-verifier SNARGs, our definition is perhaps too weak since it does not address the issue of reusing the crs for multiple proofs. In particular, if the same crs is used multiple times, then the verifier's accept/reject decisions can leak some information about the private verification state priv to an adversary, and eventually may allow it to prove false statements. However, since our focus is on *negative results*, our results are only strengthened by considering this weaker definition of security, and ignoring the issue of reusability. This issue also goes away when considering publicly verifiable SNARGs.

2.2 Falsifiable Cryptographic Assumptions

We now define our notion of falsifiable cryptographic assumptions. Although similar in spirit to that of [Nao03] our definition is significantly more inclusive (and simpler). Recall that the goal of [Nao03] was to measure how reasonable an assumption is by looking at how easy it is to falsify, while our goal is to be as inclusive as possible. We wish to capture essentially all reasonable assumptions, but do not mind capturing some unreasonable ones as well.

Definition 2.2. A falsifiable cryptographic assumption consists of an efficient interactive challenger \mathcal{C} and a constant $c \in [0, 1)$. On security parameter n , the challenger $\mathcal{C}(1^n)$ interacts with a machine $\mathcal{A}(1^n)$ and may output a special symbol win . If this occurs, we say that $\mathcal{A}(1^n)$ wins $\mathcal{C}(1^n)$.

The assumption associated with the tuple (\mathcal{C}, c) states that for any efficient \mathcal{A} , we have $\Pr[\mathcal{A}(1^n) \text{ wins } \mathcal{C}(1^n)] \leq c + \text{negl}(n)$, where the probability is over the random coins of \mathcal{C} and \mathcal{A} .

For any constant $\delta > 0$, the δ -exponential version of the assumption associated with (\mathcal{C}, c) states that for every \mathcal{A} of size $2^{O(n^\delta)}$ we have $\Pr[\mathcal{A}(1^n) \text{ wins } \mathcal{C}(1^n)] \leq c + 1/2^{\Omega(n^\delta)}$.

We say that a (possibly inefficient) machine \mathcal{A} breaks such assumption if its probability of winning exceeds that of the assumption.

Our separation result holds with respect to both standard and δ -exponential versions of falsifiable assumptions (for any constant δ). For simplicity, we only consider standard versions in the main body and discuss how to extend our proofs to the exponential case in Appendix A.

Our definition of a falsifiable assumptions captures most cryptographic assumptions used in the literature. For example, with $c = 0$, we capture various *search* assumptions such as the one-wayness of factoring, (strong) RSA, discrete-logarithm (DL), computational Diffie-Hellman (CDH), search-SVP in lattices etc. With $c = \frac{1}{2}$, we capture various *decisional* assumptions such as the decisional Diffie-Hellman (DDH), the decisional Learning with Errors (LWE) assumptions, etc.³ The definition also captures complicated interactive assumptions; for example take an arbitrary complicated proposal for a signature scheme or an identity-based encryption (IBE) scheme and just assume that it is secure – this too is falsifiable! In particular, the security definitions of most cryptographic primitives (e.g. one-way functions, trapdoor permutations, signatures, IBE schemes, fully homomorphic encryption etc.) involve an efficient challenger and hence assuming that some arbitrary construction satisfies these definitions is a falsifiable assumption.

Several examples of cryptographic assumptions are not obviously falsifiable. For example, we can take some arbitrary proof system and assume that it is *zero knowledge*. This assumption is not obviously falsifiable since it is not clear how to state it as a game between an efficient challenger and an adversary (the definition of zero-knowledge requires the existence of a simulator for every adversary). On the other hand, we have specific constructions of zero-knowledge proofs whose security is based on a simple falsifiable assumption (existence of one-way functions) in a black-box way.

Directly related to the topic of this paper, the assumption that some construction of a SNARG is sound is also not obviously falsifiable. Even though this assumption is stated as a game between a challenger and an adversary, the challenger is not necessarily efficient since it needs to decide whether the received statement x is in some \mathcal{NP} language to decide whether the adversary wins. Unlike the case of zero-knowledge proofs, the results in this paper show that the existence of SNARGs is also not *implied* by any falsifiable assumption in a black-box way.

Another example of assumptions that are not obviously falsifiable are the various “Knowledge” assumptions (e.g. Knowledge of Exponent) [Dam91, HT98, BP04], whose definitions postulate the existence of some non-black-box extractor.

2.3 Black-Box Reductions

For concreteness, we only discuss black-box reductions showing the soundness of some SNARG system $\Pi = (\mathcal{G}, \mathcal{P}, \mathcal{V})$ based on some falsifiable assumption (\mathcal{C}, c) .

Definition 2.3. A (possibly inefficient) machine $\overline{\mathcal{P}}$ is a Π -adversary if there exists a polynomial $p(\cdot)$ and infinitely many $n \in \mathbb{N}$ s.t. $\Pr \left[\mathcal{V}(\text{priv}, x, \pi) = 1 \wedge x \notin L \mid \begin{array}{l} (\text{crs}, \text{priv}) \leftarrow \mathcal{G}(1^n) \\ (x, \pi) \leftarrow \overline{\mathcal{P}}(1^n, \text{crs}) \end{array} \right] \geq 1/p(n)$.

Definition 2.4. A black-box reduction showing the soundness of Π based on a falsifiable assumption (\mathcal{C}, c) is an efficient oracle-access machine $\mathcal{R}^{(\cdot)}$ such that, for every (possibly inefficient) Π -adversary $\overline{\mathcal{P}}$ the machine $\mathcal{R}^{\overline{\mathcal{P}}}$ breaks the assumption.

In general, if an adversary is stateful, a black-box reduction can also use *rewinding* in addition to oracle access to an adversary. However, when it comes to SNARGs, we can (without loss of generality) restrict ourselves to stateless adversaries for which rewinding is useless. In particular, we will only consider stateless adversaries $\overline{\mathcal{P}}$ and assume that the reduction $\mathcal{R}^{\overline{\mathcal{P}}}(1^n)$ can only query the adversary with arbitrarily many inputs of the form $(1^m, \text{crs})$ and learn the output of $\overline{\mathcal{P}}(1^m, \text{crs})$ with fresh random coins. Note that the reduction need not set $m = n$ and can query the adversary on inputs of arbitrary (polynomial) size.

³In fact, these assumptions can also be captured with $c = 0$, using a clever idea described in [HH09]. For simplicity, we just explicitly allow arbitrary constants c in our definition.

The order of quantifiers in the definition requires that there is a single reduction \mathcal{R} that works for all adversaries $\overline{\mathcal{P}}$. However, our impossibility result will actually construct a single adversary $\overline{\mathcal{P}}$ for which no reduction can succeed. Therefore, we can even rule out the weaker order of quantifiers.

