

# Impossibilities in Succinct Arguments: Black-box Extraction and More

Matteo Campanelli<sup>1</sup>, Chaya Ganesh<sup>2</sup>, Hamidreza Khoshakhlagh<sup>3</sup>, and Janno Siim<sup>4</sup>

<sup>1</sup> Protocol Labs, [matteo@protocol.ai](mailto:matteo@protocol.ai)

<sup>2</sup> Indian Institute of Science, India, [chaya@iisc.ac.in](mailto:chaya@iisc.ac.in)

<sup>3</sup> Concordium, [hk@concordium.com](mailto:hk@concordium.com)

<sup>4</sup> Simula UiB, Bergen, Norway, [janno@simula.no](mailto:janno@simula.no)

**Abstract.** The celebrated result by Gentry and Wichs established a theoretical barrier for succinct non-interactive arguments (SNARGs), showing that for (expressive enough) hard-on-average languages, we must assume non-falsifiable assumptions. We further investigate those barriers by showing new negative and positive results related to the proof size.

1. We start by formalizing a folklore lower bound for the proof size of black-box extractable arguments based on the hardness of the language. This separates knowledge-sound SNARGs (SNARKs) in the random oracle model (that can have black-box extraction) and those in the standard model.
2. We find a positive result in the non-adaptive setting. Under the existence of non-adaptively sound SNARGs (without extractability) and from standard assumptions, it is possible to build SNARKs with black-box extractability for a non-trivial subset of **NP**.
3. On the other hand, we show that (under some mild assumptions) all **NP** languages cannot have SNARKs with black-box extractability even in the non-adaptive setting.
4. The Gentry-Wichs result does not account for the preprocessing model, under which fall several efficient constructions. We show that also, in the preprocessing model, it is impossible to construct SNARGs that rely on falsifiable assumptions in a black-box way.

Along the way, we identify a class of non-trivial languages, which we dub “trapdoor languages”, that can bypass these impossibility results.

## 1 Introduction

Proof systems have been studied extensively both in cryptography and in the theory of computation [8, 28, 34], and are a fundamental building block in various cryptographic constructions today, including delegating computation [10, 11, 23] and privacy-preserving cryptocurrencies [9] to name a few. In a *succinct* proof, it is additionally required that the communication be sublinear (ideally polylogarithmic) in the size of the non-deterministic witness used to verify the relation (*proof* succinctness). This requirement is often extended to verification complexity (*verification* succinctness).

Statistically-sound proofs are unlikely to allow for significant improvements in proof size [33, 35, 61], that is, for **NP**, statistical soundness requires the prover to communicate, roughly, as much information as the size of the witness. If we restrict ourselves to *argument systems* [14] where soundness is *computational*, then proofs can potentially be shorter than the length of the witness.

**Succinct arguments.** Succinct arguments were first studied by Kilian [46], who gave an interactive construction based on probabilistically checkable proofs (PCP) and collision-resistant hash functions. Kilian’s construction was turned into a non-interactive argument in the random oracle model using the Fiat-Shamir heuristic [27] by Micali [53]. In the standard model (i.e., without idealized primitives), non-interactivity is achieved by a trusted party generating a Common Reference String (CRS) during a setup phase. The notion of *adaptive soundness* requires soundness to hold even when a malicious prover can choose the statement after receiving CRS. Otherwise, we call soundness non-adaptive.

In this work, we are concerned with the theoretical limitations for building efficient, succinct non-interactive arguments in the standard model<sup>5</sup>. One of the best-known impossibility results on SNARGs is that of Gentry and Wichs [32] (we occasionally refer to it as “GW”), which shows that in the standard model, adaptively-sound SNARGs for (hard enough) **NP** languages cannot be proven secure via a black-box (BB) reduction to a falsifiable assumption [54]. A falsifiable assumption is an assumption where the challenger can efficiently confirm that the challenge was broken<sup>6</sup>.

A folklore way to interpret GW has been “*we cannot escape non-falsifiable assumptions to build SNARGs for NP*”. While this is essentially true, there are several caveats to this interpretation (which we discuss later in this work in section 6 and some have already been noticed in prior work). We formally explore the boundaries of this simplifying interpretation, especially motivated by the focus on (composable) extractability [6, 48] and the popular model of “preprocessing SNARGs” in recent works, e.g. [39]. We strive to provide a *modern* view of these topics, for example, by adopting the language of indexed relations from [19] (we later argue why this is a meaningful switch).

**(Black-box) knowledge soundness.** A strengthening of the soundness property is *knowledge soundness*. It requires that, whenever an efficient prover convinces the verifier, not only can we conclude that  $x \in \mathcal{L}$ , but also an **NP** witness  $w$  can be efficiently extracted from the prover. This helpful property is satisfied by many proof systems and is necessary for a lot of applications of succinct arguments. A Succinct Non-interactive ARgument of Knowledge (SNARK) is a SNARG with the knowledge soundness property.

Constructions of SNARKs for **NP** in the standard model all rely on non-falsifiable assumptions that are *knowledge-type* assumptions related to some algebraic problem, e.g., guaranteeing the existence of an extractor algorithm that can output a discrete log “from” a specific adversary. This example also hints to why these assumptions are non-falsifiable—they are *non-black-box*, that is, they require knowledge of the internal state of the adversary and an extractor aware of the concrete adversarial algorithm. This contrasts with the milder *black-box* extraction, namely the ability to extract a witness from a malicious prover only using its input/output interface.

Understanding whether we can build SNARKs with black-box extraction in the standard model is still an elusive problem. In addition to being a theoretical curiosity, if answered positively, it would allow us to construct more robust cryptographic protocols using SNARKs. Black-box extraction is required in strong notions of composition security, e.g., in universal composability (or UC-security [17]) where the “ideal-world” simulator must extract a witness without knowledge of the environment’s algorithm. (See [45] for an attempt to combine composability and knowledge-type assumptions.) If answered negatively, it would confirm the seeming incompatibility of SNARKs in the standard model and UC. In this work, we then ask the question:

*Is non-black-box extraction inherent to SNARKs?*

Addressing this question is, we believe, even more pressing because prior works [5, 6, 18, 48] have used as motivation the fact that succinctness *must* be sacrificed for black-box extraction, implying that the question had been settled (see also section 1.2). However, to the best of our knowledge, there was no formal treatment for this question prior to our work.

**OUR FIRST CONTRIBUTION:** We formally confirm the folklore belief that black-box extraction is impossible for adaptive knowledge soundness in the standard model if one requires proof-succinctness. As a consequence, this result separates the standard model and other idealized models in terms of what is possible for black-box extraction (for example, in the ROM and through the Fiat-Shamir transform, there exist black-box extractable proof-succinct non-interactive arguments [15]).

**OUR SECOND CONTRIBUTION:** We explore whether the impossibility extends to the non-adaptive case. We find out that non-adaptive black-box extractability is possible for a non-trivial subset of **NP**—which

<sup>5</sup> There exist efficient SNARKs (SNARGs of knowledge) in idealized models like ROM (random oracle model), GGM (generic group model), or AGM (algebraic group model), including constructions like Groth16 [39], Bulletproofs [15]. We later discuss the implications of our results for different models.

<sup>6</sup> For example, DLOG is a falsifiable assumption since the challenger can efficiently test if the adversary has found the correct discrete logarithm.

encompasses distributionally hard problems such as knowledge of a discrete logarithm — by assuming the existence of a non-adaptively sound SNARG and some standard assumptions (FHE). In particular, we show that a SNARG can be lifted to a SNARK with the features above for the class of languages **FewP** (roughly, **NP** statements with at most a polynomial number of valid witnesses). If the starting SNARG (non-adaptively sound) is based on falsifiable assumptions and in the standard model then so is the resulting SNARK.<sup>7</sup> Our transformation also preserves zero-knowledge of the initial SNARG.

**OUR THIRD CONTRIBUTION:** A natural question is whether the previous construction for **FewP** can be extended to **NP**. We answer this in the negative under some mild assumptions. In particular, we show that if the relation is  $y = f(w)$  where  $f$  is an  $L$ -continuous leakage-resilient one-way function (CLR-OWF, a one-way function where  $L$  bits may leak multiple times given that preimage  $w$  is updated), then the proof size must be more than  $L$  bits. There exists a CLR-OWF under the discrete logarithm assumption [1] where  $L$  is linear in the size of  $w$ . Thus, the proof cannot be succinct.

**Preprocessing and the Gentry-Wichs impossibility result.** In many applications we want to look beyond proof-succinctness and keep the verifier as efficient as possible. Ideally, we would like verification to run sublinearly in the size/time of the computation. It may seem counterintuitive that this is even possible: naturally, in circuit-based<sup>8</sup> arguments for general computations the verifier should *at least* read the statement being proven. The latter includes both the description of the computation (i.e., the circuit) and its input (i.e., the deterministic input for an **NP** statement). There exists, however, a (commonly used) way around this problem: a *preprocessing* phase. In a preprocessing SNARG, one generates a CRS, usually depending on a specific circuit  $C$ , which is constructed once and for all and can later be used to prove/verify an unbounded number of proofs for the computation of  $C$ . This CRS is structured as a *pair* of CRS'es, the prover's CRS and the verifier's CRS, used by each respective party. The verifier's CRS is morally a digest of the circuit. If the verifier's CRS is "short enough", then the online verification stage can be fast, requiring to read only the SNARG proof and a *partial* input description (the deterministic input to the circuit, without its description); thus the verifier can run in time sublinear in  $|C|$  (and in the witness size).

This preprocessing model encompasses a rich line of efficient SNARGs [13, 31, 38, 50, 51, 55]. The fact that it is a practically interesting model, as it achieves verifier-succinct SNARGs, further motivates a deeper theoretical understanding of it. A fundamental question is:

*Can we construct preprocessing SNARGs based on falsifiable assumptions?*

We argue this question has not been settled. First, none of the known preprocessing constructions rely on falsifiable assumptions. Also, known impossibility results do not inform us on the matter either. The Gentry-Wichs impossibility—which separates SNARGs and falsifiable assumptions—has long served as a justification to SNARGs for **NP** on non-falsifiable assumptions, *but it fails to shed light on the preprocessing setting*. The reason is that the GW result presumes a SNARG with a CRS with a specific pattern (we mean "prover's CRS" when we just say CRS from now on): their CRS cannot grow with the size of the instance, but should instead be bounded by a polynomial in the security parameter. In principle the question is then still open, more so because all existing preprocessing constructions, do have a CRS with the opposite pattern: it is usually as long as the instance<sup>9</sup>.

Besides GW, other existing works also fail to provide an answer. For example, the work of [12] shows how to "bootstrap" a preprocessing SNARK into one without preprocessing to obtain a *complexity-preserving* SNARK, i.e., one without expensive preprocessing. The transformation can be applied to known SNARKs with expensive preprocessing to obtain a SNARK without the costly preprocessing. This complexity-preserving

<sup>7</sup> It is still an open problem how to obtain non-adaptively secure SNARGs from polynomial-time secure falsifiable assumptions. The only candidate, the recent construction in [52], was shown to have a fundamentally flawed proof of security (see [60]).

<sup>8</sup> There are other models of computation that have a succinct description, for instance, machine computations. However, in general, the description of a computation could be as large as the computation itself.

<sup>9</sup> For example, in pairing-based constructions such as [31] it consists of at least one group element per wire in the circuit to be proven.

compilation, informally, establishes that preprocessing does not give any additional power; if preprocessing SNARKs were possible from falsifiable assumptions, one could apply the bootstrapping transformation and obtain short CRS SNARKs from falsifiable assumptions. Thus, any impossibility for SNARKs holds even for SNARKs that rely on expensive preprocessing. However, this bootstrapping crucially requires the *knowledge soundness* property and, therefore, only applies to SNARKs. The question of whether allowing a preprocessing phase allows constructing *SNARGs* based on falsifiable assumptions remains.

**OUR FOURTH CONTRIBUTION:** We fill the gap left by the GW result and show that even preprocessing SNARGs with a loosely-bounded CRS cannot be constructed from falsifiable assumptions in the standard model.

**The landscape of impossibilities for non-interactive arguments.** For our work to be as self-contained as possible, we complement the results above with an overarching view of impossibilities on non-interactive arguments (section 6). This discussion strives to give a complete picture of existing impossibility results, related key properties of positive results, and gaps between positive and negative results. Motivated by the observation that preprocessing SNARGs do not come under the GW impossibility, we articulate the assumptions behind the impossibilities, and identify settings that would bypass them. Along the way, we formalize a class of languages that does not come under the Gentry-Wichs impossibility result. We dub them *trapdoor languages* (where there exists a “trapdoor” that makes the problem feasible) and exemplify several application settings that fall under the same category. Trapdoor languages can be thought of as a generalization of witness-sampleable (algebraic) languages in the work of [24].

## 1.1 Technical Overview

**BB extraction is impossible for any hard language (adaptive case).** We show the impossibility of black-box extraction for non-interactive succinct arguments following the intuition that if an argument is too “small”, it cannot contain information about a “long” witness. This makes extraction impossible since the extractor does not have any additional power, like access to the prover’s randomness (as in non-black-box extractors for popular SNARKs) or the ability to rewind the prover (as in interactive arguments, such as Kilian’s protocol).

Our result gives a precise characterization between the hardness of guessing the witness and the size of the proof. We show that if an efficient adversary can guess the witness at most with probability  $\varepsilon(\lambda)$  and the knowledge soundness error of the argument system is  $\varepsilon_{ks}(\lambda)$ , then the proof size is at least  $-\log(\varepsilon(\lambda) + \varepsilon_{ks}(\lambda))$  bits. For example, if we consider for simplicity that  $\varepsilon_{ks}(\lambda) = 0$  and  $\varepsilon(\lambda) = 1/2^{\delta|\mathbf{w}(\lambda)|}$  for some  $\delta > 0$  and the witness size is  $|\mathbf{w}(\lambda)|$ , then the proof size will be at least  $\delta|\mathbf{w}(\lambda)|$ . In appendix B, we show how to obtain a similar result based on the hardness of leakage-resilient OWFs.

**BB extraction is possible for FewP (non-adaptive case).** We then ask if the impossibility holds if we weaken the knowledge soundness requirement to be *non-adaptive*. Indeed, the non-adaptive case escapes the GW impossibility for SNARGs as we discuss in Section 6.1, and it is natural to hope for a positive result for extraction as well. In the non-adaptive knowledge soundness definition, the adversary chooses the statement before seeing the CRS, and then outputs a proof for the chosen statement. Intuitively, an extractor for such an adversary *does have* additional power – the extractor can rewind the prover to the point after the statement is chosen, sample different CRS’es and obtain multiple proofs for the same statement. Thus, non-adaptivity makes the prover stateful allowing for rewinding to be useful for an extractor<sup>10</sup>. We give a positive result in the non-adaptive case by showing a SNARK with black-box non-adaptive extraction (for a subset of **NP**). In the construction, we take advantage of our observation that the extractor can obtain more information by seeing multiple proofs corresponding to cleverly crafted CRS’es. At a high level, we ask the prover to encrypt a bit of the witness as part of the proof, in addition to proving the underlying relation. Given the secret key of the encryption scheme as the CRS trapdoor, the extractor can recover this witness

<sup>10</sup> Contrast this with the adaptive case, where the prover is stateless and rewinding is not useful.

bit. Now, the crafted CRS'es are such that they ask for different bits of the witness to be encrypted so that with every rewinding, the extractor learns a new bit until it can completely recover the witness.

While this works for valid statements with a *unique* witness, there are some subtleties that we need to address in order to show extraction for languages that have polynomially many witnesses, that is, class **FewP**. Here, the problem is that the adversary can choose to use a different witness each time, and there is no guarantee that the extractor can collect enough bits for any one witness. We now provide an overview of our construction. Let  $\mathcal{R}$  be the relation for the language. We start with an existing SNARG for  $\mathcal{R}$  and lift it to a SNARK. We use a Fully Homomorphic Encryption (FHE) scheme in order to hide the index of the bit the prover is asked to encrypt. Intuitively this is to hide the index so that the prover cannot adversarially choose a different witness for different indices. We augment the relation the SNARG proves to include a hash of the witness. Now the extractor keeps track of which witness it is extracting by using the hash to fingerprint. The extractor still needs to collect all bits of one witness. Here, we rely on the semantic security of the FHE scheme to show that the prover cannot consistently use witness  $w_1$  for index  $i$ , and witness  $w_2$ , for index  $j$ . Since there are only polynomially many witnesses, assuming universality (non-adaptive collision resistance) of the hash function, the extractor succeeds in recovering all bits of some witness.

We also show that if the hash is encrypted and the initial SNARG has computational zero-knowledge, the resulting SNARK will also have computational zero-knowledge.

**BB extraction is impossible for all NP (non-adaptive case).** The previous result, however, cannot be extended to all NP languages. We show this by relating an extractor's existence to breaking the relation's leakage resilience. A SNARK proof can be thought of as leakage on the witness. When this leakage is small, no extractor can succeed if the NP relation is leakage resilient. This impossibility due to leakage resilience is easy to see in the adaptive case. In non-adaptive extraction, an extractor can potentially rewind the adversary and obtain multiple proofs; akin to a leakage resilience adversary obtaining leakage multiple times. We formalize this connection using *continuous* leakage resilience. In  $L$ -leakage-resilient OWF (LR-OWF), one-wayness holds even if  $L$  bits of the preimage are leaked. In  $L$ -continuous LR-OWF (CLR-OWF),  $L$  bits can be leaked multiple times with the caveat that the preimage has to be updated before each leakage. Moreover, if for an OWF  $f$  we have  $y = f(w)$  and  $w$  is updated to  $w'$ , then also  $y = f(w')$ .

We connect this primitive to the impossibility of non-adaptive black-box knowledge soundness of SNARKs. Suppose we have a SNARK for the relation  $y = f(w)$  where  $f, y$  are public, and  $w$  is the witness. We view the proof as leakage on the witness given to the adversary. If the proof is at most  $L$  bits long, then the extractor can learn at most  $L$  bits of information about the witness with each rewinding. Now if the adversary also updates its witness  $w$  between queries,  $L$ -CLR of  $f$  implies that the extractor cannot recover the witness. Thus, the SNARK proof is at least  $L$  bits long.

We can instantiate this result with  $(1 - \frac{2}{n})|w|$ -CLR-OWF from [1] which is based on the discrete logarithm assumption. The witness size  $|w| = n \log q$ , where  $q$  is the size of the discrete logarithm group and  $n$  is an input size parameter. Thus, the proof size will be asymptotically linear in  $|w|$ .

**Extending GW to preprocessing SNARKs.** The central idea in the GW proof is to show that every SNARG for an NP language has a *simulatable* adversary. An unbounded adversarial prover that breaks soundness comes with an efficient simulator such that no efficient machine can tell whether it is interacting with the prover or the simulator. A black-box reduction is an efficient oracle-access machine that breaks some falsifiable assumption when given access to a successful adversary. Suppose the reduction given oracle access to the prover breaks the assumption. In that case, the efficient machine with oracle access to the efficient simulator also breaks it since the efficient challenger of the falsifiable assumption cannot distinguish the prover from the simulator. Thus, assuming a simulatable adversary, the theorem follows.

Our proof extending the GW impossibility to preprocessing SNARKs follows the GW template. We observe that the GW proof needs the CRS to be short in constructing a simulatable adversary: the reduction that has oracle access to either the computationally unbounded prover or the efficient simulator can query the oracle with  $1^m$  where  $m$  is different from the security parameter  $n$ . If  $m$  is small enough compared to the actual security parameter  $n$ , then the reduction can distinguish the adversary from the simulator. Therefore,

the proof modifies the simulator to behave differently in answering queries with a sufficiently small  $m$ ; this is done by hardcoding a table of responses as non-uniform advice. The table has hardcoded entries  $(x, \pi)$  for every  $m$  and every CRS. Therefore, the CRS size is bounded by a polynomial in the security parameter and cannot grow with the size of the instance.

When considering security-parameter preserving reductions, the reduction queries its oracle with the same security parameter. Therefore, a hardcoded table is unnecessary, and we show how the proof goes through when the size of the CRS depends on the instance, as in indexed relations. We leave the case with non-parameter-preserving reductions as an open problem.

## 1.2 Related Work

**Succinctness vs black-box extraction.** Here we discuss works that trade succinctness for black-box extraction. *C0C0* [48] and Tiramisu [6] aim at compiling a SNARK into a UC-secure scheme. However, this transformation results in NIZK arguments whose proof size and verification time are (quasi-)linear in the witness size. This degradation in succinctness is claimed to be unavoidable if one demands black-box extraction. In [5], Bagheri et al. add black-box extraction to [39] SNARK. Although the proof size is again asymptotically linear in the witness size, the authors’ goal is to strive for concrete efficiency. In [18], Chase et al. construct controlled malleable proofs that crucially require the stronger black-box version of extractability. Even though their starting point is a SNARG, to obtain black-box extraction of the controlled malleable proof, they give up succinctness and achieve controlled malleable NIZKs.

What is common in all the works above as an idea is to perform verifiable encryption by encrypting the witness and then proving knowledge of the value inside the ciphertext in addition to the original relation. The black-box extractor works by decrypting. This is why the black-box extractor comes at the cost of succinctness: the proof includes a ciphertext and a proof of correct encryption.

**Other works.** The work in [45] proposes an alternative composability model to the UC model, which can (at least to some extent) use non-black-box extractability and knowledge-type assumptions. In this case, one can still obtain succinct UC SNARKs (under some restrictions) without needing black-box extraction. The recent work in [30] obtains witness-succinct non-interactive arguments of knowledge in UC but applying the random oracle model.

