

*Appeared in **SIAM Journal on Computing**
27(4):1116-1141, August 1998*

Computational Complexity and Knowledge Complexity*

Oded Goldreich[†] Rafail Ostrovsky[‡] Erez Petrank[§]

Appeared in SIAM Journal on Computing
27(4):1116-1141, August 1998

Abstract

We study the computational complexity of languages which have interactive proofs of logarithmic knowledge-complexity. We show that all such languages can be recognized in $\mathcal{BPP}^{\mathcal{NP}}$. Prior to this work, for languages with greater-than-zero knowledge-complexity only trivial computational complexity bounds were known. In the course of our proof, we relate statistical knowledge-complexity with perfect knowledge-complexity; specifically, we show that, for the honest verifier, these hierarchies coincide, up to a logarithmic additive term.

*An extended abstract of this paper appeared in the *26th ACM Symposium on Theory of Computing* (STOC 94), held in Montreal, Quebec, Canada, May 23-25, 1994.

[†]Department of Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot, Israel. E-mail: oded@wisdom.weizmann.ac.il. Supported by grant no. 92-00226 from the United States — Israel Binational Science Foundation, Jerusalem, Israel.

[‡]Bell Communications Research, 445 South Street, Morristown, New Jersey 07960-6438. E-mail: rafail@bellcore.com. Part of this work was done at University of California at Berkeley and International Computer Science Institute at Berkeley and supported by an NSF postdoctoral fellowship and ICSI.

[§]Department of Computer Science, University of Toronto, Toronto, Ontario, Canada M5S 3G4. E-mail: erez@cs.toronto.edu

1 Introduction

The notion of knowledge-complexity was introduced in the seminal paper of Goldwasser Micali and Rackoff [GMR-85, GMR-89]. Knowledge-complexity is intended to measure the *computational advantage* gained by interaction. A formulation of knowledge-complexity, for the case that it is not zero, has appeared in [GP-91]. A very appealing suggestion, made by Goldwasser Micali and Rackoff, is to characterize languages according to the knowledge-complexity of their interactive proof systems [GMR-89].

The lowest level of the knowledge-complexity hierarchy is the class of languages having interactive proofs of knowledge-complexity zero, better known as zero-knowledge. Actually, there are three hierarchies extending the three standard definitions of zero-knowledge; that is, *perfect*, *statistical* and *computational*. Assuming the existence of one-way functions, the third hierarchy collapses; that is, the zero level of the *computational* knowledge-complexity hierarchy contains all languages having interactive proof systems [GMW-86, IY-87, B+ 88], and thus contains all levels of the (computational) knowledge-complexity hierarchy. In this paper we will be only interested in the other two hierarchies. Previous works have provided information only concerning the zero level of these hierarchies (see, for example, Fortnow [F-89] and Aiello and Hastad [AH-87]). Our main result is an upper bound on the computational complexity of languages having logarithmic (statistical) knowledge-complexity; namely, we show that such languages are contained in $\mathcal{BPP}^{\mathcal{NP}}$.

We consider the (statistical) knowledge-complexity hierarchy to be a very natural one. Its lowest level resides in the second level of the Polynomial-Time Hierarchy (cf., [F-89, AH-87, B-85]), whereas as a whole it covers all of \mathcal{PSPACE} . Another hierarchy with a similar feature, which also deserves investigation, is the hierarchy of languages classified by the number of rounds in their (“shortest”) interactive proof system. Interestingly, the latter hierarchy has a multiplicative collapse (cf., [BM-88]) whereas no such result is known for the knowledge-complexity hierarchy.

1.1 Background on knowledge-complexity

Loosely speaking, an interactive-proof system for a language L is a two-party protocol, by which a powerful *prover* can “convince” a probabilistic polynomial-time *verifier* of membership in L , but will fail (with high probability) when trying to fool the verifier into “accepting” non-members [GMR-89]. An interactive-proof is called *zero-knowledge* if the interaction of any probabilistic polynomial-time machine with the predetermined prover, on common input $x \in L$, can be “simulated” by a probabilistic polynomial-time machine (called the *simulator*), given only x [GMR-89]. We

say that a probabilistic machine M *simulates* an interactive proof if the output distribution of M is *statistically close* to the distribution of the real interaction between the prover and the verifier.

The formulation of zero-knowledge presented above is known as *statistical* (almost-perfect) zero-knowledge. Alternative formulations of zero-knowledge are *computational* zero-knowledge and *perfect* zero-knowledge. In this paper we concentrate on statistical zero-knowledge and the knowledge-complexity hierarchy that generalizes it.

Loosely speaking, the knowledge-complexity of a protocol Π is the number of oracle bits that are needed to simulate the protocol efficiently. Namely, we say that a prover leaks $k(\cdot)$ bits of knowledge to verifier V if there is a probabilistic polynomial-time oracle machine (“simulator”) M such that on any input $x \in L$, machine M makes at most $k(|x|)$ oracle queries, and the output distribution of $M(x)$ is statistically close to the distribution of the conversations in the interaction between the prover and V . For a formal definition and further discussion, see §2.2.

The knowledge-complexity of a language is the minimum knowledge-complexity of an interactive proof system for that language. We consider the knowledge-complexity of a language to be a very natural parameter. Furthermore, the question of how does this parameter relate to the complexity of deciding the language is a very fundamental question.

1.2 Previous work

The complexity of recognizing zero-knowledge languages was first bounded by Fortnow [F-89]. He showed that any language that admits a zero-knowledge interactive-proof is in the class coAM . Subsequently, Aiello and Hastad [AH-87] showed that these languages are also in \mathcal{AM} .