2.4 Hard Subset Membership Problems

A subset membership problem consists of an \mathcal{NP} language L with a corresponding relation R along with:

- A distribution-ensemble $\mathcal{L} = \{\mathcal{L}_n\}_{n \in \mathbb{N}}$ over the language L and $\overline{\mathcal{L}} = \{\overline{\mathcal{L}}_n\}_{n \in \mathbb{N}}$ over $\overline{L} = \{0, 1\}^* \setminus L$. The latter need not be efficiently samplable.
- An efficient sampling algorithm $(x, w) \leftarrow \text{Sam}(1^n)$ whose support lies in the relation R and whose projection to the first coordinate yields the distributions $\mathcal{L} = \{\mathcal{L}_n\}_{n \in \mathbb{N}}$.

We do not put any additional requirements on the *size* of the statements $x \leftarrow \mathcal{L}_n$, but since they can be efficiently sampled via $(x, \cdot) \leftarrow \text{Sam}(1^n)$, their size must be polynomial in n . It will be easiest for us to think of all sizes and hardness-measures as functions of the security parameter n .

Definition 2.5. *Let $(\mathcal{L}, \overline{\mathcal{L}}, \text{Sam})$ be a subset-membership problem over the \mathcal{NP} language L . We say that the problem is hard if the distribution-ensembles $\mathcal{L}, \overline{\mathcal{L}}$ are computationally indistinguishable. It is $(s(n), \epsilon(n))$ -hard if the distributions $\mathcal{L}_n, \overline{\mathcal{L}}_n$ are $(s(n), \epsilon(n))$ -indistinguishable. It is sub-exponentially hard if there exists some constant $\delta > 0$ such that the problem is $(s(n), \epsilon(n))$ -hard with $s(n) = 2^{\Omega(n^\delta)}$, $\epsilon(n) = 1/2^{\Omega(n^\delta)}$. It is exponentially hard if the above occurs and the size of the statements/witnesses $(x, w) \leftarrow \text{Sam}(1^n)$ is $(|x| + |w|) = O(n^\delta)$.*

Note on Hardness. The existence of any (sub-)exponentially hard subset-membership problem $\text{Sam}, \mathcal{L}, \overline{\mathcal{L}}$ implies the existence of $(2^{n^d}, 2^{-n^d})$ -hard subset-membership problems for any (arbitrarily large) constant d , simply by defining $\text{Sam}'(n) = \text{Sam}(m(n))$, $\mathcal{L}'_n = \mathcal{L}_{m(n)}$, $\overline{\mathcal{L}}'_n = \overline{\mathcal{L}}_{m(n)}$ for a large enough polynomial $m(\cdot)$. Note that sub-exponential hardness does not attempt to relate hardness to the size of the instance/witness which may grow much faster than n^d . On the other hand, exponential hardness also guarantees that hardness grows linearly with the size of the statement/witness.

Plausibility. The existence of hard subset membership problems is a relatively mild assumption. For example, any pseudorandom generator (PRG) immediately gives us a subset membership problem by setting \mathcal{L} to be the output-distribution of the PRG and $\overline{\mathcal{L}}$ to be uniform over all other strings. Hard subset membership problems are also immediately implied by many decisional assumptions such as DDH, LWE or QR. It is also reasonable to assume that these candidates are sub-exponentially hard (it is hard to imagine any practical cryptographic constructions if sub-exponentially hard problems don't exist). Exponential hardness is somewhat stronger, and is known to *not* hold for some assumptions such as QR or DDH over subgroups of \mathbb{Z}_p^* . However, we have many other candidates such as DDH over elliptic-curve groups, lattice based assumptions, and unstructured pseudorandom generator constructions, where we do believe exponential hardness to hold.

3 Indistinguishability with Auxiliary Information

Assume that two distributions \mathcal{L} and $\overline{\mathcal{L}}$ are computationally indistinguishable. Then, for any (inefficient) short auxiliary information π that we can give about samples $x \leftarrow \mathcal{L}$, we show that there exists some “lie” $\overline{\pi}$ that we can give about samples $\overline{x} \leftarrow \overline{\mathcal{L}}$ so that (x, π) and $(\overline{x}, \overline{\pi})$ are also computationally indistinguishable. The security degrades (exponentially) with the size of the auxiliary information π .

Lemma 3.1. *There is some polynomial poly for which the following holds. Let $\mathcal{L}_n, \bar{\mathcal{L}}_n$ be two arbitrary distributions that are $(s(n), \epsilon(n))$ -indistinguishable. Let \mathcal{L}_n^* be some augmented distribution on tuples (x, π) , where x is distributed according to \mathcal{L}_n and π is some arbitrary correlated auxiliary information of length $|\pi| = \ell(n)$. Then there exists an augmented distribution $\bar{\mathcal{L}}_n^*$ on tuples $(\bar{x}, \bar{\pi})$ with \bar{x} distributed according to $\bar{\mathcal{L}}_n$, such that \mathcal{L}_n^* and $\bar{\mathcal{L}}_n^*$ are $(s^*(n), \epsilon^*(n))$ -indistinguishable for $s^*(n) = s(n)\text{poly}(\epsilon(n)/2^{\ell(n)})$, $\epsilon^*(n) = 2\epsilon(n)$.*

Remarks on Lemma 3.1. In the above lemma, the distributions $\mathcal{L}_n, \bar{\mathcal{L}}_n, \mathcal{L}_n^*, \bar{\mathcal{L}}_n^*$ need not be efficiently samplable. In fact, the lemma is trivial if the distribution \mathcal{L}_n^* allows us to sample π efficiently given x . In that case, it's clear that $\bar{\mathcal{L}}_n^*$ should just sample $\bar{x} \leftarrow \bar{\mathcal{L}}_n$ and sample π honestly using \bar{x} . When π is not efficiently samplable from x (for example, π may depend on a witness w to the fact that x is sampled from \mathcal{L}_n) it's unclear how to define the distribution $\bar{\mathcal{L}}_n^*$. In fact, our proof does not give a simple description of the distribution $\bar{\mathcal{L}}_n^*$ either and instead shows its existence non-constructively using von Neumann's min-max theorem [vN28]. Technically, the proof bears similarity to the use of the min-max theorem in proofs of seemingly unrelated results in complexity, including Nisan's proof of the Impagliazzo hardcore lemma [Imp95] and the proof of equivalence between HILL entropy and metric entropy in [BSW03]. We believe that our lemma may also be of independent interest. For example, the lemma can be used to show that ℓ bits of auxiliary information (i.e. "leakage") on the seed of a PRG reduces the HILL entropy of its output by at most ℓ bits, since the joint distribution of the PRG output and the leakage on the seed is indistinguishable from the uniform distribution and ℓ bits of correlated information. This bears similarity to the prior results of [RTTV08, DP08].