## 2 Preliminaries

PPT stands for probabilistic polynomial time. We use  $\lambda$  to denote the security parameter. We write  $x \leftarrow \$ X$  to denote that  $x$  is sampled from a distribution  $X$ . If  $X$  is a set, then  $x \leftarrow \$ X$  denotes uniform sampling. We write  $f(\lambda) = \text{negl}(\lambda)$  when  $f$  is negligible in  $\lambda$  and  $f(\lambda) = \text{poly}(\lambda)$  when  $f$  is polynomial in  $\lambda$ . For an integer  $N \geq 1$ , we define  $[N] := \{1, \dots, N\}$ .

**Indistinguishability.** We say that two distributions  $X_1$  and  $X_2$  are  $(s(\lambda), \epsilon(\lambda))$ -indistinguishable if for any circuit  $\mathcal{D}$  of size  $s(\lambda)$ , we have  $|\Pr[\mathcal{D}(X_1) = 1] - \Pr[\mathcal{D}(X_2) = 1]| \leq \epsilon(\lambda)$ .

**Hard-on-average problems.** We define a language  $\mathcal{L} \in \mathbf{NP}$  to be a hard-on-average problem if

- It has an efficient instance sampler  $\text{Samp}_{\mathcal{L}}(1^\lambda)$  that outputs  $x \in \mathcal{L}$  together with an  $\mathbf{NP}$  witness  $w$ .
- There is an efficient sampler  $\text{Samp}_{\bar{\mathcal{L}}}(1^\lambda)$  that with an overwhelming probability outputs  $x \notin \mathcal{L}$ .
- It is computationally hard to distinguish outputs of  $\text{Samp}_{\mathcal{L}}(1^\lambda)$  and  $\text{Samp}_{\bar{\mathcal{L}}}(1^\lambda)$ .

Language  $\mathcal{L}$  is  $(s(\lambda), \epsilon(\lambda))$ -hard if distributions of  $x$  from  $\text{Samp}_{\mathcal{L}}(1^\lambda)$  and  $\text{Samp}_{\bar{\mathcal{L}}}(1^\lambda)$  are  $(s(\lambda), \epsilon(\lambda))$ -indistinguishable. It is sub-exponentially hard if there exists some constant  $\delta > 0$  such that previous distributions are  $(s(\lambda), \epsilon(\lambda))$ -indistinguishable for  $s(\lambda) = 2^{\Omega(\lambda^\delta)}$  and  $\epsilon(\lambda) = 1/2^{\Omega(\lambda^\delta)}$ . Lastly,  $\mathcal{L}$  is exponentially hard if the above holds and moreover  $|x| + |w| = O(\lambda^\delta)$  for  $(x, w) \leftarrow \$ \text{Samp}_{\mathcal{L}}(1^\lambda)$ .

Simple example is the DDH language where  $\text{Samp}_{\mathcal{L}}$  outputs group elements  $g^a, g^b, g^{ab}$ , where  $a, b$  are chosen uniformly at random and  $g$  is a group generator, and  $\text{Samp}_{\bar{\mathcal{L}}}$  outputs 3 random group elements  $g^a, g^b, g^c$ . More generally, hard-on-average problem is implied by the existence of one-way-functions since it is possible to construct a PRG from a one-way function [41].

**Falsifiable assumptions.** Below we recall the notion of falsifiable assumptions.

**Definition 1 ([32]).** A falsifiable cryptographic assumption  $(\mathcal{C}, c)$  consists of a PPT challenger  $\mathcal{C}$  and a constant  $c \in [0, 1)$ . We say that  $\mathcal{A}$  wins  $(\mathcal{C}, c)$  if  $\mathcal{A}(1^\lambda)$  and  $\mathcal{C}(1^\lambda)$  interact and finally  $\mathcal{C}$  outputs 1. The assumption  $(\mathcal{C}, c)$  holds if for all non-uniform PPT  $\mathcal{A}$ ,  $\Pr[\mathcal{A} \text{ wins } (\mathcal{C}, c)] \leq c + \text{negl}(\lambda)$ . Otherwise we say that  $(\mathcal{C}, c)$  is false.

Definition 1 captures most cryptographic assumptions from the literature. In the case of *search* assumptions (e.g., discrete logarithm problem and shortest vector problem), we set  $c = 0$ . In the case of *decisional* assumptions (e.g., decisional Diffie-Hellman, decisional Learning with Errors), we set  $c = 1/2$  since the adversary can win with probability  $1/2$  by random guessing. Knowledge assumptions [25, 40] are seemingly non-falsifiable.

## 2.1 Continuous Leakage-Resilient OWFs

A leakage-resilient OWF (LR-OWF)  $f$  is a function that is one-way even when the adversary is allowed to learn arbitrary functions of  $f(x)$ 's preimage as long as this leakage is restricted to  $L$  bits. Continuous LR-OWF (CLR-OWF) in the floppy model [1, 3] is a generalization of this where leakages can happen multiple times. In short, it assumes a master secret key which is kept in a leakage-free server (e.g., on a floppy disk) and then can be used to securely update the preimage  $x$ .  $L$  bits of leakage on the preimage can occur after each update. Importantly however, updates have to preserve the output of the OWF, that is  $f(x) = f(x')$  when  $x'$  is an update of  $x$ .

More formally, a CLR-OWF consists of the following probabilistic polynomial time (PPT) algorithms: (1)  $\text{KGen}(1^\lambda)$  that outputs a public parameter  $\text{pp}$  and an update key  $\text{uk}$ . (2)  $\text{Sample}(\text{pp})$  takes as input the parameter  $\text{pp}$  and outputs a random OWF input  $x$ . (3)  $\text{Eval}(\text{pp}, x)$  is a deterministic algorithm that produces the OWF output  $y$ . (4)  $\text{Update}(\text{uk}, x)$  takes in the update key  $\text{uk}$  and  $x$ , and outputs an updated OWF input  $x'$ .

We assume that a CLR-OWF satisfies the following properties.

**Correctness.** For any  $(\text{pp}, \text{uk}) \in \text{KGen}(1^\lambda)$  and  $x \in \{0, 1\}^*$ , we have that  $\text{Eval}(\text{pp}, \text{Update}(\text{uk}, x)) = \text{Eval}(\text{pp}, x)$ .

**$L$ -Continuous leakage-resilience.** Let  $L = L(\lambda)$ . For any PPT  $\mathcal{A}$ ,

$$\Pr \left[ \begin{array}{l} (\text{pp}, \text{uk}) \leftarrow \text{KGen}(1^\lambda), x \leftarrow \text{Sample}(\text{pp}), \\ y \leftarrow \text{Eval}(\text{pp}, x), x' \leftarrow \mathcal{A}^{\text{O}_L(\cdot)}(\text{pp}, y) \end{array} : y = \text{Eval}(\text{pp}, x') \right] = \text{negl}(\lambda),$$

where  $\text{O}_L(\cdot)$  is an oracle that takes as an input a leakage function  $h : \{0, 1\}^* \rightarrow \{0, 1\}^L$ , on which  $\text{O}_L(h)$  sets  $x \leftarrow \text{Update}(\text{uk}, x)$  and then returns  $h(x)$ .

There exists CLR-OWFs [1], which can leak almost the full key. We recall this result in appendix A.

## 2.2 Argument System

We recall the notion of non-interactive argument systems.

**Definition 2 (Indexed relation [19]).** An indexed relation  $\mathcal{R}$  is a set of triples  $(i, x, w)$  where  $i$  is the index,  $x$  is the instance, and  $w$  is the NP-witness; the corresponding indexed language  $\mathcal{L}(\mathcal{R})$  is the set of pairs  $(i, x)$  for which there exists a witness  $w$  such that  $(i, x, w) \in \mathcal{R}$ . Indexed relation is associated with an efficient index sampling algorithm  $\mathcal{I}$  that outputs an index  $i$  on input  $1^\lambda$ .

For example,  $i$  can be an arithmetic circuit and  $x$  and  $w$  public and private inputs to the circuit such that the circuit outputs 1. We say that an indexed language is a hard-on-average problem if it is defined like in section 2, but additionally  $\text{Samp}_{\mathcal{L}}$  and  $\text{Samp}_{\bar{\mathcal{L}}}$  take  $i \leftarrow \mathcal{I}(1^\lambda)$  as an input.

A non-interactive argument system for an indexed relation  $\mathcal{R}$  is a tuple of PPT algorithms  $\Pi = (\text{Setup}, \text{P}, \text{V})$ . The setup algorithm  $\text{Setup}(1^\lambda, i)$  produces a common reference string  $\text{crs}$  and a trapdoor  $\text{td}$ . The prover algorithm  $\text{P}(\text{crs}, x, w)$  produces a proof  $\pi$  for the statement  $(i, x) \in \mathcal{L}$ . The verifier algorithm  $\text{V}(\text{crs}, x, \pi)$  decides if  $\pi$  is a valid proof for a statement  $(i, x)$  by outputting either 0 or 1. Notice that  $\text{P}$  and  $\text{V}$  are not directly given  $i$  as an input and instead get a  $\text{crs}$  which depends on  $i$ . This allows to potentially compress the index description by preprocessing.

We require that  $\Pi$  satisfies the following two properties.

**Completeness.** For all  $(i, x, w) \in \mathcal{R}$ ,  $\Pr[(\text{crs}, \text{td}) \leftarrow \text{Setup}(1^\lambda, i), \pi \leftarrow \text{P}(\text{crs}, x, w) : \text{V}(\text{crs}, x, \pi) = 1] = 1$ .

**Soundness.** For all non-uniform PPT adversaries  $\mathcal{A}$ ,

$$\Pr \left[ \begin{array}{l} i \leftarrow \mathcal{I}(1^\lambda), (\text{crs}, \text{td}) \leftarrow \text{Setup}(1^\lambda, i) \\ (x, \pi) \leftarrow \mathcal{A}(1^\lambda, i, \text{crs}) \end{array} : \begin{array}{l} \text{V}(\text{crs}, x, \pi) = 1 \wedge \\ (i, x) \notin \mathcal{L} \end{array} \right] = \text{negl}(\lambda) .$$

In some parts of the paper (where it does not matter), we drop  $i$  and  $\mathcal{I}(1^\lambda)$  from the definitions for simplicity. However, index plays a crucial role in section 5.

We call an argument system a SNARG (succinct non-interactive argument) if additionally the following holds.

**Proof succinctness.** [32] Exists a constant  $c < 1$  such that the length of the proof  $\pi$  is bounded by  $\text{suc}_c(\lambda, |x|, |w|) := \text{poly}(\lambda) \cdot (|x| + |w|)^c$ .

Various other succinctness definitions can be found from the literature. We occasionally discuss two other forms of succinctness.

**Verifier succinctness.** Exists a constant  $c < 1$  such that the verifier's running time is bounded by  $\text{poly}(|x| + \text{suc}_c(\lambda, |x|, |w|))$ .

**CRS succinctness.** CRS size is  $\text{poly}(\lambda)$ . Importantly, CRS size is independent of  $|i|$ .

For example, [16, 19, 39, 55, 57] are proof and verifier succinct but not CRS succinct.

### 3 On Adaptively-Secure Black-Box Extraction

A folklore understanding is that if an argument has black-box knowledge soundness (i.e., there is an efficient algorithm  $\text{Ext}$  that can recover a witness from a proof by using a trapdoor and  $\text{Ext}$  is independent of adversary's code), then the proof has to be "as long as the witness". It is easy to see that such a statement is only partially accurate. Consider an argument system for some relation  $\mathcal{R}_{\mathcal{L}}$  where  $\mathcal{L}$  is an **NP**-language. The same argument system works for a modified relation  $\mathcal{R}'_{\mathcal{L}} = \{(x, w \| 0^k) : (x, w) \in \mathcal{R}_{\mathcal{L}}\}$  where the witness is padded with  $k$  zeroes for an arbitrary number  $k$ . An extractor  $\text{Ext}$  for  $\mathcal{R}'_{\mathcal{L}}$  needs to append  $0^k$  to the witness it extracts for  $\mathcal{R}_{\mathcal{L}}$ . Notably, the proof length for  $\mathcal{R}'_{\mathcal{L}}$  remains the same as for  $\mathcal{R}_{\mathcal{L}}$  independently of witness padding length. This section correctly formalizes the folklore result about the proof size and witness length by associating the hardness of finding the witness with the size of the argument.

We begin by recalling the definition of black-box knowledge soundness.

**Black-box knowledge soundness.** An argument system is black-box  $\varepsilon_{ks}(\lambda)$ -knowledge sound for a relation  $\mathcal{R}$  if there exists a PPT extractor  $\text{Ext}$ , such that for any PPT adversary  $\mathcal{A}$ ,

$$\Pr \left[ \begin{array}{l} (\text{crs}, \text{td}) \leftarrow \text{Setup}(1^\lambda), (x, \pi) \leftarrow \mathcal{A}(\text{crs}) \\ w \leftarrow \text{Ext}(\text{crs}, \text{td}, x, \pi) \end{array} : \begin{array}{l} \text{V}(\text{crs}, x, \pi) = 1 \wedge \\ (x, w) \notin \mathcal{R} \end{array} \right] \leq \varepsilon_{ks}(\lambda) .$$

We say the argument system is black-box knowledge sound if  $\varepsilon_{ks}(\lambda) = \text{negl}(\lambda)$ .

We prove that if the witness of the language cannot be guessed, except for probability  $\varepsilon$ , then the proof size must be at least  $-\log(\varepsilon + \varepsilon_{ks})$  bits long. We start by formalizing the witness guessing probability.



**Definition 3.** Let  $\mathcal{L}$  be an NP language and  $\mathcal{R}_{\mathcal{L}}$  a corresponding relation. We say that an efficiently samplable distribution  $\mathcal{D}_{\mathcal{L}}$  over  $\mathcal{L}$  is  $\varepsilon(\lambda)$ -witness-hard for a relation  $\mathcal{R}_{\mathcal{L}}$  if for any PPT guesser  $\mathcal{M}$ , and any security parameter  $\lambda \in \mathbb{N}$ ,

$$\Pr[x \leftarrow \mathcal{D}(1^\lambda), w \leftarrow \mathcal{M}(1^\lambda, x) : (x, w) \in \mathcal{R}_{\mathcal{L}}] \leq \varepsilon(1^\lambda) .$$

**Theorem 1.** Suppose an efficiently samplable distribution  $\mathcal{D}_{\mathcal{L}}$  over some NP language  $\mathcal{L}$  is  $\varepsilon(\lambda)$ -witness-hard for a relation  $\mathcal{R}_{\mathcal{L}}$ . Let  $\Pi$  be an argument system that has (perfect) completeness and black-box  $\varepsilon_{ks}(\lambda)$ -knowledge soundness. Then the argument size of  $\Pi$  is at least  $-\log(\varepsilon(\lambda) + \varepsilon_{ks}(\lambda))$  bits.

*Proof.* Suppose that  $\Pi$  is an argument system with black-box extractor  $\text{Ext}$  and the argument size is bounded by  $p(\lambda)$  bits. We construct a witness-guesser  $\mathcal{M}^*$  (see fig. 1), which picks a  $\text{crs}$  and an extraction key  $\text{td}$  and guesses uniformly randomly a proof  $\pi$  of size  $p(\lambda)$  bits. It then returns the output of the black-box witness extractor  $\text{Ext}(\text{crs}, \text{td}, x, \pi)$ .

Let us analyze the success probability  $\varepsilon_{\mathcal{M}^*}$  of  $\mathcal{M}^*$  in the witness-hardness game against  $\mathcal{D}_{\mathcal{L}}$ . Let  $\mathcal{E}$  be the distribution  $(x, w, \text{crs}, \pi)$  obtained by running  $x \leftarrow \mathcal{D}(1^\lambda)$  and  $w \leftarrow \mathcal{M}^*(1^\lambda, x)$  ( $\text{crs}$  and  $\pi$  are generated inside  $\mathcal{M}^*$ ). Then,

$$\begin{aligned} \varepsilon_{\mathcal{M}^*} &= \Pr [(x, w, \text{crs}, \pi) \leftarrow \mathcal{E}(1^\lambda) : (x, w) \in \mathcal{R}_{\mathcal{L}}] \\ &\geq \Pr [(x, w, \text{crs}, \pi) \leftarrow \mathcal{E}(1^\lambda) : (x, w) \in \mathcal{R}_{\mathcal{L}} \wedge \mathbf{V}(\text{crs}, x, \pi) = 1] \\ &= \Pr [(x, w, \text{crs}, \pi) \leftarrow \mathcal{E}(1^\lambda) : (x, w) \in \mathcal{R}_{\mathcal{L}} \mid \mathbf{V}(\text{crs}, x, \pi) = 1] \\ &\quad \cdot \Pr [(x, w, \text{crs}, \pi) \leftarrow \mathcal{E}(1^\lambda) : \mathbf{V}(\text{crs}, x, \pi) = 1] . \end{aligned}$$

Let us separately analyze

$$\varepsilon_1 := \Pr [(x, w, \text{crs}, \pi) \leftarrow \mathcal{E}(1^\lambda) : \mathbf{V}(\text{crs}, x, \pi) = 1]$$

and

$$\varepsilon_2 := \Pr [(x, w, \text{crs}, \pi) \leftarrow \mathcal{E} : (x, w) \in \mathcal{R}_{\mathcal{L}} \mid \mathbf{V}(\text{crs}, x, \pi) = 1]$$

For  $\varepsilon_1$ : since  $x \in \mathcal{L}$ , by perfect completeness there exists at least one proof of size at most  $p(\lambda)$  bits that is accepted by the verifier. Thus,  $\varepsilon_1 \geq 1/2^{p(\lambda)}$ . In order to lower bound  $\varepsilon_2$ , we construct an adversary  $\mathcal{B}$  against black-box knowledge soundness. The adversary  $\mathcal{B}$ , described in fig. 1, outputs  $x \leftarrow \mathcal{D}_{\mathcal{L}}$  and a randomly sampled proof  $\pi \leftarrow \{0, 1\}^{p(\lambda)}$ . By inlining  $\mathcal{B}$  into the black-box knowledge soundness game, we get  $\Pr[(x, w, \text{crs}, \pi) \leftarrow \mathcal{E}(1^\lambda) : \mathbf{V}(\text{crs}, x, \pi) = 1 \wedge (x, w) \notin \mathcal{R}_{\mathcal{L}}] \leq \varepsilon_{ks}(\lambda)$ . That is

$$\begin{aligned} &\Pr[(x, w, \text{crs}, \pi) \leftarrow \mathcal{E}(1^\lambda) : \mathbf{V}(\text{crs}, x, \pi) = 1 \wedge (x, w) \notin \mathcal{R}_{\mathcal{L}}] \\ &= \Pr[(x, w, \text{crs}, \pi) \leftarrow \mathcal{E}(1^\lambda) : (x, w) \notin \mathcal{R}_{\mathcal{L}} \mid \mathbf{V}(\text{crs}, x, \pi) = 1] \\ &\quad \cdot \Pr[(x, w, \text{crs}, \pi) \leftarrow \mathcal{E}(1^\lambda) : \mathbf{V}(\text{crs}, x, \pi) = 1] \\ &\geq \Pr[(x, w, \text{crs}, \pi) \leftarrow \mathcal{E}(1^\lambda) : (x, w) \notin \mathcal{R}_{\mathcal{L}} \mid \mathbf{V}(\text{crs}, x, \pi) = 1] \cdot \frac{1}{2^{p(\lambda)}} . \end{aligned}$$

Thus,  $\Pr[(x, w, \text{crs}, \pi) \leftarrow \mathcal{E}(1^\lambda) : (x, w) \notin \mathcal{R}_{\mathcal{L}} \mid \mathbf{V}(\text{crs}, x, \pi) = 1] \leq \varepsilon_{ks}(\lambda) \cdot 2^{p(\lambda)}$ , which means that  $\varepsilon_2 > 1 - \varepsilon_{ks}(\lambda) \cdot 2^{p(\lambda)}$ .

By combining those results, we get that  $\varepsilon(\lambda) \geq \varepsilon_{\mathcal{M}^*} > \frac{1}{2^{p(\lambda)}} \cdot (1 - \varepsilon_{ks}(\lambda) \cdot 2^{p(\lambda)}) = \frac{1}{2^{p(\lambda)}} - \varepsilon_{ks}(\lambda)$ . It follows that  $\varepsilon(\lambda) + \varepsilon_{ks}(\lambda) > \frac{1}{2^{p(\lambda)}}$ , which we can rewrite as  $p(\lambda) > -\log(\varepsilon(\lambda) + \varepsilon_{ks}(\lambda))$ .  $\square$

To understand this claim better, let us consider for simplicity that  $\varepsilon_{ks}(\lambda) = 0$ . Then if  $\varepsilon = \frac{1}{2^{k(\lambda)}}$ , we obtain the lower bound  $p(\lambda) \geq -\log(\frac{1}{2^{k(\lambda)}} + 0) = k(\lambda)$ . In one extreme case, we can imagine that the best PPT witness guesser is no better than an algorithm that guesses the witness at random, i.e.,  $\varepsilon(\lambda) = 1/|w|$ . Then we would get the folklore result that  $p(\lambda) = |\pi| \geq |w|$ . In the other extreme, suppose that the language is in P, in which case  $\varepsilon(\lambda) = 1$ . Then we get that  $-\log(\varepsilon) = 0$ , which fits the intuition that there is no need to communicate a proof for languages in P. However, in a typical situation (where we have some hard language), the lower bound falls somewhere between those extremes.

A closely related but somewhat less precise result can be directly concluded from leakage-resilient OWFs by viewing a proof as leakage on the witness. We show the proof of this impossibility using LR-OWFs in appendix B.