Bellare and Petrank [BP-92] bounded the computational complexity of languages which have *short* interactive-proofs with *low* knowledge-complexity. Specifically, they showed that if a language L has a $g(n)$ -round interactive-proof which leaks at most $k(n)$ bits of knowledge and if $k(n) \cdot g(n) = O(\log n)$ then the language can be recognized in $\mathcal{BPP}^{\mathcal{NP}}$. This result does not apply to the general class of low knowledge-complexity languages, since these languages might not have interactive-proofs which are both of small round complexity and low knowledge-complexity.

1.3 This work

In this work we extend the result of [BP-92], showing that any language having an interactive proof with logarithmic knowledge-complexity, can be recognized in

BPP^{NP} . We recall that BPP^{NP} is contained in the third level of the polynomial-time hierarchy (\mathcal{PH}). It is believed that \mathcal{PH} is a proper subset of \mathcal{PSPACE} . Thus, assuming $\mathcal{PH} \subsetneq \mathcal{PSPACE}$, our result yields the first proof that there exist languages in \mathcal{PSPACE} which cannot be proven by an interactive-proof that yields $O(\log n)$ bits of knowledge. In other words, there exist languages which do have interactive proofs but only interactive proofs with super-logarithmic knowledge-complexity. We stress that prior to our work, there was no indication that would contradict the possibility that all languages in \mathcal{PSPACE} have interactive-proofs which yield only *one bit of knowledge*.

Our proof that languages of logarithmic knowledge-complexity are in BPP^{NP} consists of two parts. In the first part, we show that the BPP^{NP} procedure described by Bellare and Petrank [BP-92] is applicable for recognizing languages having interactive proofs of logarithmic *perfect* knowledge-complexity. To this end, we use a more careful analysis than the one used in [BP-92]. In the second part of our proof, we transform interactive proofs of *statistical* knowledge-complexity $k(n)$ into interactive proofs of *perfect* knowledge-complexity $k(n) + O(\log n)$. This transformation refers only to knowledge-complexity with respect to the honest verifier, but this suffices since the first part of our proof applies to the knowledge-complexity with respect to the honest verifier. Yet, the transformation is interesting for its own sake, and a few words are in place.

The question of whether statistical zero-knowledge equals perfect zero-knowledge is one of the most fundamental open problems regarding zero-knowledge. The question has been open also for the case of zero-knowledge with respect to the honest verifier. Our transformation implies, as a special case, that any statistical zero-knowledge interactive proof can be modified into an interactive proof of perfect knowledge-complexity bounded by a logarithmic function. Following the conference presentation of our work, Aiello, Bellare and Venkatesan showed that statistical zero-knowledge coincides with negligible *on the average* perfect knowledge-complexity [ABV-95]. Their result is stronger in two respects; it refer to all verifiers, not only the honest verifier, and it bounds the perfect knowledge-complexity by a negligible function rather than by a logarithmic one. On the other hand, their result is weaker as it refers to a much more liberal notion of (perfect) knowledge-complexity; that is, “average” knowledge-complexity rather than “worst case” knowledge-complexity (as considered here).

1.4 Organization

In Section 2, we present the main definitions referred to in the rest of the paper. These include the definition of an interactive proof system as well as the definition of its knowledge-complexity. Section 3 provides an overview to our proof that languages having (statistical) knowledge-complexity bounded by a logarithmic function reside in $\mathcal{BPP}^{\mathcal{NP}}$. The first part of our proof (i.e., the case of perfect knowledge-complexity) is presented in Section 4. The second part of our proof (i.e., the transformation of statistical knowledge-complexity to perfect knowledge-complexity) is presented in Section 5. Some concluding remarks appear in Section 6.

In Appendix A, we discuss a flaw in Fortnow's paper [F-89]. We stress that the main result of [F-89] as well as its main techniques remain valid.

2 Preliminaries

Let us state some of the definitions and conventions we use in the paper. Throughout this paper we use n to denote the length of the input x . A function $f : \mathbb{N} \rightarrow [0, 1]$ is called *negligible* if for every polynomial p and all sufficiently large n 's $f(n) < \frac{1}{p(n)}$.

2.1 Interactive proofs

Let us recall the concept of interactive proofs, presented by [GMR-89]. For formal definitions and motivating discussions the reader is referred to [GMR-89]. A protocol between a (computationally unbounded) *prover* P and a (probabilistic polynomial-time) *verifier* V constitutes an **interactive proof** for a language L if there exists a negligible function $\epsilon : \mathbb{N} \rightarrow [0, 1]$ such that

1. **Completeness:** If $x \in L$ then

$$\Pr[(P, V)(x) \text{ accepts}] \geq 1 - \epsilon(n)$$

2. **Soundness:** If $x \notin L$ then for any prover P^*

$$\Pr[(P^*, V)(x) \text{ accepts}] \leq \epsilon(n)$$

2.2 Knowledge Complexity

Throughout the rest of the paper, we only refer to knowledge-complexity *with respect to the honest verifier*; namely, the ability to simulate the honest verifier's view of its

interaction with the prover. (In the stronger definition, one considers the ability to simulate the point of view of *any efficient verifier* while interacting with the prover.)

We let $(P, V)(x)$ denote the random variable that represents V 's view of the interaction with P on common input x . The view contains the verifier's random tape as well as the sequence of messages exchanged between the parties.