Proof of Lemma 3.1. Define $\text{size}(m)$ be the set of all circuits of size m and let $\text{dist}(m)$ be the set of all distributions over $\text{size}(m)$. Fix the distributions $\mathcal{L}_n, \bar{\mathcal{L}}_n$, and some joint-distribution \mathcal{L}_n^* over tuples (x, π) as in the statement of the lemma. Define $\text{dist}(\bar{\mathcal{L}}_n)$ be the set of all joint-distributions on tuples $(\bar{x}, \bar{\pi})$ with the component \bar{x} distributed according to $\bar{\mathcal{L}}_n$.

Assume, by contradiction, that there does not exist any $\bar{\mathcal{L}}_n^* \in \text{dist}(\bar{\mathcal{L}}_n)$ that is $(s^*(n), \epsilon^*(n))$ -indistinguishable from \mathcal{L}_n^* . Then:

$$\begin{aligned} \epsilon^*(n) &< \min_{\bar{\mathcal{L}}_n^* \in \text{dist}(\bar{\mathcal{L}}_n)} \max_{\mathcal{D}_n \in \text{size}(s^*(n))} \left| \Pr_{(\bar{x}, \bar{\pi}) \in_R \bar{\mathcal{L}}_n^*} [\mathcal{D}_n(\bar{x}, \bar{\pi}) = 1] - \Pr_{(x, \pi) \in_R \mathcal{L}_n^*} [\mathcal{D}_n(x, \pi) = 1] \right| \\ &\leq \min_{\bar{\mathcal{L}}_n^* \in \text{dist}(\bar{\mathcal{L}}_n)} \max_{\mathcal{D}_n \in \text{size}(s^*(n)+1)} \Pr_{(\bar{x}, \bar{\pi}) \in_R \bar{\mathcal{L}}_n^*} [\mathcal{D}_n(\bar{x}, \bar{\pi}) = 1] - \Pr_{(x, \pi) \in_R \mathcal{L}_n^*} [\mathcal{D}_n(x, \pi) = 1] \end{aligned} \quad (1)$$

$$\begin{aligned} &= \min_{\bar{\mathcal{L}}_n^* \in \text{dist}(\bar{\mathcal{L}}_n)} \max_{\mathcal{D}_n \in \text{size}(s^*(n)+1)} \mathbb{E}_{(\bar{x}, \bar{\pi}) \in_R \bar{\mathcal{L}}_n^*} \left[\mathcal{D}_n(\bar{x}, \bar{\pi}) - \Pr_{(x, \pi) \in_R \mathcal{L}_n^*} [\mathcal{D}_n(x, \pi) = 1] \right] \\ &\leq \min_{\bar{\mathcal{L}}_n^* \in \text{dist}(\bar{\mathcal{L}}_n)} \max_{\mathbb{D}_n \in \text{dist}(s^*(n)+1)} \mathbb{E}_{\substack{(\bar{x}, \bar{\pi}) \in_R \bar{\mathcal{L}}_n^* \\ \mathcal{D}_n \in_R \mathbb{D}_n}} \left[\mathcal{D}_n(\bar{x}, \bar{\pi}) - \Pr_{(x, \pi) \in_R \mathcal{L}_n^*} [\mathcal{D}_n(x, \pi) = 1] \right] \end{aligned} \quad (2)$$

$$= \max_{\mathbb{D}_n \in \text{dist}(s^*(n)+1)} \min_{\bar{\mathcal{L}}_n^* \in \text{dist}(\bar{\mathcal{L}}_n)} \mathbb{E}_{\substack{(\bar{x}, \bar{\pi}) \in_R \bar{\mathcal{L}}_n^* \\ \mathcal{D}_n \in_R \mathbb{D}_n}} \left[\mathcal{D}_n(\bar{x}, \bar{\pi}) - \Pr_{(x, \pi) \in_R \mathcal{L}_n^*} [\mathcal{D}_n(x, \pi) = 1] \right] \quad (3)$$

$$= \max_{\mathbb{D}_n \in \text{dist}(s^*(n)+1)} \min_{\bar{\mathcal{L}}_n^* \in \text{dist}(\bar{\mathcal{L}}_n)} \mathbb{E}_{\substack{(\bar{x}, \bar{\pi}) \in_R \bar{\mathcal{L}}_n^* \\ \mathcal{D}_n \in_R \mathbb{D}_n}} [\mathcal{D}_n(\bar{x}, \bar{\pi})] - \mathbb{E}_{\substack{(x, \pi) \in_R \mathcal{L}_n^* \\ \mathcal{D}_n \in_R \mathbb{D}_n}} [\mathcal{D}_n(x, \pi)] \quad (4)$$

Where (1) follows by (possibly) negating the output of the distinguisher to make the distinguishing advantage positive, (2) follows since maximizing over distributions can only increase the advantage, (3) follows by von Neumann's min-max theorem [vN28], and (4) follows by linearity of expectation.

Let $\mathbb{D}_n \in \text{dist}(s^*(n)+1)$ be a distribution that maximizes equation (4), and let $\bar{\mathcal{L}}_n^*$ be a corresponding minimizing distribution. For each pair (x, π) , let

$$\text{Val}(x, \pi) \stackrel{\text{def}}{=} \mathbb{E}_{\mathcal{D}_n \in_R \mathbb{D}_n} [\mathcal{D}_n(x, \pi)] \quad , \quad \text{Val}_{\min}(x) \stackrel{\text{def}}{=} \min_{\pi} \text{Val}(x, \pi)$$

The minimality of $\bar{\mathcal{L}}_n^*$ implies that any $(\bar{x}, \bar{\pi})$ in the support of $\bar{\mathcal{L}}_n^*$, satisfies $\text{Val}(\bar{x}, \bar{\pi}) = \text{Val}_{\min}(\bar{x})$. Therefore, equation (4) gives us $\bar{\rho}(n) - \rho(n) \geq \epsilon^*(n)$ where:

$$\begin{aligned}\bar{\rho}(n) &\stackrel{\text{def}}{=} \mathbb{E}_{\bar{x} \in_R \bar{\mathcal{L}}_n} [\text{Val}_{\min}(x)] = \mathbb{E}_{(\bar{x}, \bar{\pi}) \in_R \bar{\mathcal{L}}_n^*} [\text{Val}(\bar{x}, \bar{\pi})] = \mathbb{E}_{\substack{(\bar{x}, \bar{\pi}) \in_R \bar{\mathcal{L}}_n^* \\ \mathcal{D}_n \in_R \mathbb{D}_n}} [\mathcal{D}_n(\bar{x}, \bar{\pi})] \\ \rho(n) &\stackrel{\text{def}}{=} \mathbb{E}_{x \in_R \mathcal{L}_n} [\text{Val}_{\min}(x)] \leq \mathbb{E}_{(x, \pi) \in_R \mathcal{L}_n^*} [\text{Val}(x, \pi)] = \mathbb{E}_{\substack{(x, \pi) \in_R \mathcal{L}_n^* \\ \mathcal{D}_n \in_R \mathbb{D}_n}} [\mathcal{D}_n(x, \pi)]\end{aligned}$$

This gives us a distinguisher $\tilde{\mathcal{D}}_n$ for the original subset-membership problem, which attempts to distinguish whether x is from \mathcal{L}_n or $\bar{\mathcal{L}}_n$ by estimating $\text{Val}_{\min}(x)$. We first start by describing an inefficient, randomized distinguisher which needs the ability to sample from \mathbb{D}_n .