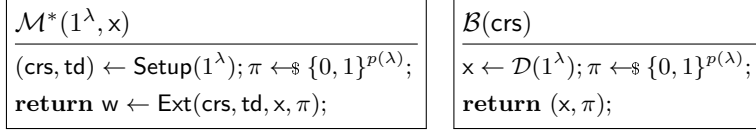


Fig. 1: A witness guessing algorithm  $\mathcal{M}^*$  for  $\mathcal{R}_{\mathcal{L}}$  and a knowledge soundness adversary  $\mathcal{B}$

## 4 Non-Adaptive Black-Box Knowledge Soundness

This section defines non-adaptive black-box knowledge soundness and shows our positive result for **FewP** and a negative result for **NP**.

Below we define non-adaptive black-box *knowledge-soundness*. To the best of our knowledge it has not appeared in prior literature.

**Definition 4 (Non-adaptive Black-box Knowledge Soundness.)** *An argument system is non-adaptive black-box  $\varepsilon_{ks}(\lambda)$ -knowledge sound for a relation  $\mathcal{R}$  if there exists a non-uniform PPT extractor  $\text{Ext}$ , such that for any non-uniform PPT adversary  $\mathcal{A} = (\mathcal{A}_{inp}, \mathcal{A}_{prf})$ ,*

$$\Pr \left[ \begin{array}{l} (x, \text{st}) \leftarrow \mathcal{A}_{inp}(1^\lambda), (\text{crs}, \text{td}) \leftarrow \text{Setup}(1^\lambda) \\ \pi \leftarrow \mathcal{A}_{prf}(\text{st}, \text{crs}), w \leftarrow \text{Ext}^{\mathcal{A}_{prf}(\text{st}, \cdot)}(\text{crs}, \text{td}, x, \pi) \end{array} : \begin{array}{l} \mathbf{V}(\text{crs}, x, \pi) = 1 \\ \wedge (x, w) \notin \mathcal{R} \end{array} \right] \leq \varepsilon_{ks}(\lambda) .$$

We say that the argument system is (non-adaptively) black-box knowledge sound if  $\varepsilon_{ks}(\lambda) = \text{negl}(\lambda)$ .

*Remark 1.* The adversary in definition 4 is stateful only between the input-challenge stage and the proof-challenge stage (through  $\text{st}$ ), but not otherwise. We also assume that on each query  $\mathcal{A}_{prf}(\text{st}, \cdot)$  gets fresh random coins.

### 4.1 A Construction for FewP

In this section, we show that, under the existence of fully homomorphic encryption, universal hash functions and SNARGs (not necessarily of knowledge) for a specific complexity class  $K$ , there exists a non-adaptively secure SNARK with black-box extraction for  $K$ <sup>11</sup>. We can obtain non-adaptive black-box knowledge soundness for a non-trivial subset of **NP** called **FewP**. The class **FewP** can be described as the class of languages admitting at most a polynomial number of witnesses. We remark that if one-way permutations exist, then  $\mathbf{P} \neq \mathbf{FewP}$ <sup>12</sup>. One example of a natural application of a SNARK for **FewP** is proving knowledge of  $w$  such that  $\mathcal{R}(w)$  is satisfied (for arbitrary relation  $\mathcal{R}$ ) and  $w$  opens a perfectly binding commitment.

Further preliminaries for this section can be found in appendix A, where we define the non-adaptive soundness of SNARG (which simply adapts definition 4 to the non-extractable case) and the standard definitions of fully homomorphic encryption (FHE) and universal hash functions (UHF) which will be tools in our construction.

We present our extractable construction in fig. 2<sup>13</sup>. As discussed in the introduction, its main intuition is that the prover provides a (ciphertext containing a) bit of the witness together with the proof. The index for which it is providing such a bit must be somehow hidden. This intuitively prevents the adversary from acting differently for different bits (e.g., using different valid witnesses). This allows us to extract by repeatedly asking the prover for a proof referring to a different index. To achieve the latter, we use an FHE scheme (see also remark 2). When extracting, we will need to keep track of what witness we are extracting (since

<sup>11</sup> This class should include FHE encryption and UHF and should be closed under conjunction. In our theorem statement, we simply require a SNARG for **NP**.

<sup>12</sup> More generally, if poly-to-one one-way functions exist then  $\mathbf{P} \neq \mathbf{FewP}$  [2].

<sup>13</sup> A slightly simpler construction for the case of **UP** (**NP** statements with a unique witness) is in appendix D.

<p><b>Setup</b>(<math>1^\lambda</math>)</p> <hr/> <p> <math>(\hat{c}rs, \hat{t}d) \leftarrow \Pi_{\exists}.\text{Setup}(1^\lambda)</math>  <math>i^* \leftarrow_{\\$} [N_w]</math>  <math>(pk, sk) \leftarrow \text{PKE.KG}(1^\lambda)</math>  <math>(pk_{\text{FHE}}, sk_{\text{FHE}}) \leftarrow \text{FHE.KG}(1^\lambda)</math>  <math>ct_{i^*} \leftarrow \text{FHE.Enc}(pk, i^*)</math>  <math>hk \leftarrow_{\\$} \mathcal{K}_{\text{UHF}}</math>  <b>return</b> <math>(crs := (\hat{c}rs, ct_{i^*}, hk, pk, pk_{\text{FHE}}), td := (sk, sk_{\text{FHE}}, \hat{t}d))</math> </p>
<p><b>P</b>(<math>crs, \mathcal{R}, x, w</math>)</p> <hr/> <p> <math>ct_{\text{bit}} \leftarrow \text{FHE.Eval}(pk_{\text{FHE}}, f_{\text{proj}}, ct_{i^*}, w)</math>          where <math>f_{\text{proj}}(i, w) := w_i</math>  <math>\pi \leftarrow \Pi_{\exists}.\text{P}(\hat{c}rs, \mathcal{R}', (x, hk, pk, pk_{\text{FHE}}, ct_h, ct_{i^*}, ct_{\text{bit}}), w)</math>          where <math>\mathcal{R}'(x, hk, pk, pk_{\text{FHE}}, ct_h, ct_{i^*}, ct_{\text{bit}}; w, r) \iff</math>  <math>\mathcal{R}(x, w) \wedge ct_h = \text{PKE.Enc}(pk, H_{hk}(w), r) \wedge ct_{\text{bit}} = \text{FHE.Eval}(pk_{\text{FHE}}, f_{\text{proj}}, ct_{i^*}, w)</math>  <b>return</b> <math>\pi^* := (\pi, ct_h, ct_{\text{bit}})</math> </p>
<p><b>V</b>(<math>crs, \mathcal{R}, x, \pi^*</math>)</p> <hr/> <p>         Parse <math>\pi^*</math> as <math>(\pi, ct_h, ct_{\text{bit}})</math>  <b>return</b> <math>\Pi_{\exists}.\text{V}(\hat{c}rs, \mathcal{R}', (x, hk, pk, pk_{\text{FHE}}, ct_h, ct_{i^*}, ct_{\text{bit}}))</math>          where <math>\mathcal{R}'</math> is defined like above       </p>

Fig. 2: Non-adaptively secure black-box extractable construction for **FewP**.  $N_w$  is a bound on the witness size.  $\Pi_{\exists}$  is the SNARG scheme.

there could be several). We do this by using a fingerprint through a universal hash function. We encrypt the hash with some (not necessarily homomorphic) cryptosystem to obtain zero-knowledge (ZK). In fact, our soundness proof requires that the prover encrypts the hash with a different public key than the witness bit.

We denote FHE encryptions of a message  $x$  (with an implicit public-key that should be clear from the context) through double brackets  $\llbracket x \rrbracket$ .

**THE EXTRACTOR FOR FewP.** The extractor is presented in fig. 3. It works by collecting different bits of the witness by decrypting  $ct_b$  (the ciphertext returned by the prover) and storing it in some table indexed by the corresponding hash. The crucial point is that there is only a polynomial number of witnesses and thus the extractor can (in the worst-case) “fingerprint” them all. Hashing the witness (through a universal hash function) keeps the proof succinct. We finally encrypt the latter hash to achieve zero-knowledge.

**Theorem 2.** *If  $\Pi_{\exists}$  is a non-adaptively sound SNARG scheme for NP, FHE is a semantically secure FHE scheme, PKE is a semantically secure cryptosystem with perfect correctness and H is a family of UHFs, then the construction in fig. 2 is a non-adaptive SNARK for FewP satisfying definition 4. If  $\Pi_{\exists}$  has additionally computational ZK, then so does the resulting SNARK.*

The proof of the theorem above is in appendix C.1.

*Remark 2 (On replacing FHE with PIR).* While our construction is described with FHE, we observe that the assumption of FHE can easily be replaced by the milder existence of Private Information Retrieval (or PIR) [20].

$\mathcal{E}(\text{crs}, \text{td}, \mathbf{x}, \pi)$	$\text{Qldx}(\mathbf{x}, j)$
Initialize empty table $W$	$\llbracket j \rrbracket \leftarrow \text{FHE.Enc}(\text{pk}, j)$
Retrieve $(\text{crs}, \text{pk}, \text{pk}_{\text{FHE}}, \text{sk}_{\text{FHE}}, \text{hk})$ from $\text{crs}, \text{td}$	Let $\text{crs}_j := (\text{crs}, \llbracket j \rrbracket, \text{hk}, \text{pk}_{\text{FHE}})$
<b>for</b> $j^* = 1, \dots, N_w$	<b>for</b> $k = 1, \dots, N_q = \text{poly}(\lambda)$
Run $\text{Qldx}(\mathbf{x}, j^*)$	Query $\mathcal{A}_{\text{prf}}$ on $(\text{crs}_j, \mathbf{x})$
<b>endfor</b>	obtaining $\pi^* = (\pi, \text{ct}_h, \text{ct}_{\text{bit}})$
Let $h^*$ s.t. $W[h^*][j] \neq \perp$ for all $j$	If proof $\pi$ accepts then
<b>return</b> $W[h^*][1] \dots W[h^*][N_w]$	$b \leftarrow \text{FHE.Dec}(\text{sk}_{\text{FHE}}, \text{ct}_{\text{bit}})$ ;
	else $b \leftarrow \perp$
	Set $W[h][j] \leftarrow b$
	where $h := \text{PKE.Dec}(\text{sk}, \text{ct}_h)$
	<b>endfor</b>

Fig. 3: Extractor for the case for **FewP**

## 4.2 Impossibility for All NP

We now show that the previous constructive result cannot be extended from **FewP** to **NP**. We mentioned at the end of section 3 that we could view the proof as a leakage of the witness and use leakage resilient (LR) cryptography to prove the impossibility of succinct black-box adaptive knowledge soundness. (appendix B) For the non-adaptive case, we can no longer view the SNARK proof as a *one-time* leakage since the extractor (LR adversary) has the ability to rewind the prover and obtain multiple proofs (leakages). Using *continuous leakage resilience*, we extend the impossibility to non-adaptive extraction.

Consider a  $L$ -CLR-OWF  $\Sigma = (\text{KGen}, \text{Sample}, \text{Eval}, \text{Update})$ . We define a relation  $\mathcal{R}_\Sigma = \{((\text{pp}, y), w) : \text{pp} \in \text{KGen}(1^\lambda), w \in \text{Sample}(\text{pp}), \text{Eval}(\text{pp}, w) = y\}$ . Suppose there is a non-adaptive black-box extractable SNARK for  $\mathcal{R}_\Sigma$ . Let us further assume that the proof size of this SNARK is less than  $L$  bits.

We construct the following adversary  $\mathcal{A}$ . First,  $\mathcal{A}$  samples  $(\text{pp}, \text{uk}) \leftarrow \text{KGen}(1^\lambda)$ , a random  $w$ , and outputs  $((\text{pp}, y = \text{Eval}(\text{pp}, w)), \text{st} = (\text{pp}, \text{uk}, w))$ . Next, the extractor can query  $\mathcal{A}(\text{st}, \cdot)$  with different CRS'es and get proofs for the statement  $(\text{pp}, y) \in \mathcal{L}_{\mathcal{R}_\Sigma}$ . Here we define  $\mathcal{A}$ 's behavior as follows: on each query,  $\mathcal{A}$  updates  $w$ , that is it computes  $w' \leftarrow \text{Update}(\text{uk}, w)$ . Then it creates a proof with  $w'$ ,  $\pi \leftarrow \text{P}(\text{crs}, (\text{pp}, y), w')$ , and returns  $\pi$ . This proof is at most size  $L$ , thus at most  $L$  bits of information about  $w'$  gets leaked. By  $L$ -CLR property it is not possible to recover a witness for  $(\text{pp}, y)$  from this amount of information. Hence, the extractor cannot extract the witness and such a SNARK cannot exist. We show this formally.

**Theorem 3.** *Let  $\Sigma = (\text{KGen}, \text{Sample}, \text{Eval}, \text{Update})$  be an  $L$ -CLR-OWF and let  $\Pi$  be a non-adaptive black-box  $\varepsilon_{ks}(\lambda)$ -knowledge sound argument for  $\mathcal{R}_\Sigma$  as defined above. If the proof size is less than  $L(\lambda)$  bits, then  $L$ -CLR-OWF can be broken with probability  $1 - \varepsilon_{ks}(\lambda)$ .*

The proof of this theorem is in appendix C.2. By combining theorem 3 and theorem 7, we obtain the following result.

**Theorem 4.** *If the discrete logarithm assumption holds in some group, then there exists an NP-language  $\mathcal{L}$  such that any non-adaptive BB knowledge sound argument system for  $\mathcal{R}_\mathcal{L}$  has a proof size  $\Omega(|w|)$  where  $|w|$  is the witness size.*

## 5 GW Impossibility for Preprocessing SNARGs

Careful study of [32] reveals that the CRS generation algorithm of a SNARG in their definition depends only on the security parameter. In other words, the proof separating SNARGs from falsifiable assumptions

assumes that the SNARG is CRS succinct and does not allow preprocessing. Many modern SNARGs have a relatively large CRS, which depends on the size of the index  $i$  (e.g., a circuit description) in some way [16, 19, 29, 31, 39, 57]. This makes it questionable if the impossibility result of Gentry and Wichs extends to such SNARGs. We reprove the impossibility theorem for SNARGs that are not necessarily CRS-succinct.

Let us recall the leakage lemma from [32]. We say that a distribution  $A$  over tuples  $(x, \pi)$  is an augmented distribution of  $X$  if  $x$  is distributed according to  $X$  and  $\pi$  is some arbitrary information, possibly correlated to  $x$ . More formally, we may write  $A$  is the distribution over  $(x, \pi)$  such that  $x \leftarrow X$  and  $\pi \leftarrow f(x)$  where  $f$  is some (randomized and possibly inefficiently computable) function.

**Lemma 1 (Leakage lemma [32]).** *There exists a polynomial  $p$  for which the following holds. Let  $X_\lambda$  and  $\bar{X}_\lambda$  be two distributions that are  $(s(\lambda), \varepsilon(\lambda))$ -indistinguishable. Let  $A_\lambda$  over  $(x, \pi)$  be an augmented distribution of  $X_\lambda$ , where  $|\pi| = \ell(\lambda)$ . Then there exist an augmented distribution  $\bar{A}_\lambda$  of  $\bar{X}_\lambda$  such that  $A_\lambda$  and  $\bar{A}_\lambda$  are  $(s^*(\lambda), \varepsilon^*(\lambda))$ -indistinguishable where  $s^*(\lambda) = s(\lambda)p(\varepsilon(\lambda)/2^{\ell(\lambda)})$  and  $\varepsilon^*(\lambda) = 2\varepsilon(\lambda)$ .*

We also present some definitions which help to prove the main result.

**Definition 5 (Breaking Adaptive Soundness [56]).** *We say that an algorithm  $\mathcal{A}$  breaks adaptive soundness of a SNARG  $\Pi$  for a language  $\mathcal{L}$  with probability  $\varepsilon(\cdot)$  if there exists an index  $i \in \mathcal{I}$  such that for every  $\lambda \in \mathbb{N}$ ,*

$$\Pr[\text{crs} \leftarrow \text{Setup}(1^\lambda, i), (x, \pi) \leftarrow \mathcal{A}(1^\lambda, \text{crs}) : (i, x) \notin \mathcal{L} \wedge \text{Verify}(\text{crs}, x, \pi) = 1] \geq \varepsilon(\lambda).$$

Note that if  $\varepsilon(\lambda)$  is non-negligible, then adaptive soundness cannot be satisfied.

**Definition 6 (Soundness Reduction [56]).** *We say that a PPT machine  $R$  is a black-box reduction for adaptive soundness of an argument  $\Pi$  based on a falsifiable assumption  $(\mathcal{C}, c)$  if there exists a polynomial  $p(\cdot, \cdot)$  such that for every  $\mathcal{A}$  that breaks adaptive soundness with probability  $\varepsilon(\cdot)$ , for every  $\lambda \in \mathbb{N}$ ,  $R^\mathcal{A}(1^\lambda)$  wins  $(\mathcal{C}, c)$  with a probability at least  $p(\varepsilon(\lambda), 1/\lambda)$ .*

We say that  $R^\mathcal{A}$  is *security-parameter preserving* if additionally there exist a polynomial  $q$  such that  $R^\mathcal{A}(1^\lambda)$  queries  $\mathcal{A}$  with inputs of the form  $(1^\lambda, x \in \{0, 1\}^*)$  and at most  $q(\lambda)$  times.

We start by stating two technical lemmas, which we prove in appendix C.

**Lemma 2.** *If an indexed languages  $\mathcal{L} \in \mathbf{NP}$  has a sub-exponentially hard-on-average problem, then for any  $d > 0$ ,  $\mathcal{L}$  also has a hard-on-average problem with  $(2^{\lambda^d}, 1/2^{\lambda^d})$ -indistinguishability.*

**Lemma 3.** *Let  $X_\lambda$  and  $\bar{X}_\lambda$  be  $(2^{\lambda^d}, 1/2^{\lambda^d})$ -indistinguishable distributions for some integer  $d \geq 2$ . Let  $A_\lambda$  over  $(x, \pi)$  be an augmented distribution of  $X_\lambda$ , where  $|\pi| = \ell(\lambda) = o(\lambda^d)$ . Then there exists an augmented distribution  $\bar{A}_\lambda$  of  $\bar{X}_\lambda$  such that  $A_\lambda$  and  $\bar{A}_\lambda$  are  $(\text{poly}(\lambda), \text{negl}(\lambda))$ -indistinguishable.*

*Remark 3.* Note that  $(s(\lambda), \varepsilon(\lambda)) = (\text{poly}(\lambda), \text{negl}(\lambda))$ -indistinguishability is not enough in the previous lemma because. Suppose  $|\pi| = \ell(\lambda) = \lambda^{d-1}$ . Then,

$$\begin{aligned} s^*(\lambda) &= s(\lambda)p(\varepsilon(\lambda)/2^{\ell(\lambda)}) = \text{poly}(\lambda)p(\text{negl}(\lambda)/2^{\lambda^{d-1}}) \\ &= \text{poly}(\lambda)p(2^{-\omega(\lambda^{d-1})}) = 2^{-\omega(\text{poly}(\lambda))} \end{aligned} ,$$

given that  $p$  is not a constant polynomial. Thus,  $A_\lambda$  and  $\bar{A}_\lambda$  would be (provably) indistinguishable only for very small circuits.

Now we are ready to restate the Gentry-Wichs impossibility result with respect to preprocessing SNARGs from section 2.2.

**Theorem 5.** *Assume that,*

- $\mathcal{L}$  is an indexed language with a sub-exponentially hard-on-average problem (see section 2).
- $\Pi$  is a SNARG for  $\mathcal{L}$ , i.e., it is complete, sound, and proof-succinct (but not necessarily verifier-succinct or CRS-succinct).

$\mathcal{A}^*(1^\lambda, \text{crs}, i)$	$\text{Emul}(1^\lambda, \text{crs}, i)$	$\text{O}_i(1^\lambda, \text{crs}, i)$ //Initially $j = 1$
<b>if</b> $i \notin \mathcal{I}(1^\lambda)$	<b>if</b> $i \notin \mathcal{I}(1^\lambda)$	<b>if</b> $j \leq i$
<b>return</b> $\perp$ ;	<b>return</b> $\perp$ ;	$(x, \pi) \leftarrow \text{Emul}(1^\lambda, \text{crs}, i)$ ;
$(\bar{x}, \bar{\pi}) \leftarrow \bar{A}_{\lambda, i, \text{crs}}$ ;	$(x, w) \leftarrow \text{Samp}_{\mathcal{L}}(1^\lambda, i)$ ;	<b>else</b>
<b>return</b> $(\bar{x}, \bar{\pi})$ ;	$\pi \leftarrow \mathbf{P}(\text{crs}, x, w)$ ;	$(x, \pi) \leftarrow \mathcal{A}^*(1^\lambda, \text{crs}, i)$ ;
	<b>return</b> $(x, \pi)$ ;	$j \leftarrow j + 1$ ;
		<b>return</b> $(x, \pi)$ ;

Fig. 4: Soundness adversary  $\mathcal{A}^*$ , its efficient emulator  $\text{Emul}$ , and hybrid adversaries  $\text{O}_i$

Then, for any falsifiable assumption  $(\mathcal{C}, c)$  either:

- $(\mathcal{C}, c)$  is false or,
- there is no security-parameter preserving black-box reduction for adaptive soundness of  $\Pi$  based on  $(\mathcal{C}, c)$ .