We begin by briefly recalling the definitions of perfect and statistical zero-knowledge. A protocol (P, V) is *perfect zero-knowledge* (resp., *statistical zero-knowledge*) over a language L if there is a probabilistic polynomial-time simulator M such that for every $x \in L$ the random variable $M(x)$ is distributed identically to $(P, V)(x)$ (resp., the statistical difference between $M(x)$ and $(P, V)(x)$ is a negligible function in $|x|$).

Next, we present the definitions of perfect (resp., statistical) knowledge-complexity which we use in the sequel. These definitions extend the definition of perfect (resp., statistical) zero-knowledge, in the sense that knowledge-complexity zero is exactly zero-knowledge. Actually, there are two alternative formulations of knowledge-complexity, called the *oracle version* and the *fraction version*. These formulations coincide at the zero level and differ by at most an additive constant otherwise [GP-91]. For further intuition and motivation see [GP-91]. It will be convenient to use both definitions in this paper.

By the *oracle formulation*, the knowledge-complexity of a protocol (P, V) is the number of oracle (bit) queries that are needed to simulate the protocol efficiently. That is

Definition 2.1 (knowledge-complexity — oracle version): *Let $k: \mathbb{N} \rightarrow \mathbb{N}$. We say that an interactive proof (P, V) for a language L has perfect (resp., statistical) knowledge-complexity $k(n)$ in the oracle sense if there exists a probabilistic polynomial-time oracle machine M and an oracle A such that:*

1. *On input $x \in L$, machine M queries the oracle A for at most $k(|x|)$ bits.*
2. *For each $x \in L$, machine M^A produces an output with probability at least $\frac{1}{2}$, and given that M^A halts with an output, $M^A(x)$ is identically distributed (resp., statistically close) to $(P, V)(x)$.*

In the *fraction formulation*, the simulator is not given any explicit help. Instead, one measures the density of the largest subspace of simulator's executions (i.e., coins) which is identical (resp., statistically close) to the (P, V) distribution.

Definition 2.2 (knowledge-complexity — fraction version): *Let $\rho: \mathbb{N} \rightarrow (0, 1]$. We say that an interactive proof (P, V) for a language L has perfect (resp., statistical) knowledge-complexity $\log_2(1/\rho(n))$ in the fraction sense if there exists a probabilistic*

polynomial-time machine M with the following good subspace property. For any $x \in L$ there is a subset of M 's possible random tapes, denoted S_x , such that:

1. The set S_x contains at least a $\rho(|x|)$ fraction of the set of all possible coin tosses of $M(x)$.
2. Conditioned on the event that $M(x)$'s coins fall in S_x , the random variable $M(x)$ is identically distributed (resp., statistically close) to $(P, V)(x)$. Namely, for the perfect case this means that for every \bar{c}

$$\text{Prob}(M(x, \omega) = \bar{c} \mid \omega \in S_x) = \text{Prob}((P, V)(x) = \bar{c})$$

where $M(x, \omega)$ denotes the output of the simulator M on input x and coin tosses sequence ω .

As mentioned above, these two measures are almost equal.

Theorem [GP-91]: *The fraction measure and the oracle measure are equal up to an additive constant.*

Since none of our results is sensitive to a difference of an additive constant in the measure, we ignore this difference in the subsequent definition as well as in the statement of our results.

Definition 2.3 (knowledge-complexity classes):

- $\mathcal{PKC}(k(\cdot)) =$ languages having interactive proofs of perfect knowledge-complexity $k(\cdot)$.
- $\mathcal{SKC}(k(\cdot)) =$ languages having interactive proofs of statistical knowledge-complexity $k(\cdot)$.

2.3 The simulation based prover

An important ingredient in our proof is the notion of a simulation based prover, introduced by Fortnow [F-89]. Consider a simulator M that outputs conversations of an interaction between a prover P and a verifier V . We define a new prover P^* , called *the simulation based prover*, which selects its messages according to the conditional probabilities induced by the simulation. Namely, on a partial history h of a conversation, P^* outputs a message α with probability

$$\text{Prob}(P^*(h) = \alpha) \stackrel{\text{def}}{=} \text{Prob}(M_{|h|+1} = h \circ \alpha \mid M_{|h|} = h)$$

where M_t denotes the distribution induced by M on t -long prefixes of conversations. (Here, the length of a prefix means the number of messages in it.)

It is important to note that the behavior of P^* is *not* necessarily close to the behavior of the original prover P . Specifically, if the knowledge-complexity is greater than 0 and we consider the simulator guaranteed by the fraction definition, then P^* and P might have quite a different behavior. Our main objective will be to show that even in this case P^* still behaves in a manner from which we can benefit.

3 Overview

Using Definition 2.3, we state the main result of this paper as

Main Theorem: $SKC(O(\log(\cdot))) \subseteq BPP^{\mathcal{N}^P}$.

We recall that all that was previously known regarding the $SKC(\cdot)$ hierarchy is $SKC(0) \subseteq \mathcal{AM} \cap \text{co}\mathcal{AM}$ and $BPP \subseteq SKC(k) \subseteq SKC(k+1) \subseteq \mathcal{PSPACE}$, for every $k : \mathbb{N} \mapsto \mathbb{N}$.

The Main Theorem is proven in two stages

1. $\mathcal{PKC}(O(\log(\cdot))) \subseteq BPP^{\mathcal{N}^P}$ (see Theorem 1).
2. $SKC(k(\cdot)) \subseteq \mathcal{PKC}(k(\cdot) + O(\log(\cdot)))$, for every $k : \mathbb{N} \mapsto \mathbb{N}$ (see Theorem 2).