Description of $\tilde{\mathcal{D}}_n$: On input x , try all possible strings π of size $|\pi| = \ell(n)$. For each such π , choose $q(n) = \frac{128}{\epsilon^*(n)^2} (\ell(n) + \ln(16/\epsilon^*(n)))$ different circuits $\mathcal{D}^{(1)}, \dots, \mathcal{D}^{(q(n))}$ at random from the distribution \mathbb{D}_n and run them all on (x, π) . Let $\tilde{\rho}_{x, \pi}$ be the fraction of the $q(n)$ circuits that output 1. Let $\tilde{\rho}_x$ be the minimum of the estimates $\tilde{\rho}_{x, \pi}$ over all the values π . Output 1 with probability $\tilde{\rho}_x$.

In any single iteration of the above algorithm with a fixed x, π , we can use Chernoff to bound the difference between the estimate $\tilde{\rho}_{x, \pi}$ and the true value $\text{Val}(x, \pi)$ with

$$\Pr[|\tilde{\rho}_{x, \pi} - \text{Val}(x, \pi)| \geq \epsilon^*(n)/8] \leq 2e^{-q(n)\epsilon^*(n)^2/128} \leq 2^{-\ell(n)}\epsilon^*(n)/8$$

where the probability is over the random coins of $\tilde{\mathcal{D}}_n$ used to compute $\tilde{\rho}_{x, \pi}$. Using the union bound over all the iterations, we derive that, with good probability, *all* estimates are close to their true values:

$$\Pr[\exists \pi \in \{0, 1\}^{\ell(n)} \text{ s.t. } |\tilde{\rho}_{x, \pi} - \text{Val}(x, \pi)| \geq \epsilon^*(n)/8] \leq \epsilon^*(n)/8$$

This also tells us that

$$\Pr[|\tilde{\rho}_x - \text{Val}_{\min}(x)| \geq \epsilon^*(n)/8] \leq \epsilon^*(n)/8 \quad (5)$$

Let us define the event E_x to be the event $|\tilde{\rho}_x - \text{Val}_{\min}(x)| \geq \epsilon^*(n)/8$. So, using the above, analysis we get:

$$\begin{aligned}\Pr_{x \in_R \mathcal{L}_n} [\tilde{\mathcal{D}}_n(x) = 1] &= \sum_x \Pr[\mathcal{D}_n(x) = 1] \Pr[\mathcal{L}_n = x] \leq \sum_x (\Pr[\mathcal{D}_n(x) = 1 | \neg E_x] + \epsilon^*(n)/8) \Pr[\mathcal{L}_n = x] \\ &\leq \sum_x (\text{Val}_{\min}(x) + \epsilon^*(n)/4) \Pr[\mathcal{L}_n = x] = \rho(n) + \epsilon^*(n)/4 \\ \Pr_{x \in_R \bar{\mathcal{L}}_n} [\tilde{\mathcal{D}}_n(x) = 1] &= \sum_x \Pr[\mathcal{D}_n(x) = 1] \Pr[\bar{\mathcal{L}}_n = x] \geq \sum_x (\Pr[\mathcal{D}_n(x) = 1 | \neg E_x] (1 - \epsilon^*(n)/8) \Pr[\mathcal{L}_n = x]) \\ &\geq \sum_x (\text{Val}_{\min}(x) - \epsilon^*(n)/4) \Pr[\mathcal{L}_n = x] \geq \bar{\rho}(n) - \epsilon^*(n)/4\end{aligned}$$

So the distinguishing advantage of $\tilde{\mathcal{D}}_n$ is

$$\Pr_{x \in_R \bar{\mathcal{L}}_n} [\tilde{\mathcal{D}}_n(x) = 1] - \Pr_{x \in_R \mathcal{L}_n} [\tilde{\mathcal{D}}_n(x) = 1] \geq \bar{\rho}(n) - \rho(n) - \frac{\epsilon^*(n)}{2} \geq \frac{\epsilon^*(n)}{2}$$

as we wanted to show. But, so far, we have only constructed a *randomized* distinguisher $\tilde{\mathcal{D}}_n$ that needs to sample from the arbitrarily complicated distribution \mathbb{D}_n over circuits. However, using an averaging argument, we can always fix the optimal coins of $\tilde{\mathcal{D}}_n$ (maximizing its distinguishing advantage) to get a *deterministic circuit* with distinguishing advantage at least $\frac{\epsilon^*(n)}{2}$. Once we fix the randomness of $\tilde{\mathcal{D}}_n$, its description just consists of $2^{\ell(n)}q(n)$ different samples from \mathbb{D}_n consisting of circuits $\mathcal{D}^{(i)}$ of size $s^*(n) + 1$

each. Therefore, there exists a *deterministic distinguisher* for the original subset-membership problem $\mathcal{L}_n, \overline{\mathcal{L}}_n$, with advantage $\epsilon(n) = \frac{\epsilon^*(n)}{2}$ and size

$$s(n) = s^*(n)O\left(2^{\ell(n)}q(n)\right) = s^*(n)\text{poly}\left(2^{\ell(n)}/\epsilon(n)\right)$$

where the exact polynomial poly is independent of the choices of $\mathcal{L}_n, \overline{\mathcal{L}}_n, \mathcal{L}_n^*, \overline{\mathcal{L}}_n^*$.

4 A Simulatable Adversary For Any SNARG

We now use our basic indistinguishability with auxiliary information lemma to prove the existence of a simulatable adversary for any SNARG.