*Proof.* Suppose there exists a security-parameter preserving black-box reduction  $R$  for adaptive soundness of  $\Pi$  based on a falsifiable assumption  $(\mathcal{C}, c)$  and that  $R$  makes at most  $q(\lambda)$  queries to its oracle, where  $q$  is some polynomial. The proof idea is to construct a computationally unbounded adversary  $\mathcal{A}^*$  that is able to break adaptive soundness. Then we show using lemma 1 that there is an efficient emulator  $\text{Emul}$  that gives outputs which are indistinguishable from outputs of  $\mathcal{A}^*$ . Thus, if  $R^{\mathcal{A}^*}(1^\lambda)$  is able to break the assumption  $(\mathcal{C}, c)$ , then so is  $R^{\text{Emul}}(1^\lambda)$  and it follows that  $(\mathcal{C}, c)$  must be false.

Since  $\Pi$  is proof-succinct there exists some  $n$  such that the proof size  $\ell$  is bounded by  $\lambda^n \cdot (|x| + |w|)^{o(1)}$ . Moreover by lemma 2, since we assume that some sub-exponentially hard-on-average problem exists for  $\mathcal{L}$ , there also exists a sub-exponentially hard-on-average problem with  $(2^{\lambda^{n+2}}, 1/2^{\lambda^{n+2}})$ -indistinguishability. Let it be defined by an index sampler  $\mathcal{I}$  and instance samplers  $\text{Samp}_{\mathcal{L}}$  and  $\text{Samp}_{\bar{\mathcal{L}}}$ . It is more convenient to start from describing the emulator  $\text{Emul}$  before we describe  $\mathcal{A}^*$ . The emulator (see also fig. 4) on input  $(1^\lambda, \text{crs}, i)$  checks that  $i$  is well-formed, samples  $(x, w) \leftarrow \text{Samp}_{\mathcal{L}}(1^\lambda, i)$ , creates a proof  $\pi \leftarrow \mathbf{P}(\text{crs}, x, w)$  and returns  $(x, \pi)$ .

Notice that since  $\text{Samp}_{\mathcal{L}}$  runs in polynomial time in  $\lambda$ , then  $|x| = \text{poly}(\lambda)$  and  $|w| = \text{poly}(\lambda)$ . Therefore, the proof size is  $\ell(\lambda) = \lambda^{o(n^{d+2})}$ .

Fix an arbitrary oracle input  $(1^\lambda, \text{crs}, i)$ . Let  $X_{\lambda, i}$  be the distribution of  $x$  from sampling  $(x, w) \leftarrow \text{Samp}_{\mathcal{L}}(1^\lambda, i)$  and  $\bar{X}_{\lambda, i}$  the distribution of  $\bar{x}$  we get by sampling  $\bar{x} \leftarrow \text{Samp}_{\bar{\mathcal{L}}}(1^\lambda, i)$ . As we established, these distributions are  $(2^{\lambda^{n+2}}, 1/2^{\lambda^{n+2}})$ -indistinguishable. Let  $A_{\lambda, i, \text{crs}}$  be the augmented distribution of  $X_{\lambda, i}$  defined as  $(x, \pi) \leftarrow \text{Emul}(1^\lambda, \text{crs}, i)$ . By lemma 3, there exists an augmented distribution  $\bar{A}_{\lambda, i, \text{crs}}$  of  $\bar{X}_{\lambda, i}$  such that  $A_{\lambda, i, \text{crs}}$  and  $\bar{A}_{\lambda, i, \text{crs}}$  are  $(\text{poly}(\lambda), \text{negl}(\lambda))$ -indistinguishable.

Now we can describe the adversary  $\mathcal{A}^*$ . On the query input  $(1^\lambda, \text{crs}, i)$  it simply returns  $(\bar{x}, \bar{\pi}) \leftarrow \bar{A}_{\lambda, i, \text{crs}}$ . Since  $\bar{A}_{\lambda, i, \text{crs}}$  is not necessarily efficiently sampleable,  $\mathcal{A}^*$  may be inefficient.

Our goal is to show that the assumption  $(\mathcal{C}, c)$  is false if  $R$  exists, i.e.,

$$\Pr[R^{\text{Emul}}(1^\lambda) \text{ wins } (\mathcal{C}, c)] > c + \text{negl}(\lambda) .$$

We show this in two parts.

1)  $R^{\mathcal{A}^*}$  wins  $(\mathcal{C}, c)$ : First, let  $\varepsilon_{\mathcal{A}^*}(\lambda)$  be the probability that  $\mathcal{A}^*$  breaks adaptive soundness of  $\Pi$ ,

$$\varepsilon_{\mathcal{A}^*}(\lambda) := \Pr \left[ \begin{array}{l} i \leftarrow \mathcal{I}(1^\lambda), \text{crs} \leftarrow \text{Setup}(1^\lambda, i) \\ (x, \pi) \leftarrow \mathcal{A}^*(1^\lambda, \text{crs}, i) \end{array} : (i, x) \notin \mathcal{L} \wedge \mathbf{V}(\text{crs}, x, \pi) = 1 \right] .$$

Let us first only consider the probability of the verifier accepting a proof,

$$\varepsilon_{\mathbf{Vf}=1}(\lambda) := \Pr \left[ \begin{array}{l} i \leftarrow \mathcal{I}(1^\lambda), \text{crs} \leftarrow \text{Setup}(1^\lambda, i), \\ (x, \pi) \leftarrow \mathcal{A}^*(1^\lambda, \text{crs}, i) \end{array} : \mathbf{V}(\text{crs}, x, \pi) = 1 \right] .$$

Due to completeness, we know that  $\varepsilon_{\text{Emul}}(\lambda) = 1$ , where

$$\varepsilon_{\text{Emul}}(\lambda) := \Pr \left[ \begin{array}{l} i \leftarrow \mathcal{I}(1^\lambda), \text{crs} \leftarrow \text{Setup}(1^\lambda, i), \\ (\mathbf{x}, \pi) \leftarrow \text{Emul}(1^\lambda, \text{crs}, i) \end{array} : \mathbf{V}(\text{crs}, \mathbf{x}, \pi) = 1 \right].$$

Since  $\mathbf{V}$  can be seen as a polynomial-size distinguisher for  $A_{\lambda, i, \text{crs}}$  and  $\bar{A}_{\lambda, i, \text{crs}}$ , we get from before that  $|\varepsilon_{\text{Emul}}(\lambda) - \varepsilon_{\mathbf{V}=1}(\lambda)| \leq \text{negl}(\lambda)$ . Therefore,  $1 - \text{negl}(\lambda) \leq \varepsilon_{\mathbf{V}=1}$ . Since  $\Pr[i \leftarrow \mathcal{I}(1^\lambda), \text{crs} \leftarrow \text{Setup}(1^\lambda, i), (\mathbf{x}, \pi) \leftarrow \mathcal{A}^*(\text{crs}, i) : (i, \mathbf{x}) \notin \mathcal{L}] = 1$ ,  $\varepsilon_{\mathcal{A}^*} = \varepsilon_{\mathbf{V}=1} \geq 1 - \text{negl}(\lambda)$ . Thus,  $\mathcal{A}^*$  breaks adaptive soundness with an overwhelming probability. Since we assumed a black-box reduction  $R$ , there must exist a polynomial  $p(\cdot, \cdot)$  such that  $R^{\mathcal{A}^*}(1^\lambda)$  breaks  $(\mathcal{C}, c)$  with probability at least  $p(1 - \text{negl}(\lambda), 1/\lambda)$ .

2)  $R^{\mathcal{A}^*}$  is indistinguishable from  $R^{\text{Emul}}$ : Suppose  $R$  makes  $q(\lambda)$  oracle queries. Let  $\mathcal{O}_i$  for  $i \in \{0, \dots, q(\lambda)\}$  denote a stateful algorithm that we describe in the following. The machine  $\mathcal{O}_i$  for the first  $i$  queries responds as  $\text{Emul}$  and for the rest of the queries  $(1^\lambda, \text{crs}, i)$  responds as  $\mathcal{A}^*$  (see fig. 4). In particular  $\mathcal{O}_0 = \mathcal{A}^*$  and  $\mathcal{O}_{q(\lambda)} = \text{Emul}$ . We denote  $\varepsilon_i := \Pr[R^{\mathcal{O}_i}(1^\lambda)$  wins  $(\mathcal{C}, c)$ ]. We can again use indistinguishability of  $A_{\lambda, i, \text{crs}}$  and  $\bar{A}_{\lambda, i, \text{crs}}$  to show that  $|\varepsilon_i - \varepsilon_{i+1}| \leq \text{negl}(\lambda)$ . Therefore, by triangle inequality  $|\varepsilon_0 - \varepsilon_{q(\lambda)}| \leq q(\lambda)\text{negl}(\lambda) = \text{negl}(\lambda)$ .

Since  $\varepsilon_0 = \varepsilon_{\mathcal{A}^*}$ , we get that  $\Pr[R^{\text{Emul}}(1^\lambda)$  wins  $(\mathcal{C}, c)] = \varepsilon_{q(\lambda)} \geq \varepsilon_{\mathcal{A}^*} - \text{negl}(\lambda) = p(1 - \text{negl}(\lambda), 1/\lambda) - \text{negl}(\lambda)$ . Thus,  $R^{\text{Emul}}(1^\lambda)$  can break the assumption  $(\mathcal{C}, c)$  with an overwhelming probability.  $\square$

## 6 Understanding SNARG Impossibilities

In this section, we attempt to provide a complete overview of known impossibilities for non-interactive arguments. This illustrates the precise assumptions behind these impossibilities in order to identify avenues for further research. The following are some of the major impossibility results.

1. Gentry-Wichs [32]: Adaptive soundness of a SNARG cannot be proven via a black-box reduction to a falsifiable assumption.
2. Pass [56]: Adaptive soundness of a statistical NIZK argument cannot be proven via a black-box reduction to a falsifiable assumption.
3. Groth [39]: Any pairing-based SNARK obtained from a NILP (a non-interactive linear proof) must contain at least 2 group elements, one in each of the pairing source groups.

Since [39] is relevant only in a very specific setting, and [56] is about general NIZKs, we will not focus on it in the rest of the paper. The proof idea of [32] is quite similar to our extension of it. We recall the proof idea of [56] in appendix F and the proof idea of [32] in appendix E. In the following, we discuss the impossibility result of [32] and then outline the landscape of positive and negative results in Table 1.

### 6.1 Impossibility of Gentry-Wichs

We recall the main result of [32].

**Theorem 6.** *Let  $\mathcal{L}$  be a sub-exponentially hard NP language and let  $\Pi$  be a SNARG for  $\mathcal{L}$ , satisfying completeness and proof succinctness. Then, for any falsifiable assumption  $(\mathcal{C}, c)$ , either  $(\mathcal{C}, c)$  is false, or there is no black-box reduction showing the (adaptive) soundness of  $\Pi$  based on  $(\mathcal{C}, c)$ .*

We take a closer look at GW impossibility and enumerate the scenarios to which it *does not* apply. While some of these are known results, they are all scattered in the literature.

- **Non-adaptive soundness:** The impossibility holds only for *adaptive soundness*. The proof technique used in GW to rule out a black-box reduction uses a stateless adversary that outputs an instance proof pair  $(\mathbf{x}, \pi)$  on input a CRS. In particular, this does not rule out our reductions that can rewind the prover and obtain different proofs for the same  $\mathbf{x}$ , which is possible in the case of non-adaptive soundness. Recent work attempted to show tightness from new albeit falsifiable assumptions, but the construction was shown to be faulty [60].

- **Low-space non-deterministic computation:** The high-level idea of the GW impossibility result is a “leakage lemma” that says the following: assuming the underlying NP language is  $2^\ell$ -hard, a reduction that breaks the assumption, cannot distinguish between pairs  $(x, \pi)$  generated by a (possibly inefficient) cheating prover, where  $x \notin L$  and  $\pi$  is a proof of length  $\ell$ , and a pair  $(\tilde{x}, \tilde{\pi})$  where  $\tilde{x} \in L$  and  $\tilde{\pi}$  is an efficiently generated proof. Therefore, for computations recognizable in  $\text{poly}(\lambda)$  time and  $S(\lambda)$  space with a non-deterministic Turing machine (the class  $NTISP(\text{poly}(\lambda), S(\lambda))$ ), the GW result does not rule out the possibility of a SNARG with proofs of length  $\text{poly}(\lambda)(S(\lambda))$ , since a computation in  $NTISP(\text{poly}(\lambda), S(\lambda))$  is in  $DTIME(\text{poly}(\lambda) \cdot 2^{S(\lambda)})$  which is not  $\text{poly}(\lambda) \cdot 2^{O(S(\lambda))}$ -hard. The work of [4] constructs a delegation scheme for non-deterministic computations with a proof length that grows only with the space of the computation.
- **Preprocessing SNARGs:** The GW separation result holds for SNARGs that have a “short” CRS. More precisely, the impossibility proof requires that the size of the CRS depends only on the security parameter, and does not grow with the size of the instance. This gap is now closed with the current work.
- **Trapdoor languages:** For some languages there are efficient proofs [42, 49] in the quasi-adaptive setting (QA-NIZK). These proofs have a constant number group elements – regardless of the instance size. The construction of [47] for languages consisting of linear subspaces of a vector space, have constant-sized proofs, achieve adaptive soundness (based on a falsifiable assumption) and perfect ZK. This seemingly contradicts the GW impossibility result (as well as the impossibility on perfect ZK [56]).<sup>14</sup> The results in the quasi-adaptive setting do not contradict the GW impossibility because the CRS hides a trapdoor for deciding membership in the language. The proof of GW rules out reductions that *cannot* efficiently detect when the soundness property is broken. We formalize this notion of *trapdoor languages* below.

**Bypassing GW: Trapdoor Languages.** Intuitively, a trapdoor language allows verifying membership in the language if one knows a trapdoor. Towards formalizing such languages, we illustrate it by taking the linear subspace language as an example. Recall the language of linear subspaces from [47]. We have a distribution  $\mathcal{D}$  that outputs a language parameter  $\text{lpar} = [M]_1 \in \mathbb{G}_1^{n \times m}$ <sup>15</sup> for some matrix  $M$  and the respective linear subspace language is defined as

$$\mathcal{L}_{[M]_1} = \{[\vec{x}]_1 \in \mathbb{G}_1^n \mid \exists \vec{w} \in \mathbb{Z}_p^m : \vec{x} = M \cdot \vec{w}\}.$$

In the proof of [32], the reduction algorithm  $R$  that picks the CRS, (and in the case of the linear subspace language, also picks the language parameter  $\text{lpar}$ ), should not efficiently distinguish between elements  $x \in \mathcal{L}$  from  $\bar{x} \in \bar{\mathcal{L}}$ . The latter condition does not hold for linear subspace languages.

In particular, we now argue that it is possible to efficiently decide if  $[\vec{x}]_1 \in \mathcal{L}_{[M]_1}$  by knowing  $M$ .

Observe that given both  $M$  and  $\vec{x}$  as integers, by Kronecker–Capelli theorem there exists  $\vec{w} \in \mathbb{Z}_p^m$  such that  $\vec{x} = M\vec{w}$  (i.e.,  $[\vec{x}]_1 \in \mathcal{L}_{[M]_1}$ ) if and only if  $\text{rank}(M) = \text{rank}(M \mid \vec{x})$ . Turns out a similar test can be used even when given only  $[x]_1$  and  $M$ , but some extra care needs to be taken to compute  $\text{rank}(M \mid \vec{x})$ . Firstly, consider a submatrix  $A = (M' \mid x') \in \mathbb{Z}_p^{d \times d}$  of  $(M \mid \vec{x})$  which includes the last column  $\vec{x}$ . By using Laplace expansion, we are able to compute  $[\det(A)]_1 = \sum_{i=1}^d (-1)^{i+d} [x'_i]_1 D_{i,d}$  where  $D_{i,d}$  is a determinant of the submatrix that we get by removing  $i$ -th row and  $d$ -th column from  $A$ . We still do not know  $\det(A)$ , but by comparing  $[\det(A)]_1$  to  $[0]_1$ , we can tell if  $A$  is a singular or a non-singular matrix. Considering that rank of a matrix is the largest order of any of its non-zero minors, we obtain the algorithm in fig. 5 for deciding elements of  $\mathcal{L}_{[M]_1}$ .

In more detail, we first compute rank  $r$  of  $M$ , which can be done efficiently. The rank of  $(M \mid \vec{x})$  can be at most  $r + 1$  since it includes only one extra column. To test this, we iterate over all the  $(r + 1) \times (r + 1)$  submatrices  $A$  of  $(M \mid \vec{x})$  that contain the  $\vec{x}$  column and compute  $[\det(A)]_1$ . If one of the determinants is non-zero, then  $\text{rank}(M \mid \vec{x}) = r + 1$  and it follows that  $[\vec{x}]_1 \notin \mathcal{L}_{[M]_1}$ . Otherwise,  $\text{rank}(M \mid \vec{x}) = r = \text{rank}(M)$  and  $[\vec{x}]_1 \in \mathcal{L}_{[M]_1}$ . In order for  $\mathcal{D}_{\mathcal{L}_{[M]_1}}$  to be efficient, we assume that  $n$  and  $m$  are small constants.

<sup>14</sup> The proof of [47] contains 1 group element and bypasses the [39] impossibility as well. This is not contradictory because the [39] impossibility only applies to pairing-based NIZKs that are compiled from NILPs

<sup>15</sup> Here,  $\mathbb{G}_1$  is an additive pairing group and  $[x]_1$  denotes a group element with a discrete logarithm  $x$ .



$\mathcal{D}_{\mathcal{L}_{[M]_1}}(M, [\vec{x}]_1)$
$r \leftarrow \text{rank}(M);$
<b>for</b> $A = (M' \mid \vec{x}') \in \mathbb{Z}_p^{(r+1) \times (r+1)}$ submatrix of $(M \mid \vec{x})$
$[\det(A)]_1 \leftarrow \sum_{i=1}^d (-1)^{i+d} [x'_i]_1 D_{i,d};$
<b>if</b> $[\det(A)]_1 \neq [0]_1$ : <b>return false</b> ;
<b>return true</b> ;

Fig. 5: Efficient decision algorithm for  $\mathcal{L}_{[M]_1}$ , given access to trapdoor  $M$

As we saw above, the linear subspace language has a trapdoor  $M$  which allows to efficiently recognize language elements and this sufficient to avoid the [32] impossibility.

We now generalize this observation by defining a trapdoor language.

**Definition 7.** Let  $\mathcal{D}(1^\lambda)$  be an efficiently sampleable distribution that outputs  $(\text{lpar}, \text{td})$  and each  $\text{lpar}$  is associated with a language  $\mathcal{L}_{\text{lpar}}$ . A family of languages  $\{\mathcal{L}_{\text{lpar}}\}_{(\text{lpar}, \text{td}) \in \mathcal{D}(1^\lambda), \lambda \in \mathbb{N}}$  are trapdoor languages if there exists a PPT decider  $\mathcal{M}$  such that for all  $\lambda \in \mathbb{N}$  and all  $(\text{lpar}, \text{td}) \in \mathcal{D}(1^\lambda)$ ,

$$x \in \mathcal{L} \Leftrightarrow \mathcal{M}(1^\lambda, \text{lpar}, \text{td}, x) = 1.$$

The security definitions from section 2.2 for non-interactive arguments slightly change in that the Setup takes the language parameter instead of index as input and outputs a CRS. The soundness definition in general is not efficiently falsifiable because checking  $x \notin \mathcal{L}$  is usually not efficient. However, with trapdoor languages it is falsifiable since  $\mathcal{M}$  is efficient. In particular, a tautological assumption “ $\Pi$  is sound” becomes a falsifiable assumption.

**Examples of useful trapdoor languages.** We are interested in “hard” trapdoor languages, i.e., trapdoor languages that are hard to decide without knowledge of  $\text{td}$ . We illustrate a few examples below.

- *Linear subspace language.* Firstly, let us observe that the linear subspace languages fits into the trapdoor language definition. We let  $\mathcal{D}(1^\lambda)$  pick a pairing description  $\text{bp}$  and sample a matrix  $M$  according to some distribution.  $\mathcal{D}(1^\lambda)$  outputs  $\text{lpar} = (\text{bp}, [M]_1)$  and  $\text{td} = M$ . Deciding if  $x \in \mathcal{L}_{[M]_1}$  can be decided efficiently given  $\text{td}$  as we argued before. For many distributions of  $M$ ,  $\mathcal{L}_{[M]_1}$  is considered to be a hard language on average. For example, if  $M = (1, x)^\top$  and  $x, w \leftarrow_{\$} \mathbb{Z}_p$ , then  $[Mw]_1 = (w, wx)$ , which is indistinguishable from a random tuple  $[u, v]_1^\top \leftarrow_{\$} \mathbb{G}_1^2$  under the decisional Diffie-Hellman assumption. More generally, hardness of such distributions is characterized by the matrix decisional Diffie-Hellman (MDDH) assumption [26].
- *Statements about encrypted values.* Many statements about ciphertexts can be naturally formalized as a trapdoor language by using the public key as  $\text{lpar}$  and the secret key as  $\text{td}$ . Consider the following example.  
Let  $(\text{KGen}, \text{Enc}, \text{Dec})$  be a public key cryptosystem for encrypting  $\ell$ -bit messages and let  $C : \{0, 1\}^\ell \rightarrow \{0, 1\}$  be an efficiently computable boolean circuit. We set  $\mathcal{D}(1^\lambda) = \text{KGen}(1^\lambda)$ , that is  $\text{lpar} = \text{pk}$  is the public key and  $\text{td} = \text{sk}$  is the corresponding secret. We define the language as  $\mathcal{L}_{\text{pk}}^C = \{c \mid C(\text{Dec}(\text{sk}, c)) = 1\}$ . In other words,  $\mathcal{L}_{\text{pk}}^C$  contains ciphertexts that encrypt a message  $m$  which satisfy some property characterized by the circuit  $C$ . For example, in range proofs we have  $C$  which checks that  $k_1 \leq m \leq k_2$  for some constants  $k_1$  and  $k_2$ . Clearly,  $\mathcal{L}_{\text{pk}}^C$  is a trapdoor language since given  $\text{sk}$  it is possible to decrypt  $c$  and efficiently check that the plaintext satisfies  $C$ .
- *Shuffle.* Popular ciphertext-based language that fits into the trapdoor language mould is the ciphertext shuffle. We set  $\mathcal{D} = \text{KGen}$ . Let  $\Sigma_n$  be the set of permutations on  $n$  elements. The shuffle language for  $n$

ciphertexts is

$$\mathcal{L}_{\text{pk}}^{n\text{-shuf}} = \{((c_1, \dots, c_n), (c'_1, \dots, c'_n)) \mid \exists \sigma \in \Sigma_n \forall i \in \{1, \dots, n\} : \text{Dec}(\text{sk}, c_i) = \text{Dec}(\text{sk}, c'_{\sigma(i)})\}$$

With the secret key as the trapdoor, statements are easy to verify since one can decrypt both ciphertext vectors, sort the resulting plaintext vectors, and then check their equality.