In the rest of this section we make several remarks regarding the above theorems and provide an overview to their proofs.

3.1 Remarks

Remark 1: Usually, the definition of interactive proofs is robust in the sense that setting the error probability to be bounded away from $\frac{1}{2}$ does not change their expressive power, since the error probability can be reduced by repetitions. However, this standard procedure is NOT applicable when knowledge-complexity is measured, since (even sequential) repetitions may increase the knowledge-complexity. The question is, thus, what is the *right* definition. The definition used in §2.1 is quite standard and natural; it is certainly less arbitrary than setting the error to be some favorite constant (e.g., $\frac{1}{3}$) or function (e.g., 2^{-n}). Yet, our techniques yield non-trivial results also in case one defines interactive proofs with non-negligible error probability (e.g., constant error probability). For example, languages having interactive proofs with

error probability $1/4$ and perfect knowledge-complexity 1 are also in $\mathcal{BPP}^{\mathcal{NP}}$. For more details see Appendix B.

Remark 2: In the definition used in §2.1 we have allowed two-sided error probability, rather than insisting on “perfect completeness” (as is sometimes done). This strengthens our Main Theorem but weakens the statistical to perfect transformation (i.e., Theorem 2), since a transformation for the case of one-sided error implies a transformation for the two-sided case¹, whereas the converse is not clear.

Remark 3: The definitions of knowledge-complexity in §2.2 refer to simulations of the honest verifier. Analogous definitions of knowledge-complexity refer to simulations of arbitrary polynomial-time verifiers (cf., [GP-91]). Let us denote the corresponding classes by $\mathcal{PKC}^*(\cdot)$ and $\mathcal{SKC}^*(\cdot)$. Clearly, $\mathcal{PKC}^*(k(\cdot)) \subseteq \mathcal{PKC}(k(\cdot))$ and $\mathcal{SKC}^*(k(\cdot)) \subseteq \mathcal{SKC}(k(\cdot))$, for every $k : \mathbb{N} \mapsto \mathbb{N}$. Thus, our Main Theorem is only strengthened by referring to the honest-verifier classes, whereas Theorem 2 is arguably weaker than an analogous statement referring to the arbitrary-verifier classes.

3.2 The Perfect Case – Overview

Our proof of Theorem 1 follows the procedure suggested in [BP-92], which in turn follows the approach of [F-89, BMO-90, Ost-91] while introducing a new “uniform generation” procedure which builds on ideas of [Si-83, St-83, GS-89, JVV-86].

Suppose that (P, V) is an interactive proof of perfect knowledge complexity $k(n) = O(\log n)$ for the languages L , and let M be the simulator guaranteed by the fraction formulation (i.e., Definition 2.2). We consider the conversations of the original verifier V with the simulation-based-prover P^* (see definition in §2.3). We show that the probability that the interaction (P^*, V) is accepting is negligible if $x \notin L$ and greater than a polynomial fraction if $x \in L$. Our proof differs from [BP-92] in the analysis of the case $x \in L$ (and thus we get a stronger result although we use the same procedure). This separation between the cases $x \notin L$ and $x \in L$ can be amplified by sequential repetitions of the protocol (P^*, V) . So it remains to observe that we can sample the (P^*, V) interactions in probabilistic polynomial-time having access to an NP-oracle. This observation originates from [BP-92] and is justified as follows. Clearly, V 's part of the interaction can be produced in polynomial-time.

¹Suppose one is given a transformation for the one-sided case. Then, given a two-sided interactive proof of some statistical knowledge-complexity one could first transform it to a one-sided error proof system of the same knowledge-complexity (cf., [GMS-87]). Applying the transformation for the one-sided case, to the resulting one-sided error proof system, yields an (one-sided error) interactive proof with the desired knowledge-complexity.

Also, by the uniform generation procedure of [BP-92] we can implement P^* by a probabilistic polynomial-time oracle machine that has access to an NP-oracle. Thus, it remains only to analyze the accepting probability of (P^*, V) on input x .

The case $x \notin L$ follows trivially from the soundness condition of V . The challenging case is when $x \in L$. If $k(n) = 0$ this case is easy since P^* behaves exactly as P and so the completeness condition guarantees that x will be accepted with very high probability. However, in case $k(n) > 0$ this argument is not valid and the simulator-based-prover may behave very differently from the prescribed prover. Note that it is possible to define a prover, P^{**} , based on the behavior of the simulator on the “good subspace” and that P^{**} will indeed behave as P . However, it is not clear if P^{**} can be implemented in a relatively efficient manner (e.g., by a probabilistic polynomial-time machine that has access to an NP-oracle). Thus, we need to analyze the behavior of (P^*, V) on $x \in L$. For sake of simplicity, we consider here only the special case in which $(P, V)(x)$ is always accepting (i.e., “perfect” completeness). Recall that the deviation of P^* from the behavior of P is due to the fact that behavior of the former is conditioned on the entire probability space of the simulator, whereas the latter is conditioned on the “good subspace”. In each case the next prover move is determined by the set of all simulator coins which match the current history of the interaction. For P^* this is the set of all coin tosses which may produce this history, whereas for P this is the set of all good coin tosses (i.e., coins in the “good subspace”) which produce this history. We first observe that the key parameter for the analysis of P^* is the ratio between the size of the residual probability space of the simulator and the size of the residual space of good coins. Actually, we consider the reciprocal of the above ratio. We observe that the *expected value* of the latter ratio may only increase as a function of the history length, where the expectation is taken over all possible histories of fixed length as produced by a (P^*, V) interaction. Finally, we observe that the expected value of the ratio for a full interaction is a lower bound on the probability that P^* makes V accept the input, whereas for the empty interaction the ratio equals 2^{-k} .