Lemma 4.1. *Let L be a language with a sub-exponentially hard subset-membership problem. Let $\Pi = (\mathcal{G}, \mathcal{P}, \mathcal{V})$ be a non-interactive proof system for the language L that satisfies the completeness and succinctness properties. Then, there is a machine $\overline{\mathcal{P}}$, called a simulatable Π -adversary satisfying the following:*

- $\overline{\mathcal{P}}$ is a stateless and computationally unbounded Π -adversary. On input $(1^m, \text{crs})$ it always outputs some (x, π) with $x \notin L$ of size $|x| = \ell_{\text{st}}(m)$, for some polynomial $\ell_{\text{st}}(\cdot)$, and:

$$\Pr[\mathcal{V}(\text{priv}, x, \pi) = 1 \mid (\text{crs}, \text{priv}) \leftarrow \mathcal{G}(1^n), (x, \pi) \leftarrow \overline{\mathcal{P}}(1^n, \text{crs})] \geq 1 - \text{negl}(n).$$

- $\overline{\mathcal{P}}$ is poly-time simulatable. That is, for every efficient distinguisher \mathcal{D} there exists some efficient simulator \mathcal{S} such that: $\Pr[\mathcal{D}^{\overline{\mathcal{P}}}(1^n) = 1] - \Pr[\mathcal{D}^{\mathcal{S}(1^n)}(1^n) = 1] \leq \text{negl}(n)$. The distinguisher $\mathcal{D}^{(\cdot)}(1^n)$ can ask its oracle any query $(1^m, \text{crs})$ and need not set $m = n$. The simulator \mathcal{S} is given 1^n as input and can run in time polynomial in n on any query.

The same conclusion also holds if we assume that L has exponentially hard subset membership problems and that Π is only slightly succinct.

Proof Intuition. Given our lemma on indistinguishability with auxiliary information, the main idea of the machines $\overline{\mathcal{P}}$ and \mathcal{S} is simple: on query $(1^m, \text{crs})$, the machine \mathcal{S} efficiently samples $(x, w) \leftarrow \text{Sam}(1^m)$ and computes an honest proof π using the SNARG, and $\overline{\mathcal{P}}$ samples from the corresponding “fake” distribution $(\bar{x}, \bar{\pi}) \leftarrow \overline{\mathcal{L}}_n^*$ defined by the lemma. The main difficulty is that, if the value m used in the query is small enough compared to the actual security parameter n , then the answers from \mathcal{S} and $\overline{\mathcal{P}}$ can be distinguished. Therefore, we modify our simulator to use a *table* of hard-coded responses (given as polynomial-sized non-uniform advice) to answer all queries with a sufficiently small m .

Proof of Lemma 4.1. Fix the language L and the argument system $\Pi = (\mathcal{G}, \mathcal{P}, \mathcal{V})$. Since Π is succinct, we can pick some sufficiently large constant d such that, the length of the common reference string $\text{crs} \leftarrow \mathcal{V}(1^n)$ is bounded by $O(n^d)$ and the length of a proof $\pi \leftarrow \mathcal{P}(1^n, \text{crs}, x, w)$ is bounded by $O(n^{d+1})(|x| + |w|)^{o(1)}$. The existence of any sub-exponentially hard subset-membership problems implies that there exists a subset-membership problem $(\text{Sam}, \mathcal{L}, \overline{\mathcal{L}})$ which is $(s(n), \epsilon(n))$ -hard with $s(n) = 2^{(n^{d+2})}$ and $\epsilon(n) = 2^{-(n^{d+2})}$.⁴ Let $\ell_{\text{pf}}(n)$ be the length of $\pi \leftarrow \mathcal{P}(1^n, \text{crs}, x, w)$ when $\text{crs} \leftarrow \mathcal{V}(1^n)$, $(x, w) \leftarrow \text{Sam}(1^n)$. Then $\ell_{\text{pf}}(n) = o(n^{d+2})$. Note that, if the subset-membership problem is $(s(n), \epsilon(n))$ -hard then it is also $(s(n), \epsilon'(n))$ hard for any $\epsilon'(n) \geq \epsilon(n)$. By applying Lemma 3.1 with a carefully chosen $\epsilon'(n)$, there exist some

$$s^*(n) = s(n)\text{poly}(\epsilon'(n)2^{\ell_{\text{pf}}(n)}) = 2^{\Omega(n^{d+2})}, \epsilon^*(n) = \epsilon'(n)/2 = 2^{-\Omega(n^{d+2})}$$

for which the following holds: for any distribution \mathcal{L}_n^* that augments \mathcal{L}_n with auxiliary-information of length $\ell_{\text{pf}}(n)$, there is a distribution $\overline{\mathcal{L}}_n^*$ that augments $\overline{\mathcal{L}}_n$ such that $\mathcal{L}_n^*, \overline{\mathcal{L}}_n^*$ are $(s^*(n), \epsilon^*(n))$ -indistinguishable. We are now ready to describe the simulatable Π -adversary $\overline{\mathcal{P}}$ and the simulator \mathcal{S} .

⁴See discussion on sub-exponential hardness. The main idea is to just replace the security parameter n with a large enough polynomial in n .

The adversary $\overline{\mathcal{P}}$: For any query $(1^m, \text{crs})$, we can define the (efficiently samplable) augmentation \mathcal{L}_m^* of \mathcal{L}_m which samples $(x, w) \leftarrow \text{Sam}(1^m)$, $\pi \leftarrow \mathcal{P}(\text{crs}, x, w)$ and outputs (x, π) . By the above discussion, there exists some (possibly inefficiently samplable) augmentation $\overline{\mathcal{L}}_m^*$ of $\overline{\mathcal{L}}_m$ which is $(s^*(m), \epsilon^*(m))$ -indistinguishable from \mathcal{L}_m^* . The machine $\overline{\mathcal{P}}$ outputs a sample $(\bar{x}, \bar{\pi}) \leftarrow \overline{\mathcal{L}}_m^*$.

The simulator $\mathcal{S}(1^n)$: The simulator has a threshold $m^*(n) = \lfloor \log^{1/(d+1)} n \rfloor$. On input $(1^m, \text{crs})$ where $m > m^*(n)$ it samples $(x, w) \leftarrow \text{Sam}(1^m)$, $\pi \leftarrow \mathcal{P}(1^m, \text{crs}, x, w)$ and outputs (x, π) .

On inputs $(1^m, \text{crs})$ with $m \leq m^*(n)$, the simulator needs some short non-uniform advice \mathcal{T}_n to answer the query. In particular, we fix a polynomial bound $q(n)$ on the number of queries that the distinguisher makes to its oracle. For each $n \in \mathbb{N}$, the simulator's advice \mathcal{T}_n is a table of tuples of the form $(i, m, \text{crs}, x, \pi)$ with one such tuple for each choice of $i \in \{1, \dots, q(n)\}$, $m \in \{1, \dots, m^*(n)\}$, $\text{crs} \in \{0, 1\}^{\ell_{\text{crs}}(m)}$. If the i th oracle query made by \mathcal{D} is of the form $(1^m, \text{crs})$ where $m \leq m^*(n)$ then the simulator finds the corresponding tuple $(i, m, \text{crs}, x, \pi)$ in the table \mathcal{T}_n and outputs the corresponding values (x, π) .