Shuffle proofs are often used to prove correct behaviour of mix-networks, which have for instance found application in e-voting systems to anonymize ciphertexts of voters [59].

- *Set membership.* Let us consider a public set  $S$  and a public key cryptosystem. A set membership language for  $S$  is defined as  $\mathcal{L}_{\text{pk}}^S = \{c : \text{Dec}(\text{sk}, c) \in S\}$ . This is clearly a trapdoor language where again  $\text{sk}$  plays the role of a trapdoor. González and Ràfols [36] show that an argument for this language (and its aggregated version for multiple ciphertexts) can be used to obtain, for example, shuffle arguments and range arguments. Of course, there are also more direct application like showing that  $c$  encrypts a valid candidate in an e-voting system, where  $S$  is the set of all candidates.

We note that trapdoor languages are interesting, arise frequently in practice and GW impossibility does not apply. Can we construct SNARGs from falsifiable assumptions for trapdoor languages? We leave resolving this as an interesting open question, and believe that our formalization is a first step in identifying the middle ground where GW does not apply but the language remains interesting.

## 6.2 A Complete Picture

In table 1, we give an overview of the impossibility results for non-interactive arguments, and positive results known under various relaxations. As already highlighted above, there are two major impossibility results. Firstly, there is no adaptively sound succinct argument for all non-deterministic computations with a black-box reduction to a falsifiable assumption [32]. This holds even with a designated verifier (a verifier that holds a private verification key). Secondly, there is no statistical zero-knowledge argument (succinct or not) for all non-deterministic computations with a black-box reduction to a falsifiable assumption [56]. Although not mentioned in the original paper, this impossibility result also extends to the designated verifier.<sup>16</sup>

On the other hand, by relaxing some of the requirements, it is possible to achieve succinct arguments and also statistical zero-knowledge arguments. Delegation schemes are adaptively sound succinct arguments for deterministic computation and they are achievable under falsifiable pairing-based and lattice-based assumptions as was shown by [22, 37, 43]. Recently Lipmaa and Pavlyk [52] showed that non-adaptivity is another possible relaxation. They construct a non-adaptively sound SNARG for non-deterministic computation that has perfect zero-knowledge based on a new, but falsifiable assumption. A non-adaptively sound SNARG under a falsifiable assumption was known even prior to [52]. Namely, Sahai and Waters [58] constructed a succinct perfect NIZK argument with non-adaptive soundness from iO. Subsequent to their work, constructions of iO have been proposed which are secure under falsifiable assumption, e.g. [62]. We give a brief overview of this SNARG construction in appendix G.

**Acknowledgements.** We thank the reviewers of CRYPTO 2022 for constructive feedback, in particular, for pointing out the connection between blackbox extractability and leakage-resilient cryptography. We thank Helger Lipmaa for comments on the paper.

## References

1. Agrawal, S., Dodis, Y., Vaikuntanathan, V., Wichs, D.: On continual leakage of discrete log representations. In: ASIACRYPT 2013, Part II. LNCS, vol. 8270, pp. 401–420

<sup>16</sup> As can be seen in appendix F, neither the inefficient soundness adversary  $\mathcal{A}_{\text{slow}}$  nor its emulator  $\mathcal{A}_{\text{fast}}$  need to run the verifier internally and thus the same impossibility proof applies for the designated verifier setting.

Table 1: (Im)possibility results for non-interactive arguments under falsifiable assumptions. BB stands for black-box.

adaptive soundness	public verifier	succinct argument	language class	statistical ZK	notes/citation
+	+/-	+	<b>NP</b>	+/-	No BB reduction [32]
+	+/-	+/-	<b>NP</b>	+	No BB reduction [56]
-	+	+	<b>NP</b>	+	[52] and [58]
+	+	+	<b>P</b>	trivial	[22, 37, 43]
+	+	+	linear subspace	+	[47]
(quasi-adaptive)					
-	+	+	batch <b>NP</b>	-	[21]
+	+	+	non-deterministic bounded space	-	[44]

2. Allender, E.W.: The complexity of sparse sets in  $p$ . In: Structure in Complexity Theory, pp. 1–11
3. Alwen, J., Dodis, Y., Wichs, D.: Leakage-resilient public-key cryptography in the bounded-retrieval model. In: CRYPTO 2009. LNCS, vol. 5677, pp. 36–54
4. Badrinarayanan, S., Kalai, Y.T., Khurana, D., Sahai, A., Wichs, D.: Succinct delegation for low-space non-deterministic computation. In: 50th ACM STOC, pp. 709–721
5. Bagheri, K., Kohlweiss, M., Siim, J., Volkhov, M.: Another look at extraction and randomization of groth’s zk-snark. In: Financial Cryptography and Data Security, pp. 457–475
6. Bagheri, K., Sedaghat, M.: Tiramisu: black-box simulation extractable nizks in the updatable crs model. In: International Conference on Cryptology and Network Security, pp. 531–551
7. Barta, O., Ishai, Y., Ostrovsky, R., Wu, D.J.: On succinct arguments and witness encryption from groups. In: CRYPTO 2020, Part I. LNCS, vol. 12170, pp. 776–806
8. Ben-Or, M., Goldreich, O., Goldwasser, S., Håstad, J., Kilian, J., Micali, S., Rogaway, P.: Everything provable is provable in zero-knowledge. In: CRYPTO’88. LNCS, vol. 403, pp. 37–56
9. Ben-Sasson, E., Chiesa, A., Garman, C., Green, M., Miers, I., Tromer, E., Virza, M.: Zerocash: Decentralized anonymous payments from bitcoin. In: 2014 IEEE Symposium on Security and Privacy, pp. 459–474
10. Ben-Sasson, E., Chiesa, A., Genkin, D., Tromer, E., Virza, M.: SNARKs for C: Verifying program executions succinctly and in zero knowledge. In: CRYPTO 2013, Part II. LNCS, vol. 8043, pp. 90–108
11. Ben-Sasson, E., Chiesa, A., Tromer, E., Virza, M.: Succinct non-interactive zero knowledge for a von neumann architecture. In: USENIX Security 2014, pp. 781–796
12. Bitansky, N., Canetti, R., Chiesa, A., Tromer, E.: Recursive composition and bootstrapping for SNARKS and proof-carrying data. In: 45th ACM STOC, pp. 111–120
13. Bitansky, N., Chiesa, A., Ishai, Y., Ostrovsky, R., Paneth, O.: Succinct non-interactive arguments via linear interactive proofs. In: TCC 2013. LNCS, vol. 7785, pp. 315–333
14. Brassard, G., Chaum, D., Crépeau, C.: Minimum disclosure proofs of knowledge. Journal of computer and system sciences **37**(2) (1988) pp. 156–189
15. Bünz, B., Bootle, J., Boneh, D., Poelstra, A., Wulle, P., Maxwell, G.: Bulletproofs: Short proofs for confidential transactions and more. In: 2018 IEEE Symposium on Security and Privacy, pp. 315–334
16. Campanelli, M., Faonio, A., Fiore, D., Querol, A., Rodríguez, H.: Lunar: a toolbox for more efficient universal and updatable zkSNARKs and commit-and-prove extensions. Cryptology ePrint Archive, Report 2020/1069 (2020) <https://eprint.iacr.org/2020/1069>.
17. Canetti, R.: Universally composable security: A new paradigm for cryptographic protocols. In: 42nd FOCS, pp. 136–145
18. Chase, M., Kohlweiss, M., Lysyanskaya, A., Meiklejohn, S.: Succinct malleable NIZKs and an application to compact shuffles. In: TCC 2013. LNCS, vol. 7785, pp. 100–119
19. Chiesa, A., Hu, Y., Maller, M., Mishra, P., Vesely, N., Ward, N.P.: Marlin: Preprocessing zkSNARKs with universal and updatable SRS. In: EUROCRYPT 2020, Part I. LNCS, vol. 12105, pp. 738–768
20. Chor, B., Goldreich, O., Kushilevitz, E., Sudan, M.: Private information retrieval. In: 36th FOCS, pp. 41–50
21. Choudhuri, A.R., Jain, A., Jin, Z.: Non-interactive batch arguments for NP from standard assumptions. In: CRYPTO 2021, Part IV. LNCS, vol. 12828, pp. 394–423

22. Choudhuri, A.R., Jain, A., Jin, Z.: Snargs for  $\mathcal{P}$  from lwe. Cryptology ePrint Archive, Report 2021/808 (2021) <https://ia.cr/2021/808>.
23. Costello, C., Fournet, C., Howell, J., Kohlweiss, M., Kreuter, B., Naehrig, M., Parno, B., Zahur, S.: Geppetto: Versatile verifiable computation. In: 2015 IEEE Symposium on Security and Privacy, pp. 253–270
24. Couteau, G., Hartmann, D.: Shorter non-interactive zero-knowledge arguments and ZAPs for algebraic languages. In: CRYPTO 2020, Part III. LNCS, vol. 12172, pp. 768–798
25. Damgård, I.: Towards practical public key systems secure against chosen ciphertext attacks. In: CRYPTO'91. LNCS, vol. 576, pp. 445–456
26. Escala, A., Herold, G., Kiltz, E., Ràfols, C., Villar, J.: An algebraic framework for Diffie-Hellman assumptions. In: CRYPTO 2013, Part II. LNCS, vol. 8043, pp. 129–147
27. Fiat, A., Shamir, A.: How to prove yourself: Practical solutions to identification and signature problems. In: CRYPTO'86. LNCS, vol. 263, pp. 186–194
28. Fortnow, L.: The complexity of perfect zero-knowledge (extended abstract). In: 19th ACM STOC, pp. 204–209
29. Gabizon, A., Williamson, Z.J., Ciobotaru, O.: PLONK: Permutations over lagrange-bases for oecumenical non-interactive arguments of knowledge. Cryptology ePrint Archive, Report 2019/953 (2019) <https://eprint.iacr.org/2019/953>.
30. Ganesh, C., Kondi, Y., Orlandi, C., Pancholi, M., Takahashi, A., Tschudi, D.: Witness-succinct universally-composable snarks. Cryptology ePrint Archive (2022)
31. Gennaro, R., Gentry, C., Parno, B., Raykova, M.: Quadratic span programs and succinct NIZKs without PCPs. In: EUROCRYPT 2013. LNCS, vol. 7881, pp. 626–645
32. Gentry, C., Wichs, D.: Separating succinct non-interactive arguments from all falsifiable assumptions. In: 43rd ACM STOC, pp. 99–108
33. Goldreich, O., Håstad, J.: On the complexity of interactive proofs with bounded communication. Inf. Process. Lett. **67**(4) (1998) pp. 205–214
34. Goldreich, O., Micali, S., Wigderson, A.: Proofs that yield nothing but their validity and a methodology of cryptographic protocol design (extended abstract). In: 27th FOCS, pp. 174–187
35. Goldreich, O., Vadhan, S., Wigderson, A.: On interactive proofs with a laconic prover. Computational Complexity **11**(1) (2002) pp. 1–53
36. González, A., Ràfols, C.: New techniques for non-interactive shuffle and range arguments. In: ACNS 16. LNCS, vol. 9696, pp. 427–444
37. González, A., Zacharakis, A.: Fully-succinct publicly verifiable delegation from constant-size assumptions. Cryptology ePrint Archive, Report 2021/353 (2021) <https://eprint.iacr.org/2021/353>.
38. Groth, J.: Short pairing-based non-interactive zero-knowledge arguments. In: ASIACRYPT 2010. LNCS, vol. 6477, pp. 321–340
39. Groth, J.: On the size of pairing-based non-interactive arguments. In: EUROCRYPT 2016, Part II. LNCS, vol. 9666, pp. 305–326
40. Hada, S., Tanaka, T.: On the existence of 3-round zero-knowledge protocols. In: CRYPTO'98. LNCS, vol. 1462, pp. 408–423
41. Håstad, J., Impagliazzo, R., Levin, L.A., Luby, M.: A pseudorandom generator from any one-way function. SIAM J. Comput. **28**(4) (1999) p. 1364–1396
42. Jutla, C.S., Roy, A.: Shorter quasi-adaptive NIZK proofs for linear subspaces. In: ASIACRYPT 2013, Part I. LNCS, vol. 8269, pp. 1–20
43. Kalai, Y.T., Paneth, O., Yang, L.: How to delegate computations publicly. In: 51st ACM STOC, pp. 1115–1124
44. Kalai, Y.T., Vaikuntanathan, V., Zhang, R.Y.: Somewhere statistical soundness, post-quantum security, and SNARGs. Cryptology ePrint Archive, Report 2021/788 (2021) <https://eprint.iacr.org/2021/788>.
45. Kerber, T., Kiayias, A., Kohlweiss, M.: Composition with knowledge assumptions. In: CRYPTO 2021, Part IV. LNCS, vol. 12828, pp. 364–393
46. Kilian, J.: A note on efficient zero-knowledge proofs and arguments. In: Proceedings of the twenty-fourth annual ACM symposium on Theory of computing, pp. 723–732
47. Kiltz, E., Wee, H.: Quasi-adaptive NIZK for linear subspaces revisited. In: EUROCRYPT 2015, Part II. LNCS, vol. 9057, pp. 101–128
48. Kosba, A., Zhao, Z., Miller, A., Qian, Y., Chan, H., Papamanthou, C., Pass, R., abhi shelat, Shi, E.: C0C0: A framework for building composable zero-knowledge proofs. Cryptology ePrint Archive, Report 2015/1093 (2015) <https://ia.cr/2015/1093>.
49. Libert, B., Peters, T., Joye, M., Yung, M.: Non-malleability from malleability: Simulation-sound quasi-adaptive NIZK proofs and CCA2-secure encryption from homomorphic signatures. In: EUROCRYPT 2014. LNCS, vol. 8441, pp. 514–532

50. Lipmaa, H.: Progression-free sets and sublinear pairing-based non-interactive zero-knowledge arguments. In: TCC 2012. LNCS, vol. 7194, pp. 169–189
51. Lipmaa, H.: Succinct non-interactive zero knowledge arguments from span programs and linear error-correcting codes. In: ASIACRYPT 2013, Part I. LNCS, vol. 8269, pp. 41–60
52. Lipmaa, H., Pavlyk, K.: Gentry-wichs is tight: a falsifiable non-adaptively sound snarg. In: Advances in Cryptology – ASIACRYPT 2021, pp. 34–64
53. Micali, S.: CS proofs. In: Proceedings 35th Annual Symposium on Foundations of Computer Science, pp. 436–453
54. Naor, M.: On cryptographic assumptions and challenges (invited talk). In: CRYPTO 2003. LNCS, vol. 2729, pp. 96–109
55. Parno, B., Howell, J., Gentry, C., Raykova, M.: Pinocchio: Nearly practical verifiable computation. In: 2013 IEEE Symposium on Security and Privacy, pp. 238–252
56. Pass, R.: Unprovable security of perfect NIZK and non-interactive non-malleable commitments. In: TCC 2013. LNCS, vol. 7785, pp. 334–354
57. Ràfols, C., Zapico, A.: An algebraic framework for universal and updatable SNARKs. In: CRYPTO 2021, Part I. LNCS, vol. 12825, pp. 774–804
58. Sahai, A., Waters, B.: How to use indistinguishability obfuscation: deniable encryption, and more. In: 46th ACM STOC, pp. 475–484
59. Sako, K., Kilian, J.: Receipt-free mix-type voting scheme - a practical solution to the implementation of a voting booth. In: EUROCRYPT'95. LNCS, vol. 921, pp. 393–403
60. Waters, B., Wu, D.J.: Batch arguments for np and more from standard bilinear group assumptions. Cryptology ePrint Archive, Paper 2022/336 (2022) <https://eprint.iacr.org/2022/336>.
61. Wee, H.: On round-efficient argument systems. In: ICALP 2005. LNCS, vol. 3580, pp. 140–152
62. Wee, H., Wichs, D.: Candidate obfuscation via oblivious LWE sampling. In: EUROCRYPT 2021, Part III. LNCS, vol. 12698, pp. 127–156

## A Further Preliminaries

### A.1 Argument’s Security

We say an argument system satisfies *non-adaptive* soundness for a relation  $\mathcal{R}$  if for any PPT adversary  $\mathcal{A} = (\mathcal{A}_{\text{inp}}, \mathcal{A}_{\text{prf}})$ ,

$$\Pr \left[ \begin{array}{l} (x, \text{st}) \leftarrow \mathcal{A}_{\text{inp}}(1^\lambda) \quad \vee(\text{crs}, x, \pi) = 1 \\ (\text{crs}, \text{td}) \leftarrow \text{Setup}(1^\lambda), \pi \leftarrow \mathcal{A}_{\text{prf}}(\text{st}, \text{crs}) \quad \wedge \forall w(x, w) \notin \mathcal{R} \end{array} \right] = \text{negl}(\lambda).$$

We say that an argument system for a relation  $\mathcal{R}$  satisfies computational zero-knowledge if there exists a PPT simulator  $\text{Sim}$ , such that for any PPT  $\mathcal{A}$ ,  $|\varepsilon_0 - \varepsilon_1| = \text{negl}(\lambda)$ , where

$$\varepsilon_b := \Pr \left[ (\text{crs}, \text{td}) \leftarrow \text{Setup}(1^\lambda) : \mathcal{A}^{\text{O}_b(\text{crs}, \text{td}, \cdot)}(\text{crs}) = 1 \right].$$

The oracle  $\text{O}_b$  takes an input  $(\text{crs}, \text{td}, x, w)$ . It returns  $\perp$ , when  $(x, w) \notin \mathcal{R}$ . Otherwise, if  $b = 0$  it returns  $\pi \leftarrow \text{P}(\text{crs}, x, w)$  and if  $b = 1$  it returns  $\pi \leftarrow \text{Sim}(\text{crs}, \text{td}, x)$ .

### A.2 Fully Homomorphic Encryption (FHE)

A FHE scheme consists of a tuple of algorithms  $(\text{KG}, \text{Enc}, \text{Dec}, \text{Eval})$  with the following syntax:

$\text{KG}(1^\lambda) \rightarrow (\text{pk}, \text{sk})$ : generates a key pair (the algorithm is randomized).

$\text{Enc}(\text{pk}, m) \rightarrow \text{ct}$ : produces a ciphertext corresponding to a message  $m$  through the public key (the algorithm is randomized).

$\text{Dec}(\text{sk}, \text{ct}) \rightarrow m$ : decrypts a ciphertext through the secret key (the algorithm is deterministic).

$\text{Eval}(\text{pk}, \text{ct}_m, F) \rightarrow \text{ct}_F$ : produces an encryption of  $F(m)$  from an encryption of  $m$  through the public key. Occasionally we will overload this notation for functions with arity higher than 1 or that take as input plaintexts, which can be seen as dummy ciphertexts (the algorithm is deterministic).

Occasionally, we need to explicitly write  $\text{Enc}(\text{pk}, m; r)$ , where  $r \leftarrow_{\$} \text{Rnd}$  are the random coins sampled from the randomness space  $\text{Rnd}$ .

An FHE scheme should satisfy correctness and semantic security.

**Correctness.** For any  $\lambda$ , plaintext  $m$  and function  $F$

$$\Pr [\text{Dec}(\text{sk}, \text{ct}) = m] = 1$$

and

$$\Pr [\text{Dec}(\text{sk}, \text{Eval}(\text{pk}, \text{ct}, F)) = F(m)] = 1$$

where  $(\text{pk}, \text{sk}) \leftarrow \text{KG}(1^\lambda)$  and  $\text{ct} \leftarrow \text{Enc}(\text{pk}, m)$ .