3.3 The Transformation – Overview

Our proof of Theorem 2 refers to the oracle formulation of knowledge-complexity (see Definition 2.1). Suppose we are given a simulator which produces output that is *statistically close* to the real prover–verifier interaction. We will change both the interactive proof and its simulation *so that they produce exactly the same distribution space*. We will take advantage of the fact that the prover in the interactive proof and the oracle that “assists” the simulator are both infinitely powerful. Thus, the

modification to the prover’s program and the augmentation to the oracle need not be efficiently computable. We stress that the modification to the simulator itself will be efficiently computable. Also, we maintain the original verifier (of the interactive proof), and thus the resulting interactive proof is still sound. Furthermore, the resulting interaction will be statistically close to the original one (on any $x \in L$) and therefore the completeness property of the original interactive proof is maintained (although the error probability here may increase by a negligible amount).

The key question is how can we modify the two relevant distributions so that they become identical rather than statistically close. The easy case is when some conversation is more likely in the simulation (than in the original prover–verifier interaction). This case is handled by providing the oracle with a candidate conversation and having the oracle decide probabilistically whether we should output this conversation or not. Thus, we can use one additional oracle query in order to lower the probability of conversations produced by the original simulator. However, the challenging case is when some conversation is less probable in the simulation (than in the original interaction). Using the oracle to produce such conversations is too costly, in terms of query complexity, unless we consider average-case query complexity (as in [ABV-95]). Thus, we need a different approach. Our approach is to modify the original prover so that it truncates conversations at a point they become less probable in the simulation. This truncation is also probabilistic. A new simulator, with the help of an augmented oracle, will have to detect the truncation point and produce truncated conversations with the same probability as they are produced in interaction with the new prover. In order to specify the truncation point we need to get a t -ary value from the oracle, where t is the total number of bits in the interaction. This is implemented using $\log_2 t$ queries giving rise to the additive logarithmic factor in the result of the theorem.

4 The Perfect Case

In this section we prove that the Main Theorem holds for the special case of *perfect* knowledge-complexity. Namely,

Theorem 1 $\mathcal{PKC}(O(\log n)) \subseteq \mathcal{BPP}^{\mathcal{NP}}$

As stated above, our proof follows the procedure suggested in [BP-92]. Suppose that (P, V) is an interactive proof of perfect knowledge-complexity $k(n) = O(\log n)$ for the languages L , and let M be the simulator guaranteed by Definition 2.2. Let us denote by P^* the simulation-based-prover (for M); see §2.3. Then,

Lemma 4.1 [BP-92]: P^* can be implemented by a probabilistic polynomial-time oracle machine that has access to an NP-oracle.

Lemma 4.2 (Analysis of the behavior of P^*):

1. If $x \in L$ then the probability that (P^*, V) outputs an accepting conversation is at least $\frac{1}{2} \cdot 2^{-k(|x|)}$.
2. If $x \notin L$ then the probability that (P^*, V) outputs an accepting conversation is negligible (in $|x|$).

Remark: In [BP-92], a weaker lemma is proven. Specifically, they show that the probability that (P^*, V) outputs an accepting conversation on $x \in L$ is related to $2^{-k(|x|) \cdot t(|x|)}$, where $t(\cdot)$ is the number of rounds in the protocol. We stress that our lemma does not refer to the number of rounds which may be polynomial in $|x|$, whereas the weaker form of [BP-92] is meaningful only for logarithmic number of rounds (i.e., $t(n) = O(\log n)$).

4.1 Proof of Theorem 1

Combining Lemma 4.1 with the fact that V is probabilistic polynomial-time and using Lemma 4.2, we obtain a probabilistic polynomial-time oracle machine, A , that when given access to an NP-oracle satisfies, for some polynomial p ,

1. If $x \in L$ then $\text{Prob}(A^{NP}(x) = 1) \geq \frac{1}{p(|x|)}$
2. If $x \notin L$ then $\text{Prob}(A^{NP}(x) = 1) \leq \frac{1}{2^{p(|x|)}}$

Using standard amplification, we conclude that $L \in \mathcal{BPP}^{NP}$. ■

4.2 Proof of Lemma 4.2

The second part of the lemma follows from the soundness property of V ; namely, the probability that V accepts $x \notin L$ is negligible no matter what the prover does. We thus concentrate on the first part. We fix an arbitrary $x \in L$ for the rest of the proof and allow ourselves not to mention it in the sequel discussion and notation. Let $k = k(|x|)$ and q be the number of coin tosses made by M . We denote by $\Omega \stackrel{\text{def}}{=} \{0, 1\}^q$ the set of all possible coin tosses, and by S the “good subspace” of M (i.e., S has density 2^{-k} in Ω and for ω chosen uniformly in S the simulator outputs exactly the distribution of the interaction (P, V)).

4.2.1 Motivation

Consider the conversations that are output by the simulator on coins $\omega \in S$. The probability to get such a conversation when the simulator is run on ω uniformly selected in Ω , is at least 2^{-k} . We claim that the probability to get these conversations in the interaction (P^*, V) is also at least 2^{-k} . This is not obvious, since the distribution produced by (P^*, V) may not be identical to the distribution produced by M on a uniformly selected $\omega \in \Omega$. Nor is it necessarily identical to the distribution produced by M on a uniformly selected $\omega \in S$. However, the prover's moves in (P^*, V) are distributed as in the case that the simulator selects ω uniformly in Ω , whereas the verifier's moves (in (P^*, V)) are distributed as in the case that the simulator selects ω uniformly in S . Thus, it should not be too surprising that the above claim can be proven.