The table \mathcal{T}_n is defined by considering a (stateful) machine $\overline{\mathcal{P}}(\mathcal{T}_n)$ which answers i th query $(1^m, \text{crs})$ in the same way as $\overline{\mathcal{P}}$ when $m > m^*(n)$ and using the table \mathcal{T}_n when $m \leq m^*(n)$. By an averaging argument, there exists some choice of \mathcal{T}_n such that $\Pr[\mathcal{D}^{\overline{\mathcal{P}}(\mathcal{T}_n)}(1^n) = 1] \geq \Pr[\mathcal{D}^{\overline{\mathcal{P}}}(1^n) = 1]$ and we pick this table as the advice.⁵ Its size is

$$|\mathcal{T}_n| = q(n) \sum_{m=1}^{m^*(n)} 2^{\ell_{\text{crs}}(m)} (\ell_{\text{pf}}(m) + \ell_{\text{st}}(m)) = \text{poly}(n). \quad (6)$$

Therefore the simulator $\mathcal{S}(1^n)$ is a $\text{poly}(n)$ -time machine with $\text{poly}(n)$ -sized advice.

To prove the first part of the Lemma, it is clear that $\overline{\mathcal{P}}$ only outputs tuples (x, π) with $x \notin L$. Also,

$$\begin{aligned} \Pr \left[\begin{array}{c} \mathcal{V}(\text{priv}, x, \pi) \\ = 1 \end{array} \middle| \begin{array}{c} (\text{crs}, \text{priv}) \leftarrow \mathcal{G}(1^n) \\ (x, \pi) \leftarrow \overline{\mathcal{P}}(1^n, \text{crs}) \end{array} \right] &\geq \Pr \left[\begin{array}{c} \mathcal{V}(\text{priv}, x, \pi) \\ = 1 \end{array} \middle| \begin{array}{c} (\text{crs}, \text{priv}) \leftarrow \mathcal{G}(1^n) \\ (x, w) \leftarrow \text{Sam}(1^n) \\ \pi \leftarrow \mathcal{P}(\text{crs}, x, \pi) \end{array} \right] - \epsilon^*(n) \\ &\geq 1 - \text{negl}(n) \end{aligned}$$

The first inequality follows by thinking of the verifier as a distinguisher and recalling that, for a fixed setting of crs , the distribution $\overline{\mathcal{L}}_n^*$ of the tuples (x, π) output by $\overline{\mathcal{P}}$ is $(s^*(n), \epsilon^*(n))$ -indistinguishable from the distribution \mathcal{L}_n^* defined by $(x, w) \leftarrow \text{Sam}(1^n)$, $\pi \leftarrow \mathcal{P}(\text{crs}, x, w)$. The second inequality follows by completeness.

To prove the second part of the Lemma, we first define a hybrid game where we replace $\overline{\mathcal{P}}$ with the (stateful) machine $\overline{\mathcal{P}}(\mathcal{T}_n)$ that uses the table \mathcal{T}_n to answer queries $(1^m, \text{crs})$ with $m \leq m^*(n)$. The way we constructed \mathcal{T}_n was exactly to ensure that $\Pr[\mathcal{D}^{\overline{\mathcal{P}}(\mathcal{T}_n)}(1^n)] \geq \Pr[\mathcal{D}^{\overline{\mathcal{P}}}(1^n) = 1]$. We now do a hybrid argument over the number of queries $q(n)$ that the distinguisher $\mathcal{D}(1^n)$ makes to its oracle. That is, consider (stateful) machines $\mathcal{O}_0, \dots, \mathcal{O}_{q(n)}$ where $\mathcal{O}_0 = \overline{\mathcal{P}}(\mathcal{T}_n)$, $\mathcal{O}_{q(n)} = \mathcal{S}(1^n)$, and, in general, \mathcal{O}_i answers the first i queries using the strategy of $\mathcal{S}(1^n)$ and the last $(q(n) - i)$ queries using the strategy of $\overline{\mathcal{P}}(\mathcal{T}_n)$. Then we claim that, for all $i \in \{0, \dots, q(n) - 1\}$,

$$|\Pr[\mathcal{D}^{\mathcal{O}_i}(1^n) = 1] - \Pr[\mathcal{D}^{\mathcal{O}_{i+1}}(1^n) = 1]| \leq \text{negl}(n) \quad (7)$$

To show this, fix any index i and fix any *preamble* of the experiment, consisting of the coins of \mathcal{D} as well as the first i queries/answers which were made/received by \mathcal{D} . The preamble therefore also fixes the $(i + 1)$ st query $(1^m, c)$ made by \mathcal{D} , but *not* the response. Conditioned on this preamble, if $m \leq m^*(n)$

⁵In particular, the advice to \mathcal{S} depends on the distinguisher \mathcal{D} . This is allowed since we are showing that for every \mathcal{D} there exists an \mathcal{S} .

then experiments i and $i + 1$ are identical. On the other hand, if $m \geq m^*(n)$, then the experiments only differ in the distribution of the $(i + 1)$ st response being either $\bar{\mathcal{L}}_m^*$ or \mathcal{L}_m^* . But these distributions are $(\tilde{s}(n), \tilde{\epsilon}(n))$ -indistinguishable where

$$\tilde{s}(n) = \Omega(s^*(m^*(n))) = 2^{\Omega(\log^{(d+2)/(d+1)} n)} = n^{\omega(1)} \quad , \quad \tilde{\epsilon}(n) = \Omega(\epsilon^*(m^*(n))) = n^{-\omega(1)} \quad (8)$$

Therefore, we see that equation (7) holds (even when conditioned on any fixed preamble) and, using the hybrid argument, we get $\Pr[\mathcal{D}^{\bar{\mathcal{P}}}(1^n) = 1] - \Pr[\mathcal{D}^{\mathcal{S}(1^n)}(1^n) = 1] \leq \text{negl}(n)$ as we wanted to show.

To prove the second version of the lemma, for exponential assumptions and slightly succinct arguments, we only need to slightly modify the first paragraph of the proof. In particular, we choose a constant d such that the length of the crs is bounded by $O(n^d)$ and the length of the proof is bounded by $O(n^{d+1})o(|x| + |w|)$. Then, the existence of any exponentially hard subset-membership problem implies that there also exists a subset-membership problem $(\text{Sam}, \mathcal{L}, \bar{\mathcal{L}})$ which is $(s(n), \epsilon(n))$ -hard where $s(n) = 2^{(n^{d+2})}$, $\epsilon(n) = 2^{-(n^{d+2})}$ and the statement/witness size for $(x, w) \leftarrow \text{Sam}(1^n)$ is bounded by $|x| = O(n^{d+2})$, $|w| = O(n^{d+2})$. But this means that the length of an argument π for $(x, w) \leftarrow \text{Sam}(1^n)$ is given by $\ell_{\text{pf}}(n) = o(n^{d+2})$. The rest of the proof proceeds the same way as before.