**Semantic security.** For all  $\lambda$ , for any PPT adversary  $\mathcal{A} = (\mathcal{A}^1, \mathcal{A}^2)$

$$\Pr \left[ \begin{array}{l} (\text{pk}, \text{sk}) \leftarrow \text{KG}(1^\lambda), (\text{st}, m_0, m_1) \leftarrow \mathcal{A}^1(\text{pk}) \\ b \leftarrow_{\$} \{0, 1\}, \text{ct} \leftarrow \text{Enc}(m_b), b' \leftarrow \mathcal{A}^2(\text{st}, \text{ct}) \end{array} : b = b' \right] = \text{negl}(\lambda)$$

### A.3 Universal Hash Functions

We say that a family of functions  $\mathcal{H} = \{\text{H}_{\text{hk}} : M \rightarrow Y\}_{\text{hk} \in K}$  for  $|M| > |Y|$  is a universal hash function (UHF) family if for all distinct  $m_1, m_2 \in M$ ,

$$\Pr [\text{H}_{\text{hk}}(m_1) = \text{H}_{\text{hk}}(m_2) : \text{hk} \leftarrow_{\$} K] = 1/|Y| .$$

Therefore, UHF must only satisfy non-adaptive collision resistance compared to a standard cryptographic hash function. More precisely, in a cryptographic hash function, the adversary first gets a key  $\text{hk}$  and then has to find distinct  $m_1, m_2$  that produce a collision. In the case of UHF,  $m_1, m_2$  are fixed beforehand, and then  $\text{hk}$  is picked independently of  $m_1$  and  $m_2$ . Although satisfying a weaker property, UHFs are very efficient to compute and information-theoretically secure.

Consider, for example, the following construction. Let  $\mathbb{F}$  be a prime order finite field. For some  $n \geq 0$ , we define  $M = K = \mathbb{F}^n$  and  $Y = \mathbb{F}$ . For  $\text{hk} = (k_1, \dots, k_n) \leftarrow_{\$} \mathbb{F}^n$  the hash function is defined as  $\text{H}_{\text{hk}}(m_1, \dots, m_n) = \sum_{i=1}^n k_i m_i$ . Let us confirm that this is indeed a UHF. Suppose for distinct  $\vec{a}, \vec{b} \in \mathbb{F}^n$ ,  $\text{H}_{\text{hk}}(\vec{a}) = \text{H}_{\text{hk}}(\vec{b})$ . Thus,  $\text{H}_{\text{hk}}(\vec{a}) - \text{H}_{\text{hk}}(\vec{b}) = \sum_{i=1}^n k_i (a_i - b_i) = 0$ . Since at least for some  $j$ ,  $a_j - b_j \neq 0$ , there are  $|\mathbb{F}|^{n-1}$  keys that satisfy this equation. Thus, the probability of collision is  $|\mathbb{F}|^{n-1}/|\mathbb{F}|^n = 1/|\mathbb{F}|$ . This probability will be negligible by picking a suitably large field  $\mathbb{F}$ .

One can also compress the key size in the previous construction by defining the key space to be  $K = \mathbb{F}$  and  $\text{H}_{\text{hk}}(\vec{m}) = \sum_{i=1}^n m_i k^{i-1}$ , where  $k \leftarrow_{\$} K = \mathbb{F}$ . Since,  $\text{H}_{\text{hk}}(\vec{a}) - \text{H}_{\text{hk}}(\vec{b}) = \sum_{i=1}^n (a_i - b_i) k^{i-1} = 0$  if  $k$  is a root of a non-zero at most degree  $n-1$  polynomial, it follows that the collision probability is bounded by  $(n-1)/|\mathbb{F}|$ . This is not formally a UHF anymore (collision probability is not  $1/|Y|$ ), but still serves our purposes.

### A.4 Construction of CLR-OWF

Agrawal et al. [1] propose and prove the security of the following CLR-OWF.  $\text{KGen}(1^\lambda)$  picks a discrete logarithm secure group  $\mathbb{G}$  of order  $p$  with a generator  $g$ . It samples  $\vec{\alpha} = (\alpha_1, \dots, \alpha_n) \leftarrow_{\$} \mathbb{Z}_p^n$  and sets  $g_i \leftarrow g^{\alpha_i}$  for  $i = 1, \dots, n$ . The public parameter is  $\text{pp} = (\mathbb{G}, g, g_1, \dots, g_n)$  and the update key is  $\text{uk} = \vec{\alpha}$ . The sampling algorithm  $\text{Sample}(\text{pp})$  outputs  $\vec{x} \leftarrow_{\$} \mathbb{Z}_p^n$ .  $\text{Eval}(\text{pp}, \vec{x})$  returns  $y \leftarrow \prod_{i=1}^n g_i^{x_i}$ .  $\text{Update}(\text{uk}, \vec{x})$  chooses a random vector  $\vec{\beta}$  that is orthogonal to  $\vec{\alpha}$  and returns  $\vec{x}' \leftarrow \vec{x} + \vec{\beta}$ .

The correctness holds since  $\prod_{i=1}^n g_i^{x'_i} = g^{\sum_{i=1}^n \alpha_i x_i + \sum_{i=1}^n \alpha_i \beta_i} = g^{\sum_{i=1}^n \alpha_i x_i} = \prod_{i=1}^n g_i^{x_i}$ .

**Theorem 7 ([1]).** *If the discrete logarithm assumption holds in group  $\mathbb{G}$ , then there exists a L-CLR-OWF in the floppy model, with  $L(\lambda) < (n-2) \log p - \omega(\log \lambda)$ .*

## B Impossibility of Adaptive BB Extraction Using LR-OWF

The impossibility shown in theorem 1 can be interpreted as a consequence of leakage-resilience; a SNARK proof is leakage on the witness. For an **NP**-relation that is leakage-resilient, recovering the entire witness is impossible for an extractor even given the leakage, if this leakage is small.

**Definition 8** ( $(\ell, \varepsilon)$ -LR-OWF). A function family  $\mathcal{F} = \{f_i : D_i \rightarrow R_i\}$  is  $(\ell, \varepsilon)$ -LR one-way if:

- There exists efficient algorithms (i)  $\text{KGen}(1^\lambda)$  to sample an index  $i$  (ii)  $\text{Sample}(i)$  for sampling an input  $x \leftarrow_{\$} D_i$  (iii)  $\text{Eval}(i, x)$  for computing  $y = f_i(x)$ .
- For any PPT  $\mathcal{A}$ ,

$$\Pr \left[ \begin{array}{l} i \leftarrow \text{KGen}(1^\lambda), x \leftarrow \text{Sample}(i), \\ y \leftarrow \text{Eval}(i, x), x' \leftarrow \mathcal{A}^{\text{O}_\ell(\cdot)}(i, y) : y = \text{Eval}(i, x') \end{array} \right] \leq \varepsilon$$

where  $\text{O}_\ell(\cdot)$  is an oracle that takes as an input a leakage function  $h : \{0, 1\}^* \rightarrow \{0, 1\}^\ell$ , on which  $\text{O}_\ell(h)$  returns  $h(x)$ . Adversary can query  $\text{O}_\ell(\cdot)$  only once.

Let  $\mathcal{F}$  be a family of  $(\ell, \varepsilon)$ -LR OWFs. For  $f \in \mathcal{F}$ , consider the relation  $\mathcal{R}_f := \{((x, i), w) \mid i \in \text{KGen}(1^\lambda), w \in \text{Sample}(i), x = \text{Eval}(i, w)\}$ .

**Theorem 8.** A non-interactive argument system  $\Pi$  for  $\mathcal{R}_f$  with argument size at most  $\ell$  bits, has black-box knowledge soundness error  $\varepsilon_{ks} \geq 1 - \varepsilon$ .

*Proof.* By LR one-wayness of  $f$ , we have

$$\Pr[i \leftarrow \text{KGen}(1^\lambda), x \leftarrow_{\$} D_i, w \leftarrow \mathcal{A}^{\text{O}_\ell(\cdot)}(1^\lambda, x) : ((x, i), w) \in \mathcal{R}_f] \leq \varepsilon.$$

Consider an argument system  $\Pi$  for  $\mathcal{R}_f$  with argument size bounded by  $\ell$  bits. Let  $\text{Ext}$  be the black-box extractor guaranteed by  $\Pi$ . We construct an adversary  $\mathcal{A}^{\text{O}_\ell(\cdot)}$  that breaks the  $\ell$ -leakage resilience of  $f$ .  $\mathcal{A}$  receives as challenge  $x$ , picks a  $\text{crs}$  together with an extraction key  $\text{td}$ , sets  $h(X) := \text{P}(\text{crs}, x, X)$ , and receives  $\pi \leftarrow \text{O}_\ell(h) = \text{P}(\text{crs}, x, w)$ . It then invokes the black-box witness extractor  $\text{Ext}(\text{crs}, \text{td}, x, \pi)$  to receive  $w$ , and returns  $w$  as preimage. Assuming perfect correctness, we have that  $\mathcal{A}$  succeeds in breaking one-wayness of  $f$  with the probability that the extractor succeeds.  $\Pr[\mathcal{A} \text{ succeeds}] \geq 1 - \varepsilon_{ks}(\lambda)$ . Thus,  $\varepsilon_{ks} \geq 1 - \varepsilon$ .  $\square$

Since an adversary can always guess the correct leakage with probability  $1/2^\ell$ , all OWFs are LR with  $\ell = O(\log |w|)$ . We therefore obtain as a corollary, that an argument for  $\mathcal{R}_f$  must be at least of logarithmic size.

**Corollary 1.** Assuming the existence of OWFs, a SNARK with negligible black-box knowledge soundness error must have argument size at least  $\Omega(\log |w|)$ .

With concrete LR-OWFs even better lower bounds can be achieved. Consider for example the discrete logarithm based (C)LR-OWF from theorem 7. We obtain the following result.

**Corollary 2.** If the discrete logarithm assumption holds, then there exist a LR-OWF family  $\mathcal{F}$  such that any non-interactive black-box knowledge sound argument for the relation  $\mathcal{R}_f$  must have size  $\Omega(|w|)$ .

The latter also shows that black-box extractable SNARKs for all **NP** do not exist.

## C Additional Proofs

### C.1 Proof of Theorem 2

**Theorem 9 (Restatement of theorem 2).** If  $\Pi_{\exists}$  is a non-adaptively sound SNARG scheme for **NP**, FHE is a semantically secure FHE scheme, PKE is a semantically secure cryptosystem and  $\text{H}$  is a family of UHFs, then the construction in fig. 2 is a SNARK for **FewP** satisfying definition 4. If  $\Pi_{\exists}$  has additionally computational zero-knowledge, then so does the resulting SNARK.

*Proof.* We prove zero-knowledge in lemma 10. In the rest, we will focus on proving knowledge soundness. We use the extractor in fig. 3. The extractor can have embedded  $N_w$ , a bound on the length of witnesses, since it is non-uniform. In the remainder, we define  $S(j)$ , for index  $j$ , as the set of strings  $h$  for which  $W[h][j] \neq \perp$  after running  $\text{Qldx}(x, j)$ . That is,

$$S(j) := \{h : W[h][j] \neq \perp \text{ after running } \text{Qldx}(x, j)\} .$$

Later in lemma 4 we show that for  $h \in S(j)$  there exists  $w$  such that  $R(x, w) \wedge H(w) = h$  (with high probability). Therefore,  $h \in S(j)$  intuitively means “the extractor holds the  $j$ -th bit of a witness  $w$  such that  $H(w) = h$ ”.

In order to argue black-box knowledge soundness, we should be able to successfully extract from an adversary that returns an accepting proof with noticeable probability, i.e. with probability  $\lambda^{-c}$  for some positive constant  $c^{17}$ . We show that for this type of adversary it holds with noticeable probability that  $\exists h \in \bigcap_{j=1}^{N_w} S(j)$  (this is key for extraction; see last line in extractor definition). We argue this is the case by combining two facts:

- $S(j) = S(j')$  with overwhelming probability for all  $j, j'$  (lemma 6);
- If  $\Pr[\text{adversary returns an accepting proof}]$  is non-negligible then  $\Pr[S(j) \neq \emptyset]$  is non-negligible (lemma 7).

If  $\exists h \in \bigcap_{j=1}^{N_w} S(j)$ , then the string returned by the extractor is a witness with overwhelming probability because  $W[h][j]$  is a bit of a witness for the relation with overwhelming probability (by lemma 4) and because, except with negligible probability, there exists a unique  $w$  such that  $H_{\text{hk}}(w) = h$  (by lemma 8). This concludes the proof.  $\square$

The following auxiliary lemma shows that an element in the table constructed by the extractor actually captures a bit of the witness with high probability.

**Lemma 4.** *For any PPT adversary  $\mathcal{A}$ , for each  $j \in \{1, \dots, N_w\}$ , for each  $h \in S(j)$  (where  $S$  is defined in the proof of theorem 9) the following probability  $p$  is overwhelming:*

$$p := \Pr[\exists w : R(x, w) \wedge H_{\text{hk}}(w) = h \wedge W[h][j] = w_j] .$$

*Proof.* Consider  $h \in S(j)$ . Notice that the event above is implied by the event  $R(x, w) \wedge \text{ct}_h = \text{PKE.Enc}(\text{pk}, H_{\text{hk}}(w); r) \wedge \text{ct}_{\text{bit}} = \text{FHE.Eval}(\text{pk}_{\text{FHE}}, f_{\text{proj}}, \text{ct}_{i^*}, w)$ . This is because the extractor in fig. 3 decrypts  $\text{ct}_h$  to recover  $h$  and sets  $W[h][j]$  to the decryption of  $\text{ct}_{\text{bit}}$ . We can argue that the probability of such event is overwhelming because by definition of the extractor if  $h \in S(j)$  then the adversary provided a corresponding SNARG proof for the relation  $\mathcal{R}'$ . Invoking correctness of the encryption schemes and soundness of the SNARG concludes the proof.  $\square$

The following auxiliary lemma observes that the probability of an adversary returning a valid proof with good probability for a CRS containing a randomly sampled index  $i^*$  should also hold when we provide them with a CRS “referring to” an arbitrary index  $i^*$ . This is useful to ensure that we can apply our extraction strategy. Otherwise we could for example conceive an adversary returning a valid proof for all indices except a few. Such an adversary would return a valid proof with high probability for a honestly generated CRS but we would not be able to extract from it.

**Lemma 5.** *For any PPT adversary  $\mathcal{A}$ , if  $\Pr[\mathcal{A} \text{ returns an accepting proof}]$  in the black-box knowledge-soundness experiment (definition 4) is non-negligible then for any  $i^* \in [N_w]$  the following probability is non-negligible:*

$$p_{\text{acc}}^{(i^*)} := \Pr[(\text{crs}, \text{td}) \leftarrow \overline{\text{Setup}}_{i^*}(1^\lambda), (x, \pi) \leftarrow \mathcal{A}(\text{crs}) : V(\text{crs}, x, \pi) = 1]$$

where  $\overline{\text{Setup}}_{i^*}$  is defined in fig. 6.

<sup>17</sup> This simplifies the proof, but we can argue with minor modifications the case for an adversary returning an accepting proof with only non-negligible probability.



$\overline{\text{Setup}}_{i^*}(1^\lambda)$ <hr style="border: 0.5px solid black;"/> $(\text{crs}, \hat{\text{td}}) \leftarrow \Pi_{\exists}.\text{Setup}(1^\lambda)$ $(\text{pk}_{\text{FHE}}, \text{sk}_{\text{FHE}}) \leftarrow \text{FHE.KG}(1^\lambda)$ $(\text{pk}, \text{sk}) \leftarrow \text{PKE.KG}(1^\lambda)$ $\text{ct}_{i^*} \leftarrow \text{FHE.Enc}(\text{pk}, i^*)$ $\text{hk} \leftarrow \$_\mathcal{K}_{\text{UHF}}$ <b>return</b> $(\text{crs} := (\text{crs}, \text{ct}_{i^*}, \text{hk}, \text{pk}_{\text{FHE}}, \text{pk}), \text{td} := (\text{sk}_{\text{FHE}}, \text{sk}, \hat{\text{td}}))$
---

Fig. 6: Modified setup with fixed index in lemma 5.

*Proof.* First observe that for any adversary  $\mathcal{A}$ , for any  $j, j' \in [N_w]$ . The probabilities  $p_{\text{acc}}^{(j)}$  and  $p_{\text{acc}}^{(j')}$  must be negligibly close. If they were not then we could build an adversary breaking IND-CPA of the FHE since intuitively we could distinguish ciphertexts of  $j$  from those of  $j'$  (a formal description of this adversary would be a simpler variant of the one we build in the proof of lemma 6).

Next, we observe that we can write  $p_{\text{acc}}^{(\text{avg})} \Pr[\mathcal{A} \text{ returns an accepting proof}]$  in the black-box knowledge-soundness experiment (definition 4) as a function of  $p_{\text{acc}}^{(i^*)}$  for  $i^* = 1, \dots, N_w$  through a simple marginalization and bound it as follows

$$p_{\text{acc}}^{(\text{avg})} = \frac{1}{|N_w|} \sum_{i^* \in [N_w]} p_{\text{acc}}^{(i^*)} \leq \min_{i^* \in [N_w]} p_{\text{acc}}^{(i^*)} + \epsilon ,$$

where  $\epsilon$  is a negligible. We can argue the bound by simple algebra and by applying our previous observation. As a consequence of the above, it is easy to see that, if  $p_{\text{acc}}^{(\text{avg})}$  is non-negligible, so must be each  $p_{\text{acc}}^{(i^*)}$ .  $\square$

**Lemma 6.** *For any PPT adversary  $\mathcal{A}_{\text{ksnd}} = (\mathcal{A}_{\text{inp}}, \mathcal{A}_{\text{prf}})$ , for all  $j \neq j'$  the sets  $S(j), S(j')$  are equal except with negligible probability (where  $S$  is defined in the proof of theorem 9).*

*Proof.* Assume by contradiction that it is not the case. We show we can break semantic security of FHE (appendix A.2) with the adversary  $\mathcal{A}_{\text{CPA}}$  in fig. 7.

Intuitively the adversary  $\mathcal{A}_{\text{CPA}}$  does the following. After receiving a public key  $\text{pk}_{\text{FHE}}$  from the FHE challenger, it uses it to “emulate” the extractor invoking a variant of  $\text{QIdx}$  in fig. 3 ( $\overline{\text{QIdx}}$  in fig. 7). That is,  $\mathcal{A}_{\text{CPA}}$  constructs set  $S(j)$  exactly as the extractor (implicitly) does, but without storing the decrypted bits in  $W[h][j]$  (which it cannot due to not having the secret key). More precisely, after receiving a valid proof for index  $j$  which includes hash  $h^*$  (after decrypting  $\text{ct}_h$  with  $\text{sk}$ ),  $\mathcal{A}_{\text{CPA}}$  will simply set  $W[h^*][j]$  to a dummy “check” value ( $\checkmark$  in fig. 7). This is enough to define the sets  $S(j)$  which are our concern below. By hypothesis there exist with non-negligible probability indices  $j_0, j_1$  such that  $S(j_0) \neq S(j_1)$ . Adversary  $\mathcal{A}_{\text{CPA}}^1$  finds such indices and returns  $j_0$  and  $j_1$  to the FHE challenger as challenge plaintexts. Once received a ciphertext  $\text{ct}_?$   $\mathcal{A}_{\text{CPA}}^2$  will query polynomially many times  $\mathcal{A}_{\text{prf}}$  with a CRS that uses  $\text{ct}_?$  as encrypted index. Call the set of response hash ciphertexts from these queries  $S_?$ .

By setting  $N_q$ —the number of queries—to an appropriately high value in  $\overline{\text{QIdx}}$  and  $\overline{\text{QIdx}}'$  we can claim the following fact. By invoking lemma 9<sup>18</sup>, except with non negligible probability, the set  $S_?$  is equal to either  $S(j_0)$  or  $S(j_1)$ .  $\mathcal{A}_{\text{CPA}}$  compares  $S_?$  to them and outputs the bit corresponding to which one it is equal to. From this we can conclude that the advantage of  $\mathcal{A}_{\text{CPA}}$  in breaking semantic security is negligibly close to  $\Pr[\exists j, j' : S(j) \neq S(j')]$ . If the latter is non-negligible so is the advantage of  $\mathcal{A}_{\text{CPA}}$ . Absurd.  $\square$

**Lemma 7.** *For any PPT adversary  $\mathcal{A}$ , if  $\Pr[\mathcal{A} \text{ returns an accepting proof}]$  in the black-box knowledge-soundness experiment (definition 4) is non-negligible then for any  $j \in [N_w]$  the probability  $\Pr[S(j) \neq \emptyset]$  is non-negligible.*

<sup>18</sup> Which guarantees that  $N_q = \text{poly}(\lambda)$  is sufficient.

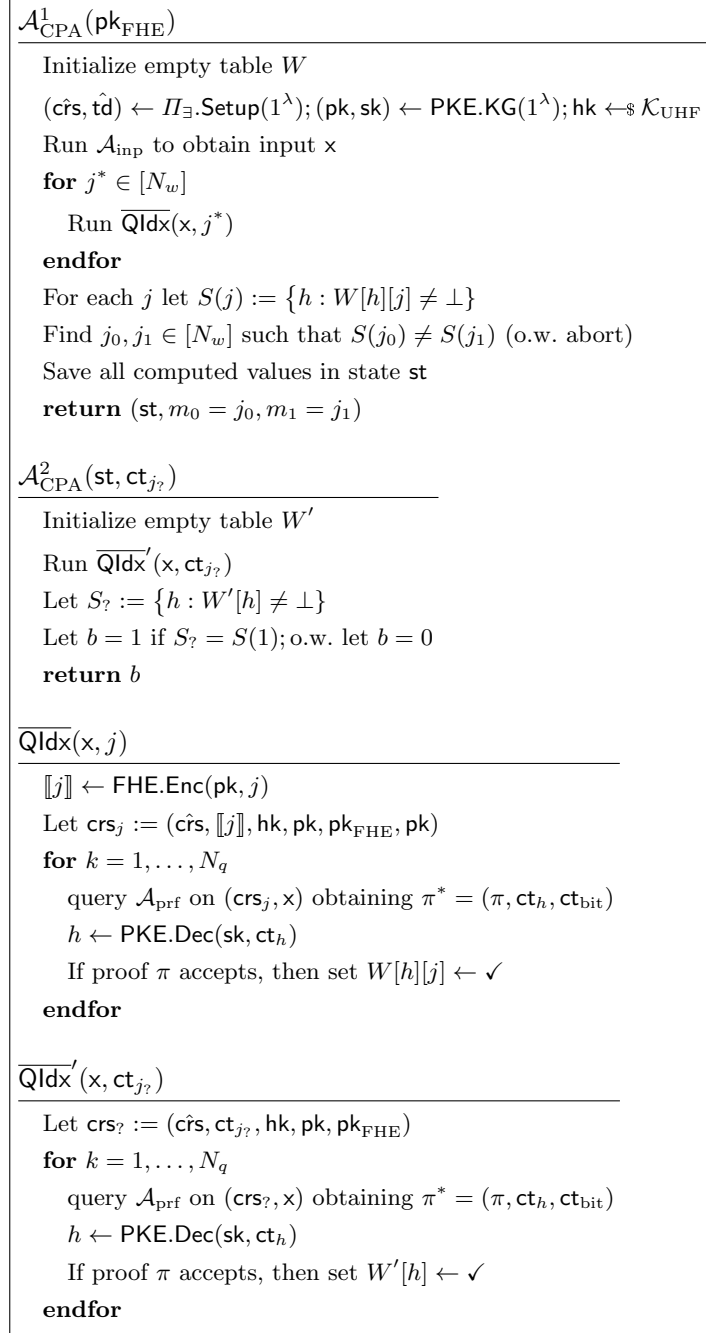


Fig. 7: IND-CPA adversary and auxiliary algorithms for proof in lemma 6.