However, we need more than the above claim. It is not enough that the (P^*, V) conversations have an origin in S , they must be *accepting* as well. (Note that this is not obvious since M simulates an interactive proof that may have two-sided error.) Again, the density of the accepting conversations in the “good subspace” of M is high, yet we need to show that this is the case also for the (P^*, V) interaction. Actually, we will show that the probability that an (P^*, V) conversation is accepting and “has an origin” in S is at least $\frac{1}{2} \cdot 2^{-k}$.

4.2.2 Preliminaries

Let us begin the formal argument with some notations. For each possible history of the interaction, h , we define subsets of the random tapes of the simulator (i.e., subsets of Ω) as follows. Ω_h is the set of $\omega \in \Omega$ which cause the simulator to output a conversation with prefix h . S_h is the subset of ω 's in Ω_h which are also in S . A_h is the set of ω 's in S_h which are also accepting. Thus, letting $M_t(\omega)$ denote the t -message long prefix output by the simulator M on coins ω , we get

$$\begin{aligned} \Omega_h &\stackrel{\text{def}}{=} \{\omega : M_{|h|}(\omega) = h\} \\ S_h &\stackrel{\text{def}}{=} \Omega_h \cap S \\ A_h &\stackrel{\text{def}}{=} \{\omega \in S_h : M(\omega) \text{ is accepting}\} \end{aligned}$$

Let C be a random variable representing the (P^*, V) interaction, and χ be an indicator so that $\chi(\bar{c}) = 1$ if the conversation \bar{c} is accepting and $\chi(\bar{c}) = 0$ otherwise. Our aim is to prove that $\text{Prob}(\chi(C) = 1) \geq \frac{1}{2} \cdot 2^{-k}$. Note that

$$\text{Prob}(\chi(C) = 1) = \sum_{\bar{c}} \text{Prob}(C = \bar{c}) \cdot \chi(\bar{c})$$

$$\geq \sum_{\bar{c}} \text{Prob}(C = \bar{c}) \cdot \frac{|A_{\bar{c}}|}{|\Omega_{\bar{c}}|}$$

The above expression is exactly the expectation value of $\frac{|A_c|}{|\Omega_c|}$. Thus, we need to show that:

$$\text{Exp}_{\bar{c}} \left(\frac{|A_{\bar{c}}|}{|\Omega_{\bar{c}}|} \right) > \frac{1}{2} \cdot 2^{-k} \quad (1)$$

where the expectation is over the possible conversations \bar{c} as produced by the interaction (P^*, V) . Once Equation (1) is proven, we are done. Denote the empty history by λ . To prove Equation (1) it suffices to prove that

$$\text{Exp}_{\bar{c}} \left(\frac{|A_{\bar{c}}|}{|\Omega_{\bar{c}}|} \cdot \frac{|A_{\bar{c}}|}{|S_{\bar{c}}|} \right) \geq \frac{|A_{\lambda}|}{|\Omega_{\lambda}|} \cdot \frac{|A_{\lambda}|}{|S_{\lambda}|} \quad (2)$$

since using $\frac{|A_{\lambda}|}{|S_{\lambda}|} > \sqrt{\frac{1}{2}}$ and $\frac{|S_{\lambda}|}{|\Omega_{\lambda}|} \geq 2^{-k}$, we get

$$\begin{aligned} \text{Exp}_{\bar{c}} \left(\frac{|A_{\bar{c}}|}{|\Omega_{\bar{c}}|} \right) &\geq \frac{|A_{\lambda}|}{|\Omega_{\lambda}|} \cdot \frac{|A_{\lambda}|}{|S_{\lambda}|} \\ &= \left(\frac{|A_{\lambda}|}{|S_{\lambda}|} \right)^2 \cdot \frac{|S_{\lambda}|}{|\Omega_{\lambda}|} \\ &\geq \frac{1}{2} \cdot 2^{-k} \end{aligned}$$

The proof of Equation (2) is by induction on the number of rounds. Namely, for each round i , we show that the expected value of $\frac{|A_h|}{|\Omega_h|} \cdot \frac{|A_h|}{|S_h|}$ over all possible histories h of i rounds (i.e., length i) is greater or equal to the expected value of this expression over all histories h' of $i - 1$ rounds. In order to show the induction step we consider two cases:

1. the current step is by the prover (i.e., P^*); and
2. the current step is by the verifier (i.e., V).

In both cases we show, for any history h ,

$$\text{Exp}_m \left(\frac{|A_{hom}|}{|\Omega_{hom}|} \cdot \frac{|A_{hom}|}{|S_{hom}|} \right) \geq \frac{|A_h|}{|\Omega_h|} \cdot \frac{|A_h|}{|S_h|} \quad (3)$$

where the expectation is over the possible current moves m , given history h , as produced by the interaction (P^*, V) .

4.2.3 A Technical Claim

The following technical claim is used for deriving the inequalities in both cases.