5 Black-Box Separation of SNARGs From Falsifiable Assumptions

We are now ready to state and prove the main result of the paper.

Theorem 5.1. *Assume that an \mathcal{NP} language L has a sub-exponentially hard subset-membership problem and let $\Pi = (\mathcal{G}, \mathcal{P}, \mathcal{V})$ be a candidate SNARG for L , satisfying the completeness and succinctness properties. Then, for any falsifiable assumption (\mathcal{C}, c) , one of the following must hold:*

- *The assumption (\mathcal{C}, c) is false.*
- *There is no black-box reduction showing the soundness of Π based on the assumption (\mathcal{C}, c) .*

The same conclusion holds if we assume that L has an exponentially hard subset-membership problem, and we allow Π to only be slightly succinct.

Proof. Assume there is a black-box reduction \mathcal{R} showing the soundness of Π based on the assumption (\mathcal{C}, c) . Let $\bar{\mathcal{P}}$ be the simulatable Π -adversary as defined in Lemma 4.1 (we can apply this lemma with either sub-exponential hardness and full succinctness or exponential hardness and slight succinctness). Then there exists some polynomial $q(\cdot)$ and infinitely many $n \in \mathbb{N}$ for which $\Pr[\mathcal{R}^{\bar{\mathcal{P}}}(1^n) \text{ wins } \mathcal{C}(1^n)] \geq c + 1/q(n)$. Since the challenger \mathcal{C} is efficient, we can think of \mathcal{R} and \mathcal{C} together as a single efficient oracle-access distinguisher $\mathcal{D}^{(\cdot)}$ and, using Lemma 4.1, there is therefore a poly-time simulator \mathcal{S} such that $\Pr[\mathcal{R}^{\mathcal{S}(1^n)}(1^n) \text{ wins } \mathcal{C}(1^n)] \geq c + 1/q(n) - \text{negl}(n)$. So the efficient attacker $\mathcal{R}^{\mathcal{S}(1^n)}(1^n)$ shows that the assumption (\mathcal{C}, c) is false. □

Corollary 5.2. *Assume that there exists some sub-exponentially hard subset-membership problem in \mathcal{NP} . Then any candidate SNARG construction for any \mathcal{NP} -complete language L cannot have a black-box proof of security based on a falsifiable assumption (\mathcal{C}, c) , unless the assumption is false.*

Proof. We use the fact that a SNARG for an \mathcal{NP} -complete language yields a SNARG for every \mathcal{NP} language. This is not the case if we consider *slightly succinct* arguments since Cook-Levin reductions can increase the instance sizes by polynomial factors. □

6 Conclusions

In this work, we give a broad black-box separation result showing that one cannot prove the security of SNARGs under any falsifiable assumption via a black-box reduction. The major open problem left by this work is to come up with *non-black-box* techniques that could overcome this negative result. It would also be interesting to see if our result can be extended to two-round and three-round interactive arguments with static soundness, or if one can construct such arguments and prove them secure via a black-box reduction from some falsifiable assumption. Our result can be used to show that any such reduction must necessarily use *rewinding* to get many different proofs π , under many different challenges, and for the same statement x .⁶ Lastly, it would be interesting to see if similar broad separation results using the idea of a simulatable adversary can be applied to other primitives.

References

- [ABOR00] William Aiello, Sandeep N. Bhatt, Rafail Ostrovsky, and Sivaramakrishnan Rajagopalan. Fast verification of any remote procedure call: Short witness-indistinguishable one-round proofs for np. In *ICALP*, pages 463–474, 2000.
- [AF07] Masayuki Abe and Serge Fehr. Perfect nizk with adaptive soundness. In *TCC*, pages 118–136, 2007.
- [ALM⁺98] Sanjeev Arora, Carsten Lund, Rajeev Motwani, Madhu Sudan, and Mario Szegedy. Proof verification and the hardness of approximation problems. *J. ACM*, 45(3):501–555, 1998.
- [Bar01] Boaz Barak. How to go beyond the black-box simulation barrier. In *FOCS*, pages 106–115, 2001.
- [BFLS91] László Babai, Lance Fortnow, Leonid A. Levin, and Mario Szegedy. Checking computations in polylogarithmic time. In *STOC*, pages 21–31. ACM, 1991.
- [BP04] Boaz Barak and Rafael Pass. On the possibility of one-message weak zero-knowledge. In *TCC*, pages 121–132, 2004.
- [BPR⁺08] Dan Boneh, Periklis A. Papakonstantinou, Charles Rackoff, Yevgeniy Vahlis, and Brent Waters. On the impossibility of basing identity based encryption on trapdoor permutations. In *FOCS*, pages 283–292, 2008.
- [BSW03] Boaz Barak, Ronen Shaltiel, and Avi Wigderson. Computational analogues of entropy. In *RANDOM-APPROX*, pages 200–215, 2003.
- [CL08] Giovanni Di Crescenzo and Helger Lipmaa. Succinct np proofs from an extractability assumption. In Arnold Beckmann, Costas Dimitracopoulos, and Benedikt Löwe, editors, *CiE*, volume 5028 of *Lecture Notes in Computer Science*, pages 175–185. Springer, 2008.
- [Dam91] Ivan Damgård. Towards practical public key systems secure against chosen ciphertext attacks. In *CRYPTO*, pages 445–456, 1991.
- [DLN⁺04] Cynthia Dwork, Michael Langberg, Moni Naor, Kobbi Nissim, and Omer Reingold. Succinct proofs for NP and spooky interactions. Manuscript, 2004.
- [DOP05] Yevgeniy Dodis, Roberto Oliveira, and Krzysztof Pietrzak. On the generic insecurity of the full domain hash. In *CRYPTO*, pages 449–466, 2005.

⁶For example, this gives further evidence that the natural approach of [ABOR00] for constructing such two-round arguments using PIR will not work, even if we make a stronger falsifiable assumption on the PIR, since it does not rely on rewinding.