*Proof.* This follows easily by the definition of set  $S$  and how  $\overline{\text{Qldx}}$  works (fig. 3). In fact, by inspecting the last two lines in the loop of  $\overline{\text{Qldx}}$  and the definition of  $S$  (see proof of theorem 9) we can see that  $S(j)$  becomes a non-empty set as long as the adversary returns at least one accepting proof referring to index  $j$ . This can be bounded by the probability  $p_{\text{acc}}^{(j)}$  as defined in lemma 5. Applying that same lemma concludes the proof.  $\square$

The following fact is useful in the proofs of the lemmas above. It states that we should not expect hash collisions among witnesses.

**Lemma 8.** *For any PPT adversary, for all witnesses  $w, w'$  such that  $R(x, w)$  and  $R(x, w')$  it holds that  $\Pr [H_{\text{hk}}(w) = H_{\text{hk}}(w')] is negligible, where the probability is over the randomness of the adversary and the sampling of  $\text{hk}$ .$*

*Proof.* The statement follows directly from the likely absence of collisions in the hash function (appendix A.3) and from the fact that the instance  $x$  is selected independently of the hash key (this is implied by non-adaptive security, i.e. definition 4).

The following lemma essentially states that the set  $S(j)$  “converges” after sufficiently many queries  $N_q = \text{poly}(\lambda)$ .

**Lemma 9.** *For any PPT adversary, for each index  $j$ , there exists constant  $c$  such that for all constants  $c' > c$  the sets  $S^{(\lambda^{c'})}(j) = S^{(\lambda^c)}(j)$  except with negligible probability, where  $S^{(t)}(j)$  denotes the set  $S(j)$  after  $t$  queries to  $\text{QIdx}(x, j)$ .*

*Proof.* For an adversary returning a valid proof with negligible probability, the result follows immediately. Let us then consider the case of an adversary returning a valid proof with non-negligible probability. We proceed by contradiction: assume that for any polynomial number of invocations  $t$  to  $\text{QIdx}(x, j)$ , the size of the set  $S(j)$  increases with non-negligible probability after a polynomial number of steps. Call  $N$  the number of witnesses of  $x$  and recall that  $N = \text{poly}(|x|) = \text{poly}(\lambda)$  since the language is in **FewP**. Then observe that there exists polynomial number of steps  $t^*$  after which the set  $S$  became larger than  $N$ , i.e.,  $|S^{(t^*)}(j)| > N$  with non-negligible probability. By our hypothesis, after a polynomial number of steps, at least one hash image will be added to  $S(j)$  with non-negligible probability. By soundness of the underlying SNARG, for each of these hashes, there is a preimage that is also a witness with non-negligible probability, implying that the number of witnesses is greater than  $N$ . Absurd.  $\square$

**Lemma 10.** *If  $\Pi_{\exists}$  has computational zero-knowledge and FHE and PKE are semantically secure cryptosystems, then the construction in fig. 2 has computational zero-knowledge.*

*Proof.* Recall the computational zero-knowledge definition in appendix A.1. Let  $\mathcal{A}$  be a PPT adversary, and suppose the number of queries that  $\mathcal{A}$  can make is bounded by  $q(\lambda)$ . Since  $\mathcal{A}$  runs in polynomial time, the function  $q(X)$  is bounded by some polynomial. We define the simulator  $\text{Sim}$  in fig. 8. In short,  $\text{Sim}$  encrypts arbitrary values (e.g., zeros) inside  $\text{ct}_h$  and  $\text{ct}_{\text{bit}}$  and then simulates  $\Pi_{\exists}$  proof  $\pi$ .

Let  $\text{Game}_0$  be the game where the zero-knowledge oracle  $\mathcal{O}_0$  responds with honestly constructed proofs. In  $\text{Game}_1$  we change the oracle to  $\mathcal{O}_1$ , which simulates the proofs  $\pi$  of  $\Pi_{\exists}$ .

Next, we consider a sequence of games  $\text{Game}_{2;i}$  for  $i = 1, \dots, q(\lambda)$ . On the first  $i$  queries, in  $\text{Game}_{2;i}$ , we run  $\mathcal{O}_2$  which works just as  $\mathcal{O}_1$ , but encrypts 0 in  $\text{ct}_h$  instead of the hash of the witness. Subsequent queries will still use  $\mathcal{O}_1$ . Similarly, we will have games  $\text{Game}_{3;i}$ , which on the first  $i$  queries will run  $\mathcal{O}_3$  that additionally encrypts 0 in  $\text{ct}_{\text{bit}}$  and on the subsequent queries will run  $\mathcal{O}_2$ . Oracles  $\mathcal{O}_1, \mathcal{O}_2, \mathcal{O}_3$  are described in full detail in fig. 8. Observe that  $\text{Game}_{3;q(\lambda)}$  is equivalent to the zero-knowledge game, where  $\mathcal{A}$  receives a simulated proof (that is, proofs are generated by  $\text{Sim}$ ) for all queries.

Let us now analyze the success probability of  $\mathcal{A}$  in the previous games. Let  $\varepsilon_i$  denote  $\mathcal{A}$ 's success probability in  $\text{Game}_i$ . We need to show that  $|\varepsilon_0 - \varepsilon_{3;q(\lambda)}| = \text{negl}(\lambda)$ .

$\text{Game}_0 \rightarrow \text{Game}_1$ . We construct a reduction  $\mathcal{B}_{zk}$  for computational zero-knowledge of  $\Pi_{\exists}$ . The reduction is straightforward.  $\mathcal{B}_{zk}$  has access to an oracle  $\mathcal{O}_{b, \Pi_{\exists}}$  which returns an honest  $\Pi_{\exists}$  proof when  $b = 0$  and a simulated proof when  $b = 1$ .  $\mathcal{B}_{zk}$  runs  $\mathcal{A}$  internally and creates proofs for  $\mathcal{A}$ 's queries. However, when it comes to creating  $\pi$ ,  $\mathcal{B}_{zk}$  uses the oracle  $\mathcal{O}_{b, \Pi_{\exists}}$  to create it.  $\mathcal{B}_{zk}$  outputs the same value as  $\mathcal{A}$ . Detailed construction of  $\mathcal{B}_{zk}$  can be seen in fig. 9. It follows that  $|\varepsilon_0 - \varepsilon_1| = \text{negl}(\lambda)$  since  $\mathcal{B}_{zk}$ 's advantage is negligible in the zero-knowledge game for  $\Pi_{\exists}$ .

$\text{Game}_{2;i-1} \rightarrow \text{Game}_{2;i}$ . We construct a reduction  $\mathcal{B}_{sem}$  to the semantic security of PKE. Full details are described in fig. 10. Essentially,  $\mathcal{B}_{sem}^1$  runs  $\mathcal{A}$  for  $i$  queries and outputs  $m_0 = H_{\text{hk}}(w^{(i)})$  and  $m_1 = 0$ , where

$O_{\text{lbl}}(\text{crs}, \text{td}, x, w)$	$\text{Sim}(\text{crs}, \text{td}, x)$
$r \leftarrow \text{PKE.Rnd}$	$r \leftarrow \text{PKE.Rnd}$
$h \leftarrow \text{H}_{\text{hk}}(w); \text{ct}_h \leftarrow \text{PKE.Enc}(\text{pk}, h; r)$	$\text{ct}_h \leftarrow \text{PKE.Enc}(\text{pk}, 0; r)$
$\text{ct}_h \leftarrow \text{PKE.Enc}(\text{pk}, 0; r)$	$\text{ct}_{\text{bit}} \leftarrow \text{FHE.Enc}(\text{pk}_{\text{FHE}}, 0)$
$\text{ct}_{\text{bit}} \leftarrow \text{FHE.Eval}(\text{pk}_{\text{FHE}}, f_{\text{proj}}, \text{ct}_{i^*}, w)$	$\hat{x} \leftarrow (x, \text{hk}, \text{pk}_{\text{FHE}}, \text{pk}, \text{ct}_h, \text{ct}_{i^*}, \text{ct}_{\text{bit}})$
$\text{ct}_{\text{bit}} \leftarrow \text{FHE.Enc}(\text{pk}_{\text{FHE}}, 0)$	$\pi \leftarrow \Pi_{\exists}.\text{Sim}(\text{crs}, \text{td}, \hat{x})$
$\hat{x} \leftarrow (x, \text{hk}, \text{pk}_{\text{FHE}}, \text{pk}, \text{ct}_h, \text{ct}_{i^*}, \text{ct}_{\text{bit}})$	<b>return</b> $\pi^* := (\pi, \text{ct}_h, \text{ct}_{\text{bit}})$
$\pi \leftarrow \Pi_{\exists}.\text{Sim}(\text{crs}, \hat{\text{td}}, \hat{x})$	
<b>return</b> $\pi^* := (\pi, \text{ct}_h, \text{ct}_{\text{bit}})$	

Fig. 8: Oracles  $O_1, O_2, O_3$  and the simulator  $\text{Sim}$  for lemma 10. When  $\text{lbl} = 1, 2, 3$ , then we respectively consider only lines with  $\boxed{AAA}$ ,  $\boxed{AAA}$ , and  $\boxed{AAA}$ . Recall  $f_{\text{proj}}(i, w) := w_i$

$\mathcal{B}_{zk}^{O_{b, \Pi_{\exists}}(\text{crs}, \hat{\text{td}}, \cdot)}(\text{crs})$	$O_{\mathcal{B}_{zk}}^{O_{b, \Pi_{\exists}}(\text{crs}, \hat{\text{td}}, \cdot)}(\text{crs}, x, w)$
$(\text{pk}_{\text{FHE}}, \text{sk}_{\text{FHE}}) \leftarrow \text{FHE.KG}(1^\lambda)$	<b>if</b> $(x, w) \notin \mathcal{R}$ <b>then return</b> $\perp$
$\text{hk} \leftarrow \mathcal{K}_{\text{UHF}}$	$h \leftarrow \text{H}_{\text{hk}}(w)$
$(\text{pk}, \text{sk}) \leftarrow \text{PKE.KG}(1^\lambda)$	$r \leftarrow \text{PKE.Rnd}; \text{ct}_h \leftarrow \text{PKE.Enc}(\text{pk}, h; r)$
$i^* \leftarrow [N_w]; \text{ct}_{i^*} \leftarrow \text{FHE.Enc}(\text{pk}_{\text{FHE}}, i^*)$	$\text{ct}_{\text{bit}} \leftarrow \text{FHE.Eval}(\text{pk}_{\text{FHE}}, f_{\text{proj}}, \text{ct}_{i^*}, w)$
$\text{crs} \leftarrow (\text{crs}, \text{ct}_{i^*}, \text{hk}, \text{pk}_{\text{FHE}}, \text{pk})$	$\hat{x} \leftarrow (x, \text{hk}, \text{pk}_{\text{FHE}}, \text{pk}, \text{ct}_h, \text{ct}_{i^*}, \text{ct}_{\text{bit}})$
<b>return</b> $\mathcal{A}^{O_{\mathcal{B}_{zk}}(\text{crs}, \cdot)}(\text{crs})$	$\hat{w} \leftarrow (w, r)$
	Query $\hat{\pi} \leftarrow O_{b, \Pi_{\exists}}(\text{crs}, \hat{\text{td}}, \hat{x}, \hat{w})$
	<b>return</b> $\pi \leftarrow (\hat{\pi}, \text{ct}_h, \text{ct}_{\text{bit}})$

Fig. 9:  $\text{Game}_0 \rightarrow \text{Game}_1$  reduction in lemma 10

$w^{(i)}$  is the witness  $\mathcal{A}$  submits on the  $i$ -th query. Then,  $\mathcal{B}_{\text{sem}}^2$  gets an encryption  $\text{ct}$ , which encrypts either  $m_0$  or  $m_1$ .  $\mathcal{B}_{\text{sem}}^2$  runs internally  $\mathcal{A}$ . For queries  $j = 1, \dots, i-1$  adversary  $\mathcal{B}_{\text{sem}}^2$  returns the answers of  $O_1$  and for queries  $j > i$  it returns the answer of  $O_2$  to  $\mathcal{A}$ . However, on the  $i$ -th query it behaves as  $O_1$  but  $\text{ct}_h = \text{ct}$ . So if  $\text{ct}_h$  encrypts  $\text{H}_{\text{hk}}(w^{(i)})$  it will be oracle  $O_1$  and if  $\text{ct}_h$  encrypts 0 it will be the oracle  $O_2$ .  $\mathcal{B}_{\text{sem}}^2$  returns the same answer as  $\mathcal{A}$ .

It follows from the semantic security that  $|\varepsilon_{2:i-1} - \varepsilon_{2:i}| = \text{negl}(\lambda)$ .

$\text{Game}_{3:i-1} \rightarrow \text{Game}_{3:i}$ . This is identical to the previous case, with the only difference being that we rely on the semantic security of FHE.

Since  $q(\lambda)$  is bounded by a polynomial and at every game transition  $\mathcal{A}$ 's advantage increases at most by a negligible amount, then it follows that our construction has computational zero-knowledge.  $\square$

## C.2 Proof of Theorem 3

**Theorem 10 (Restatement of theorem 3).** *Let  $\Sigma = (\text{KGen}, \text{Sample}, \text{Eval}, \text{Update})$  be an  $L$ -CLR-OWF and let  $\Pi$  be a non-adaptive black-box  $\varepsilon_{ks}(\lambda)$ -knowledge sound argument for  $\mathcal{R}_{\Sigma}$  as defined above. If the proof size is less than  $L(\lambda)$  bits, then  $L$ -CLR-OWF can be broken with probability  $1 - \varepsilon_{ks}(\lambda)$ .*

*Proof.* Assume that the proof size of  $\Pi$  is upper bounded by  $L(\lambda)$ . We show that in this case it is possible to break  $L$ -continuous leakage-resilient one-wayness of  $\Sigma$ . Firstly, we describe an adversary  $\mathcal{A} = (\mathcal{A}_{\text{inp}}, \mathcal{A}_{\text{prf}})$  for non-adaptive black-box knowledge soundness in fig. 11. This exactly follows the intuition we discussed above.

$\mathcal{B}_{sem}^1(\text{pk})$	$\mathcal{B}_{sem}^2(\text{st} = (\text{crs}, \text{td}, r), \text{ct})$
$(\text{pk}_{\text{FHE}}, \text{sk}_{\text{FHE}}) \leftarrow \text{FHE.KG}(1^\lambda)$ $\text{hk} \leftarrow \mathcal{K}_{\text{UHF}}$ $(\hat{\text{crs}}, \hat{\text{td}}) \leftarrow \Pi_{\exists}.\text{Setup}(1^\lambda)$ $i^* \leftarrow \mathcal{N}_w$ ; $\text{ct}_{i^*} \leftarrow \text{FHE.Enc}(\text{pk}_{\text{FHE}}, i^*)$ $\text{crs} \leftarrow (\hat{\text{crs}}, \text{ct}_{i^*}, \text{hk}, \text{pk}_{\text{FHE}}, \text{pk})$ $\text{td} \leftarrow (\text{sk}_{\text{FHE}}, \hat{\text{td}})$ Sample random coins $r$ for $\mathcal{A}$ $b \leftarrow \mathcal{A}^{\text{O}_2}(\text{crs}; r)$ Recover $\mathcal{A}$ 's $i$ -th query $(x^{(i)}, w^{(i)})$ <b>return</b> $(\text{st} \leftarrow (\text{crs}, \text{td}, r), m_0 \leftarrow \text{H}_{\text{hk}}(w^{(i)}), m_1 \leftarrow 0)$	<b>return</b> $\mathcal{A}^{\text{O}_{sem}(\text{st}, \text{ct}, \cdot)}$ , where For queries $j = 1, \dots, i-1$ : $\text{O}_{sem} = \text{O}_2$ On $i$ -th query $\text{O}_{sem}(\text{st}, \text{ct}, x, w)$ : $\text{ct}_h \leftarrow \text{ct}$ $\text{ct}_{\text{bit}} \leftarrow \text{FHE.Eval}(\text{pk}_{\text{FHE}}, f_{\text{proj}}, \text{ct}_{i^*}, w)$ $\hat{x} \leftarrow (x, \text{hk}, \text{pk}_{\text{FHE}}, \text{pk}, \text{ct}_h, \text{ct}_{i^*}, \text{ct}_{\text{bit}})$ $\pi \leftarrow \Pi_{\exists}.\text{Sim}(\hat{\text{crs}}, \hat{\text{td}}, \hat{x})$ <b>return</b> $\pi^* := (\pi, \text{ct}_h, \text{ct}_{\text{bit}})$ For queries $j > i$ : $\text{O}_{sem} = \text{O}_1$

Fig. 10:  $\text{Game}_{2:i-1} \rightarrow \text{Game}_{2:i}$  reduction in lemma 10

$\mathcal{A}_{inp}(1^\lambda)$ $(\text{pp}, \text{uk}) \leftarrow \text{KGen}(1^\lambda)$ ; $w \leftarrow \text{Sample}(\text{pp})$ ; $y \leftarrow \text{Eval}(\text{pp}, w)$ ; <b>return</b> $(x = (\text{pp}, y), \text{st} = (\text{pp}, \text{uk}, y, w))$	$\mathcal{A}_{prf}(\text{st}, \text{crs})$ Parse $\text{st} = (\text{pp}, \text{uk}, y, w)$ ; $w' \leftarrow \text{Update}(\text{uk}, w)$ ; $\pi \leftarrow \text{P}(\text{crs}, (\text{pp}, y), w')$ ; <b>return</b> $\pi$
---	---

Fig. 11: Non-adaptive black-box knowledge soundness adversary for theorem 3

$\mathcal{B}^{\text{O}_L(\cdot)}(\text{pp}, y)$ $(\text{crs}, \text{td}) \leftarrow \text{KGen}(1^\lambda)$ $\pi \leftarrow \text{Sim}_{\text{pp}, y}^{\text{O}_L(\cdot)}(\text{crs})$ Run $\text{Ext}^{\text{Sim}_{\text{pp}, y}^{\text{O}_L(\cdot)}(\cdot)}(\text{crs}, \text{td}, (\text{pp}, y), \pi)$	$\text{Sim}_{\text{pp}, y}^{\text{O}_L(\cdot)}(\text{crs})$ Define $h_{\text{crs}, y, \text{pp}}(X) := \text{P}(\text{crs}, (\text{pp}, y), X)$ Query $\pi \leftarrow \text{O}_L(h_{\text{crs}, y, \text{pp}})$ <b>return</b> $\pi$
---	--

Fig. 12: Continuous leakage-resilience adversary  $\mathcal{B}$

Next, we construct an adversary  $\mathcal{B}$  against CLR in fig. 12. Idea is that we want to use the extractor  $\text{Ext}$  of  $\Pi$  to recover the OWF preimage. To do so,  $\mathcal{B}$  must provide  $\text{Ext}$  with proofs which are created by  $\text{crs}$  chosen by  $\text{Ext}$ . Since proofs depend on the OWF preimage  $w$ ,  $\mathcal{B}$  can use the leakage query oracle  $\text{O}_L$  with a function  $h_{\text{crs}, y, \text{pp}}(X) := \text{P}(\text{crs}, (\text{pp}, y), X)$ . This is possible only because the proof size is  $\leq L(\lambda)$  bits. In fig. 12,  $\mathcal{B}$  runs a subroutine  $\text{Sim}_{\text{pp}, y}^{\text{O}_L(\cdot)}(\text{crs})$  for creating proofs.

In the following, we will analyze the success probability of  $\mathcal{B}$ . Essentially, we show that if  $\text{Ext}$  succeeds in extracting with high probability, then also  $\mathcal{B}$  will succeed in breaking  $L$ -continuous leakage-resilience of OWF with high probability.

Game<sub>0</sub>. This is the original  $L$ -CLR game with the adversary  $\mathcal{B}$  from section 2.1.

Game<sub>1</sub>. This is the same game with  $\mathcal{B}$  being in-lined. See fig. 13. Obviously the probability that  $y = \text{Eval}(\text{pp}, x')$  is the same in both games.