Claim 4.3 *Let $x_i, y_i, 1 \leq i \leq n$ be positive reals. Then,*

$$\sum_{i=1}^n \frac{x_i^2}{y_i} \geq \frac{(\sum_{i=1}^n x_i)^2}{\sum_{i=1}^n y_i}$$

Proof: The Cauchy-Schwartz Inequality asserts:

$$\left(\sum_{i=1}^n a_i^2 \right) \cdot \left(\sum_{i=1}^n b_i^2 \right) \geq \left(\sum_{i=1}^n a_i \cdot b_i \right)^2$$

Setting $a_i \stackrel{\text{def}}{=} \sqrt{y_i}$ (we can do this since y_i is positive) and $b_i \stackrel{\text{def}}{=} \frac{x_i}{a_i}$, and rearranging the terms, we get the desired inequality. \square

4.2.4 Prover Step – denoted α

Using the fact that P^* is a simulation-based-prover for M , we observe that given history h , the prover P^* sends α as its next message with probability exactly $\frac{|\Omega_{h\circ\alpha}|}{|\Omega_h|}$. Thus,

$$\begin{aligned} \text{Exp}_\alpha \left(\frac{|A_{h\circ\alpha}|}{|\Omega_{h\circ\alpha}|} \cdot \frac{|A_{h\circ\alpha}|}{|S_{h\circ\alpha}|} \right) &= \sum_\alpha \frac{|\Omega_{h\circ\alpha}|}{|\Omega_h|} \cdot \frac{|A_{h\circ\alpha}|}{|\Omega_{h\circ\alpha}|} \cdot \frac{|A_{h\circ\alpha}|}{|S_{h\circ\alpha}|} \\ &= \frac{1}{|\Omega_h|} \cdot \sum_\alpha \frac{|A_{h\circ\alpha}|^2}{|S_{h\circ\alpha}|} \\ &\geq \frac{|A_h|}{|\Omega_h|} \cdot \frac{|A_h|}{|S_h|} \end{aligned}$$

The inequality is justified by using Claim 4.3 and noting that $\sum_\alpha |A_{h\circ\alpha}| = |A_h|$ and $\sum_\alpha |S_{h\circ\alpha}| = |S_h|$.

4.2.5 Verifier Step – denoted β

Using the perfectness of the simulation, when restricted to the good subspace S , we observe that given history h , the verifier V sends β as its next message with probability exactly $\frac{|S_{h\circ\beta}|}{|S_h|}$. Thus,

$$\text{Exp}_\beta \left(\frac{|A_{h\circ\beta}|}{|\Omega_{h\circ\beta}|} \cdot \frac{|A_{h\circ\beta}|}{|S_{h\circ\beta}|} \right) = \sum_\beta \frac{|S_{h\circ\beta}|}{|S_h|} \cdot \frac{|A_{h\circ\beta}|}{|\Omega_{h\circ\beta}|} \cdot \frac{|A_{h\circ\beta}|}{|S_{h\circ\beta}|}$$

$$\begin{aligned}
&= \frac{1}{|S_h|} \cdot \sum_{\beta} \frac{|A_{h \circ \beta}|^2}{|\Omega_{h \circ \beta}|} \\
&\geq \frac{|A_h|}{|\Omega_h|} \cdot \frac{|A_h|}{|S_h|}
\end{aligned}$$

The inequality is justified by using Claim 4.3 and noting that $\sum_{\beta} |A_{h \circ \beta}| = |A_h|$ and $\sum_{\beta} |\Omega_{h \circ \beta}| = |\Omega_h|$.

Having proven Equation (3) for both cases, Equation (2) follows and so does the lemma. \square

5 The Transformation

In this section we show how to transform *statistical* knowledge-complexity into *perfect* knowledge-complexity, incurring only a logarithmic additive term.

Theorem 2 *For every (poly-time computable) $k : \mathbb{N} \mapsto \mathbb{N}$,*

$$SKC(k(\cdot)) \subseteq PKC(k(\cdot) + O(\log(\cdot)))$$

We stress again that these knowledge-complexity classes refer to the honest verifier and that we don't know whether such a result holds for the analogous knowledge-complexity classes referring to arbitrary (polynomial-time) verifiers.

The rest of this section is devoted to proving the above theorem. All the numbered claims appearing below are quite evident from the corresponding definitions and so the reader may skip their proofs (which are provided for sake of completeness). This holds also with respect to Claim 5.3.

5.1 Preliminaries

Let $L \in SKC(k(\cdot))$, and (P, V) be the guaranteed interactive proof. Without loss of generality, we may assume that all messages are of length 1. Here we use the oracle formulation of knowledge-complexity (see Definition 2.1). Recall that Definition 2.1 only guarantees that the simulator produces output with probability at least $\frac{1}{2}$. Yet, employing Proposition 3.8 of [GP-91], we get that there exists an oracle machine M , that after asking $k(n) + 2 \log \log n$ queries, *always* produces an output so that the output is statistically close to the interaction of (P, V) . Let A denote the associated oracle and let $M' \stackrel{\text{def}}{=} M^A$. When we talk in the sequel of modifying M' , what we

actually mean is modifications to the code of M and augmentations of the oracle A . All the modifications in the code correspond to operations that can be performed in probabilistic polynomial-time.

Let P' be the simulation-based prover induced by M' . Similarly, let V' be the simulator-based verifier induced by M' . A simulator-based verifier is defined analogously to the simulator-based prover. It is a fictitious entity which does not necessarily coincide with V . Thus, $M'(x)$ and $(P', V')(x)$ are identically distributed. In the rest of the presentation, we fix a generic input $x \in L$ and omit it from the notation.

Notations: Let $[A, B]_i$ be a random variable representing the i -message (i -bit) long prefix of the interaction between A and B (the common input x is implicit in the notation). We denote by $A(h)$ the random variable representing the message sent by A after interaction-history h . Thus, if the i^{th} message is sent by A , we can write $[A, B]_{i-1} \circ A([A, B]_{i-1}) = [A, B]_i$. By $X \stackrel{s}{=} Y$ we denote the fact that the random variables X and Y are statistically close.