- [DP08] Stefan Dziembowski and Krzysztof Pietrzak. Leakage-resilient cryptography. In *FOCS*, pages 293–302, 2008.
- [FGL⁺91] Uriel Feige, Shafi Goldwasser, László Lovász, Shmuel Safra, and Mario Szegedy. Approximating clique is almost np-complete (preliminary version). In *FOCS*, pages 2–12, 1991.
- [FS86] Amos Fiat and Adi Shamir. How to prove yourself: Practical solutions to identification and signature problems. In *CRYPTO*, pages 186–194, 1986.
- [GH98] Oded Goldreich and Johan Håstad. On the complexity of interactive proofs with bounded communication. *Inf. Process. Lett.*, 67(4):205–214, 1998.
- [GKM⁺00] Yael Gertner, Sampath Kannan, Tal Malkin, Omer Reingold, and Mahesh Viswanathan. The relationship between public key encryption and oblivious transfer. In *FOCS*, pages 325–335, 2000.
- [GMR85] Shafi Goldwasser, Silvio Micali, and Charles Rackoff. The knowledge complexity of interactive proof-systems (extended abstract). In *STOC*, pages 291–304, 1985.
- [GMR01] Yael Gertner, Tal Malkin, and Omer Reingold. On the impossibility of basing trapdoor functions on trapdoor predicates. In *FOCS*, pages 126–135, 2001.
- [Gro10] Jens Groth. Short pairing-based non-interactive zero-knowledge arguments. Asiacrypt, 2010. To Appear.
- [GVW02] Oded Goldreich, Salil P. Vadhan, and Avi Wigderson. On interactive proofs with a laconic prover. *Computational Complexity*, 11(1-2):1–53, 2002.
- [HH09] Iftach Haitner and Thomas Holenstein. On the (im)possibility of key dependent encryption. In Omer Reingold, editor, *TCC*, volume 5444 of *Lecture Notes in Computer Science*, pages 202–219. Springer, 2009.
- [HT98] Satoshi Hada and Toshiaki Tanaka. On the existence of 3-round zero-knowledge protocols. In Hugo Krawczyk, editor, *CRYPTO*, volume 1462 of *Lecture Notes in Computer Science*, pages 408–423. Springer, 1998.
- [Imp95] Russell Impagliazzo. Hard-core distributions for somewhat hard problems. In *FOCS*, pages 538–545, 1995.
- [IR89] Russell Impagliazzo and Steven Rudich. Limits on the provable consequences of one-way permutations. In *STOC*, pages 44–61, 1989.
- [Kil92] Joe Kilian. A note on efficient zero-knowledge proofs and arguments (extended abstract). In *STOC*, pages 723–732, 1992.
- [LFKN90] Carsten Lund, Lance Fortnow, Howard J. Karloff, and Noam Nisan. Algebraic methods for interactive proof systems. In *FOCS*, pages 2–10, 1990.
- [Mic94] Silvio Micali. Cs proofs (extended abstracts). In *FOCS*, pages 436–453, 1994.
- [Mie08] Thilo Mie. Polylogarithmic two-round argument systems. *Journal of Mathematical Cryptology*, 2(4):343–363, 2008.
- [Nao03] Moni Naor. On cryptographic assumptions and challenges. In *CRYPTO*, pages 96–109, 2003.
- [Pas11] Rafael Pass. Limits of security reductions from standard assumptions. In *STOC*, 2011.

- [RTTV08] Omer Reingold, Luca Trevisan, Madhur Tulsiani, and Salil P. Vadhan. Dense subsets of pseudorandom sets. In *FOCS*, pages 76–85, 2008.
- [RTV04] Omer Reingold, Luca Trevisan, and Salil P. Vadhan. Notions of reducibility between cryptographic primitives. In *TCC*, pages 1–20, 2004.
- [RV10] Guy N. Rothblum and Salil P. Vadhan. Are pcps inherent in efficient arguments? *Computational Complexity*, 19(2):265–304, 2010.
- [Sha90] Adi Shamir. IP=PSPACE. In *FOCS*, pages 11–15, 1990.
- [Sim98] Daniel R. Simon. Finding collisions on a one-way street: Can secure hash functions be based on general assumptions? In *EUROCRYPT*, pages 334–345, 1998.
- [vN28] J. von Neumann. Zur theorie der gesellschaftsspiele. In *Math. Annalen*, pages 100:295–320, 1928.
- [Wee05] Hoeteck Wee. On round-efficient argument systems. In *ICALP*, pages 140–152, 2005.

A Separation from Exponential Assumptions

In this section, we consider δ -exponential versions of falsifiable assumptions (Definition 2.2) for constants $\delta > 0$ and show that one cannot prove the security of a SNARG even under such assumptions via a black-box reduction. In particular, we can prove a variant of Theorem 5.1 showing that either (1) the δ -exponential version of the assumption is false, or (2) there is no black-box reduction from the δ -exponential version of the assumption to the soundness of a SNARG. Perhaps counter-intuitively, this separation result is incomparable to our original one since, although conclusion (2) is stronger than before, conclusion (1) is weaker than before. Our proof changes in only a few places.

Firstly, we need to modify our requirements on the *simulatable adversary* in Lemma 4.1. In particular, we show that for every $2^{O(n^\delta)}$ -sized distinguisher \mathcal{D} there exists some $2^{O(n^\delta)}$ -sized simulator \mathcal{S} such that $\Pr[\mathcal{D}^{\overline{\mathcal{P}}}(1^n) = 1] - \Pr[\mathcal{D}^{\mathcal{S}(1^n)}(1^n) = 1] \leq 1/2^{\omega(n^\delta)}$. To show this, we can essentially reuse the proof of Lemma 4.1, but modify the simulator’s threshold for when to give true/false statement to $m^*(n) = n^{\delta/(d+1)}$. This way, the size of the simulator \mathcal{S} given by equation (6) becomes $2^{O(n^\delta)}$. On each query, the responses of the simulator \mathcal{S} and $\overline{\mathcal{P}}$ are $(\tilde{s}(n), \tilde{e}(n))$ -indistinguishable where $\tilde{s}(n), \tilde{e}(n)$ given in equation (8) now become $\tilde{s}(n) = 2^{n^{\delta(d+2/d+1)}} = 2^{\omega(n^\delta)}$ and $\tilde{e}(n) = 2^{-\omega(n^\delta)}$. Using the hybrid argument over all $2^{O(n^\delta)}$ queries that \mathcal{D} can make, we get the modified lemma.

Now, to prove our variant of Theorem 5.1 for δ -exponential assumptions, we start with a reduction \mathcal{R} of size $2^{O(n^\delta)}$ so that $\Pr[\mathcal{R}^{\overline{\mathcal{P}}}(1^n) \text{ wins } \mathcal{C}(1^n)] \geq c + \nu(n)$ for some $\nu(n) \notin 2^{-\Omega(n^\delta)}$. But that means that $\Pr[\mathcal{R}^{\mathcal{S}}(1^n) \text{ wins } \mathcal{C}(1^n)] \geq c + \nu(n) - 1/2^{\omega(n^\delta)}$ which shows that the δ -exponential version of the assumption is false.

One interesting additional note is that, if we consider δ -exponential assumptions, we can also allow the challenger \mathcal{C} of such assumption to run in time $2^{O(n^\delta)}$.