Game<sub>2</sub>. This game is again just a slight rewrite of the previous  $\text{Game}_1$ . Note that the first three lines of  $\text{Game}_1$  in fig. 13 are equivalent to  $\mathcal{A}_{inp}(1^\lambda)$ . So we instead write  $(x = (\text{pp}, y), \text{st} = (\text{pp}, \text{uk}, y, w)) \leftarrow \mathcal{A}_{inp}(1^\lambda)$  in  $\text{Game}_2$ . Moreover,  $\text{Sim}_{\text{pp}, y}^{\text{O}_L(\cdot)}(\text{crs})$  and  $\mathcal{A}_{prf}(\text{st}, \text{crs})$  produce the exact same proof  $\pi$ . We change  $\text{Sim}_{\text{pp}, y}^{\text{O}_L(\cdot)}(\text{crs})$  to  $\mathcal{A}_{prf}(\text{st}, \text{crs})$  in  $\text{Game}_2$ . Clearly again the probability of  $y = \text{Eval}(\text{pp}, x')$  is the same as before.

<p><b>Game<sub>1</sub>(1<sup>λ</sup>)</b></p> <hr/> <p>(pp, uk) ← KGen(1<sup>λ</sup>)  w ← Sample(pp)  y ← Eval(pp, w)  (crs, td) ← KGen(1<sup>λ</sup>)  π ← Sim<sub>pp,y</sub><sup>O<sub>L</sub>(·)</sup>(crs)  x' ← Ext<sup>Sim<sub>pp,y</sub><sup>O<sub>L</sub>(·)</sup>(·)</sup>(crs, td, (pp, y), π)  <b>return</b> y = Eval(pp, x')</p>	<p><b>Game<sub>2</sub>(1<sup>λ</sup>)</b></p> <hr/> <p>(x = (pp, y), st = (pp, uk, y, w)) ← <math>\mathcal{A}_{inp}(1^\lambda)</math>  (crs, td) ← KGen(1<sup>λ</sup>)  π ← <math>\mathcal{A}_{prf}(st, crs)</math>  x' ← Ext<sup><math>\mathcal{A}_{prf}(st, crs)</math></sup>(crs, td, (pp, y), π)  <b>return</b> y = Eval(pp, x')</p>
---	--

Fig. 13: Security games for theorem 3

Note that Game<sub>2</sub> with the winning condition  $V(\text{crs}, (\text{pp}, y), \pi) = 1 \wedge y \neq \text{Eval}(\text{pp}, x')$  is the non-adaptive black-box knowledge soundness game. We know that this probability is bounded by  $\varepsilon_{ks}(\lambda)$ . Thus,  $V(\text{crs}, (\text{pp}, y), \pi) \neq 1 \vee y = \text{Eval}(\text{pp}, x')$  happens with a probability  $> 1 - \varepsilon_{ks}(\lambda)$ . However,  $V(\text{crs}, (\text{pp}, y), \pi) \neq 1$  is not possible given the construction of  $\mathcal{A}_{prf}$ . Thus,  $V(\text{crs}, (\text{pp}, y), \pi) \neq 1 \vee y = \text{Eval}(\text{pp}, x')$  is equivalent to  $y = \text{Eval}(\text{pp}, x')$ , which is the Game<sub>2</sub> winning condition. It follows that  $\mathcal{B}$  can break  $L$ -CLR with the probability  $1 - \varepsilon_{ks}(\lambda)$ .  $\square$

### C.3 Proofs of Lemmas 2 and 3

**Lemma 11 (Restatement of lemma 2).** *If an indexed languages  $\mathcal{L} \in \mathbf{NP}$  has a sub-exponentially hard-on-average problem, then for any  $d > 0$ ,  $\mathcal{L}$  also has a hard-on-average problem with  $(2^{\lambda^d}, 1/2^{\lambda^d})$ -indistinguishability.*

*Proof.* Let us recall that  $X_\lambda$  and  $\bar{X}_\lambda$  are said to be sub-exponentially indistinguishable if there exists a  $\delta > 0$  such that  $X_\lambda$  and  $\bar{X}_\lambda$  are  $(2^{\Omega(\lambda^\delta)}, 1/2^{\Omega(\lambda^\delta)})$ -indistinguishable. Let us define  $n(\lambda) = \lceil \lambda^{d/\delta} \rceil$ . Then  $Y_\lambda := X_{n(\lambda)}$  and  $\bar{Y}_\lambda := \bar{X}_{n(\lambda)}$  are  $(s(\lambda), \varepsilon(\lambda))$ -indistinguishable, where  $s(\lambda) = 2^{\Omega(n^\delta)} = 2^{\Omega(\lceil \lambda^{d/\delta} \rceil^\delta)}$  and  $\varepsilon(\lambda) = 1/2^{\Omega(n^\delta)} = 1/2^{\Omega(\lceil \lambda^{d/\delta} \rceil^\delta)}$ . Firstly, since the circuit of size  $s(\lambda) = 2^{\Omega(\lceil \lambda^{d/\delta} \rceil^\delta)}$  grows faster than the circuit of size  $2^{\lambda^d}$ , then, for a sufficiently large  $\lambda$ ,  $Y_\lambda$  and  $\bar{Y}_\lambda$  are also  $(2^{\lambda^d}, \varepsilon(\lambda))$ -indistinguishable. Conversely,  $Y_\lambda$  and  $\bar{Y}_\lambda$  are  $(2^{\lambda^d}, \varepsilon'(\lambda))$ -indistinguishable if  $\varepsilon'(\lambda) \geq \varepsilon(\lambda)$ . This is the case for  $\varepsilon'(\lambda) = 1/2^{\lambda^d}$  if  $\lambda$  is again sufficiently large. It follows that for a large enough  $\lambda$ , exists  $Y_\lambda$  and  $\bar{Y}_\lambda$  that are  $(2^{\lambda^d}, 1/2^{\lambda^d})$ -indistinguishable.

Let  $i$  be the index sampler and  $(\text{Samp}_{\mathcal{L}}, \text{Samp}_{\bar{\mathcal{L}}})$  instance samplers for sub-exponentially hard-on-average problem. Then we can always define  $\mathcal{I}'(1^\lambda)$  as  $\mathcal{I}(1^{n(\lambda)})$ ,  $\text{Samp}'_{\mathcal{L}}(1^\lambda, i)$  as  $\text{Samp}'_{\mathcal{L}}(1^{n(\lambda)}, i)$ , and  $\text{Samp}_{\bar{\mathcal{L}}}(1^\lambda, i)$  as  $\text{Samp}_{\bar{\mathcal{L}}}(1^{n(\lambda)}, i)$ , which gives the desired hard-on-average problem with  $(2^{\lambda^d}, 1/2^{\lambda^d})$ -indistinguishability.  $\square$

**Lemma 12 (Restatement of lemma 3).** *Let  $X_\lambda$  and  $\bar{X}_\lambda$  be  $(2^{\lambda^d}, 1/2^{\lambda^d})$ -indistinguishable distributions for some integer  $d \geq 2$ . Let  $A_\lambda$  over  $(x, \pi)$  be an augmented distribution of  $X_\lambda$ , where  $|\pi| = \ell(\lambda) = o(\lambda^d)$ . Then there exists an augmented distribution  $\bar{A}_\lambda$  of  $\bar{X}_\lambda$  such that  $A_\lambda$  and  $\bar{A}_\lambda$  are  $(\text{poly}(\lambda), \text{negl}(\lambda))$ -indistinguishable.*

*Proof.* Let  $X_\lambda$  and  $\bar{X}_\lambda$  be  $(s(\lambda), \varepsilon(\lambda))$ -indistinguishable, where  $s(\lambda) = 2^{\lambda^d}$  and  $\varepsilon(\lambda) = 1/2^{\lambda^d}$ . Then  $X_\lambda$  and  $\bar{X}_\lambda$  are  $(s(\lambda), \varepsilon'(\lambda))$ -indistinguishable for any  $\varepsilon'(\lambda) \geq 1/2^{\lambda^d}$ . Let us take  $\varepsilon'(\lambda) = 1/2^{\lambda^{d-1}}$ . According to the leakage lemma, there exists a polynomial  $p$  and an augmented distribution  $\bar{A}_\lambda$  of  $\bar{X}_\lambda$  such that  $A_\lambda$  and  $\bar{A}_\lambda$  are  $(s^*(\lambda), \varepsilon^*(\lambda))$ -indistinguishable where  $s^*(\lambda) = s(\lambda)p(\varepsilon'(\lambda)/2^{\ell(\lambda)})$  and  $\varepsilon^*(\lambda) = 2\varepsilon'(\lambda)$ .

Clearly  $\varepsilon^*(\lambda)$  is negligible since  $\varepsilon^*(\lambda) = 2\varepsilon'(\lambda) = 2/2^{\lambda^{d-1}} = 1/2^{\lambda^{d-1}-1} = \text{negl}(\lambda)$ . The distinguisher circuit size  $s^*(\lambda)$  is  $s^*(\lambda) = s(\lambda)p(\varepsilon'(\lambda)/2^{\ell(\lambda)}) = 2^{\lambda^d} p(2^{-\lambda^{d-1}-o(\lambda^d)})$ . Here,  $p(2^{-\lambda^{d-1}-o(\lambda^d)}) = 2^{-o(\lambda^{d-1})}$  and therefore  $s^*(\lambda) = 2^{\lambda^d - o(\lambda^{d-1})} = 2^{\Omega(\lambda^d)}$ . This means that indistinguishability holds also for all polynomial size circuits, as  $s^*(\lambda)$  grows faster than any polynomial.  $\square$

## D Non-Adaptive BB Extractable SNARK for UP

In fig. 14 we show a slightly simpler version of the construction in fig. 2. The following constructs a SNARK for a language in **UP** (restriction of **NP** where a statement in the language has exactly one accepting witness). In this variant we do not need hashing to fingerprint the witnesses. Our SNARK works as follows.

<p><b>Setup</b>(<math>1^\lambda</math>)</p> <hr/> <p> <math>(\hat{c}rs, \hat{t}d) \leftarrow II_{\exists}.Setup(1^\lambda)</math>  <math>i^* \leftarrow_{\\$} [N_w]</math>  <math>(pk_{FHE}, sk_{FHE}) \leftarrow FHE.KG(1^\lambda)</math>  <math>ct_{i^*} \leftarrow FHE.Enc(pk, i^*)</math>  <b>return</b> <math>(crs := (\hat{c}rs, ct_{i^*}, pk_{FHE}), td := (sk_{FHE}, \hat{t}d))</math> </p>
<p><b>P</b>(<math>crs, R, x, w</math>)</p> <hr/> <p> <math>ct_{bit} \leftarrow FHE.Eval(pk_{FHE}, f_{proj}, ct_{i^*}, w)</math>          where <math>f_{proj}(i, w) := w_i</math>  <math>\pi \leftarrow II_{\exists}.P(\hat{c}rs, R', (x, pk_{FHE}, ct_{i^*}, ct_{bit}), w)</math>          where <math>R'(x, pk_{FHE}, ct_{i^*}, ct_{bit}; w) \iff</math>  <math>R(x, w) \wedge ct_{bit} = FHE.Eval(pk_{FHE}, f_{proj}, ct_{i^*}, w)</math>  <b>return</b> <math>\pi^* := (\pi, ct_{bit})</math> </p>
<p><b>V</b>(<math>crs, R, x, \pi^*</math>)</p> <hr/> <p>         Parse <math>\pi^*</math> as <math>(\pi, ct_{bit})</math>  <b>return</b> <math>II_{\exists}.V(\hat{c}rs, R', (x, pk_{FHE}, ct_{i^*}, ct_{bit}))</math>          where <math>R'</math> is defined like above       </p>

Fig. 14: Non-adaptively secure black-box extractable construction for **UP**.  $N_w$  is a bound on the witness size.  $II_{\exists}$  is the SNARG scheme.

$\mathcal{A}_{slow}(crs)$	$\mathcal{A}_{fast}(crs)$
<b>if</b> $crs \in \text{img}(\text{Setup})$	$(x, w) \leftarrow \text{Samp}_{\mathcal{L}};$
Find $td$ for $crs$ ;	$\pi \leftarrow \mathbf{P}(crs, x, w);$
$x \leftarrow \text{Samp}_{\bar{\mathcal{L}}};$	<b>return</b> $(x, \pi);$
$\pi \leftarrow \mathcal{S}(crs, td, x);$	
<b>else</b>	
$(x, w) \leftarrow \text{Samp}_{\mathcal{L}};$	
$\pi \leftarrow \mathbf{P}(crs, x, w);$	
<b>return</b> $(x, \pi)$	

Fig. 15: Inefficient adversary  $\mathcal{A}_{slow}$  against adaptive soundness and its efficient emulator  $\mathcal{A}_{fast}$

## E A Summary of the GW Impossibility Proof [32]

The main technique used in the proof is showing the existence of a *simulatable adversary*  $\bar{\mathbf{P}}$  for any SNARG for an NP complete language. A simulatable adversary is an *inefficient* adversary that, given a CRS, outputs a false statement  $x \notin \mathcal{L}$  with a valid proof  $\pi$  for it. While the existence of such adversary is trivial for any SNARG, a simulatable adversary also comes with an *efficient* simulator  $\mathcal{S}$  such that no efficient distinguisher can distinguish them. To show the existence of a simulatable adversary, a lemma in [32] shows that for any two computationally indistinguishable distributions respectively over a set  $\mathcal{L}$  and its complement  $\bar{\mathcal{L}} = \{0, 1\}^* \setminus \mathcal{L}$ , and for any leakage information  $\pi$  on  $x \in \mathcal{L}$ , there exists some leakage information  $\bar{\pi}$  on  $\bar{x} \in \bar{\mathcal{L}}$  such that  $(x, \pi)$  and  $(\bar{x}, \bar{\pi})$  are also computationally indistinguishable. What is important in this lemma is that the security degrades exponentially with the size of the leakage  $\pi$  and this is the reason why the underlying SNARG  $\Pi$  should have succinctness property.

Now, given this simulatable adversary, the result can be concluded as follows. Assume there exists a black-box reduction  $R$  that shows the soundness of  $\pi$  based on  $(\mathcal{C}, c)$ . This means that the efficient reduction  $R^{\bar{\mathbf{P}}}$ , given black-box access to successful adversary  $\bar{\mathbf{P}}$  can break  $(\mathcal{C}, c)$ . But if (inefficient)  $R^{\bar{\mathbf{P}}}$  can break  $(\mathcal{C}, c)$ , then (efficient)  $R^{\mathcal{S}}$  can also break it since no efficient distinguisher (including the challenger of  $(\mathcal{C}, c)$ ) can distinguish  $\bar{\mathbf{P}}$  from  $\mathcal{S}$ . Thus, if this black-box reduction exists, then the assumption  $(\mathcal{C}, c)$  should be false.

## F A Summary of the Pass Impossibility Proof [56]

We briefly summarize the impossibility result of [56]. Let  $\mathcal{L}$  be a hard-on-average language as defined in section 2. Now we can give an informal statement for the result in [56].

**Theorem 11.** *Let  $\Pi$  be an adaptively sound and perfectly zero-knowledge non-interactive argument for an hard-on-average problem  $\mathcal{L}$ . Suppose that there exists an efficient black-box reduction  $R$  that can reduce adaptive soundness of  $\Pi$  to some falsifiable assumption  $C$ . Then  $C$  can be broken in polynomial time.*

The intuition behind the result is as follows. First, we construct an *inefficient* adversary  $\mathcal{A}_{slow}$  (the first algorithm in fig. 15) that can break adaptive soundness. If  $\mathcal{A}_{slow}$  gets a valid CRS as an input ( $crs$  is in the image of  $\text{Setup}$ ), then it brute-force computes a trapdoor  $td$ , samples a false statement  $x$ , and runs the simulator with  $td$  to produce a proof  $\pi$ . If the CRS is invalid (outside of the image of  $\text{Setup}$ ), then it just tries to compute a proof for an honest statement  $x$ . Note that in the adaptive soundness game the CRS will always be valid and thus only the first branch will matter. Since  $\mathcal{L}$  is a hard-on-average language, then false and true statements are indistinguishable, and therefore  $\mathcal{S}$  will produce a proof which is accepted by a verifier with an overwhelming probability. So indeed  $\mathcal{A}_{slow}$  does break adaptive soundness.

Let us suppose that there exists an efficient reduction  $R$  that given black-box access to any adaptive soundness adversary  $\mathcal{A}$ , can break some falsifiable assumption  $C$ . The problem is that although  $\mathcal{A}_{slow}$  does



P'	V'
<b>Constants:</b> PRF key $K$	<b>Constants:</b> PRF key $K$
<b>Input:</b> $(x, w)$	<b>Input:</b> $(x, \pi)$
<b>if</b> $(x, w) \in \mathcal{R}$	<b>if</b> $f(\pi) = f(\text{PRF}_K(x))$
<b>return</b> $\text{PRF}_K(x)$	<b>return</b> 1
<b>else</b>	<b>else</b>
<b>return</b> $\perp$	<b>return</b> 0

Fig. 16: Programs P' and V'

break adaptive soundness, it is not efficient. Therefore also the reduction  $R^{\mathcal{A}_{slow}}$  will be inefficient. We solve this issue by constructing an efficient emulator  $\mathcal{A}_{fast}$  (the second algorithm in fig. 15) for  $\mathcal{A}_{slow}$ .

The emulator  $\mathcal{A}_{fast}$  simply samples an honest  $(x, w)$  and generates an honest proof  $\pi$ . Now let us compare  $\mathcal{A}_{fast}$  and  $\mathcal{A}_{slow}$ . If CRS is invalid, then  $\mathcal{A}_{slow}$  and  $\mathcal{A}_{fast}$  are identical. However, if the CRS is valid, then it needs a bit more work to show that outputs are indistinguishable. Intuitively,  $x$  is indistinguishable due to the hard-on-average property and  $\pi$  is indistinguishable due to zero-knowledge. However, here it is important that zero-knowledge property holds even with respect to a fixed CRS since we do not know how the distinguishing adversary may pick the valid CRS. Moreover, with computational zero-knowledge it may be even possible to extract the witness from a proof  $\pi$  which would make distinguishing  $\mathcal{A}_{slow}$  and  $\mathcal{A}_{fast}$  trivial. This is the reason why zero-knowledge has to be perfect (or statistical). It follows now that outputs of  $\mathcal{A}_{fast}$  and  $\mathcal{A}_{slow}$  are computationally indistinguishable.

Since  $R^{\mathcal{A}_{slow}}$  can break the assumption  $C$ , then so does  $R^{\mathcal{A}_{fast}}$  which means that the assumption  $C$  is insecure. Hence, it is impossible to base adaptively sound perfect zero-knowledge argument on a falsifiable assumption using a black-box reduction.

## G Non-Adaptive SNARGs With Perfect ZK Based on iO

We show how for the non-adaptive case, none of [56] and [32] results hold. We do so by the following observation: assuming that indistinguishability obfuscation (iO) can be built from falsifiable assumptions (see [62]), the perfect NIZK arguments of Sahai and Waters [58], instantiated with a puncturable PRF (PPRF) that satisfies *succinctness* property, is a non-adaptive SNARG with perfect ZK for all NP languages in the CRS model. While this can be seen as a feasibility result, proposing a construction with more standard assumptions (i.e., without iO) is still an interesting open question.

We now recall the NIZK arguments of Sahai and Waters [58].

**NIZK arguments of Sahai and Waters.** The idea is very simple: the proof system consists of two obfuscated programs put in the CRS. The first program is the proving algorithm that inputs a statement  $x$  and witness  $w$  and outputs a signature on  $x$  if  $(x, w) \in \mathcal{R}$ . The signature is realized by a PRF in the construction. The second program is the verification algorithm that is just the signature verification and verifies the proof by checking the validity of the signature on  $x$ .

Let PRF be a puncturable PRF that inputs  $\ell$ -bit long strings and outputs  $\lambda$  bits (where  $\lambda$  is the security parameter). Let  $f(\cdot)$  be a one way function. The NIZK argument  $\Pi = (\text{Setup}, P, V)$  for language  $\mathcal{L}$  with relation  $\mathcal{R}$  is as follows:

- $\text{Setup}(1^\lambda)$  first selects a puncturable PRF key  $K$  for PRF. Next, it creates an obfuscation of programs P' and V' as depicted in Figure 16. The CRS  $\text{crs}$  consists of the two obfuscated programs.
- $P(\text{crs}, x, w)$  runs the obfuscated program P' on input  $(x, w)$  and returns the proof  $\pi$  if  $(x, w) \in \mathcal{R}$ .
- $V(\text{crs}, x, \pi)$  runs the obfuscated program V' on input  $(x, \pi)$  and returns a bit indicating accept or reject.

**Theorem 12.** [58] *The argument system  $\Pi$  is perfectly zero-knowledge. Moreover, if the obfuscation scheme is indistinguishably secure, PRF is a secure punctured PRF with succinctness property, and  $f(\cdot)$  is an injective one way function, then  $\Pi$  is a non-adaptive SNARG.*

*Remark 4.* While SNARGs with non-adaptive security can be seen as interactive two-message arguments by thinking of the CRS as the verifier’s message, the type of non-adaptivity in the resulting argument is still “strong” in the sense that the verifier’s message does not depend on the prover’s (fixed) statement. One can also define a weaker notion of non-adaptivity for two-message arguments where the first message is *statement-dependent* (See [7] for example). We note that while the above iO-based construction satisfies the stronger notion, giving a construction for the weaker notion of non-adaptivity based on seemingly weaker tools is not a hard task. Namely, the verifier can use a witness encryption scheme to encrypt a succinct random value  $r$  under the prover’s statement and ask the prover to return  $r$ .