Using these notations we may write for every $h \in \{0, 1\}^i$ and $\sigma \in \{0, 1\}$:

$$\text{Prob}(P'(h) = \sigma) = \text{Prob}([M']_{i+1} = h \circ \sigma \mid [M']_i = h)$$

and similarly,

$$\text{Prob}(V'(h) = \sigma) = \text{Prob}([M']_{i+1} = h \circ \sigma \mid [M']_i = h).$$

Claim 5.1 (Analysis of the behavior of (P', V)): *The distribution induced by (P', V) is statistically close to the distributions induced by both $M' = (P', V')$ and (P, V) .*

Proof: By definition, the distributions produced by $M' = (P', V')$ and (P, V) are statistically close. Thus, we have

$$[P, V]_i \stackrel{s}{=} [P', V']_i, \quad \text{for every } i \tag{4}$$

We prove that $[P', V]$ is statistically close to $[P', V']$ by induction on the length of the interaction. We stress that since the induction hypothesis is used only once in our proof of the induction step, the statistical distance grows linearly with the number of induction steps. Assuming that $[P', V]_i \stackrel{s}{=} [P', V']_i$, we wish to prove it for $i + 1$. We distinguish two cases. In case the $i + 1^{\text{st}}$ move is by the prover, we get

$$\begin{aligned} [P', V]_{i+1} &= [P', V]_i \circ P'([P', V]_i) \\ &\stackrel{s}{=} [P', V']_i \circ P'([P', V']_i) \\ &= [P', V']_{i+1} \end{aligned}$$

where $\stackrel{s}{=}$ follows by the induction hypothesis. (Actually, we use the fact that the statistical distance can only decrease when the same probabilistic process is applied to two random variables; specifically, the process here is $R(x) \stackrel{\text{def}}{=} x \circ P'(x)$.) In case the $i + 1^{\text{st}}$ move is by the verifier, we get

$$\begin{aligned}
[P', V]_{i+1} &= [P', V]_i \circ V([P', V]_i) \\
&\stackrel{s}{=} [P', V']_i \circ V([P', V']_i) \\
&\stackrel{s}{=} [P, V]_i \circ V([P, V]_i) \\
&= [P, V]_{i+1} \\
&\stackrel{s}{=} [P', V']_{i+1}
\end{aligned}$$

where the first $\stackrel{s}{=}$ is justified by the induction hypothesis and the two others by Eq. (4). \square

5.2 Motivating Discussion

Note that the statistical difference between the interaction (P', V) and the simulation $M' = (P', V')$ is due solely to the difference between the proper verifier (i.e., V) and the verifier induced by the simulator (i.e., V'). This difference is due to V' putting too much probability weight on certain moves and thus also too little weight on their sibling messages (recall that a message in the interaction consists of a single bit). In what follows we deal with two cases.

The first case is when this difference between the behavior of V' (induced by M') and the behavior of the verifier V is “more than tiny”. This case receives most of our attention. We are going to use the oracle in order to move weight from a verifier message β that gets too much weight (after a history h) to its sibling message $\beta \oplus 1$ that gets too little weight in the simulation. Specifically, when the new simulator M'' invokes M' and comes up with a conversation that has $h \circ \beta$ as a prefix, the simulator M'' (with the help of the oracle) will output a conversation with the prefix $h \circ (\beta \oplus 1)$ instead of outputting the original conversation. The simulator M'' will do this with probability that exactly compensates for the difference between V' and V . This leaves one problem. How does the new simulator M'' come up with a conversation that has a prefix $h \circ (\beta \oplus 1)$? The cost of letting the oracle supply the rest of the conversation (after the known prefix $h \circ (\beta \oplus 1)$) is too high. We adopt a “brutal” solution in which we truncate all conversations that have $h \circ (\beta \oplus 1)$ as a prefix. The truncation takes place both in the interaction (P'', V) , where P'' stops the conversation after $\beta \oplus 1$ (with a special STOP message) and in the simulation where the oracle recognizes cases in which the simulator M'' should output a truncated conversation. These changes

make M'' and V behave exactly the same on messages for which the difference between V' and V is more than tiny. Naturally, V immediately rejects when P'' stops the interaction abruptly, so we have to make sure that this change does not foil the ability of P'' to convince V on an input $x \in L$. It turns out that these truncations happen with negligible probability since such truncation is needed only when the difference between V and V' is more than tiny. Thus, P'' convinces V on $x \in L$ almost with the same probability as P' does.

The second possible case is that the difference between the behavior of V and V' is tiny. In this case, looking at a full conversation \bar{c} , we get that the tiny differences sum up to a small difference between the probability of \bar{c} in the distributions of M' and of (P', V) . We correct these differences by lowering the probabilities of all conversations in the new simulator. The probability of each conversation is lowered so that its relative weight (relatively to all other conversations) is equal to its relative weight in the interaction (P'', V) . Technically, this is done by M'' not producing an output in certain cases that M' did produce an output.

Technical Remark: The oracle can be used to allow the simulator to toss bias coins even when the simulator does not “know” the bias. Suppose that the simulator needs to toss a coin so that it comes-up **head** with probability $\frac{N}{2^m}$, where $N < 2^m$ and both N and m are integers. The simulator supplies the oracle with a uniformly chosen $r \in \{0, 1\}^m$ and the oracle answers **head** if r is among the first N strings in $\{0, 1\}^m$ and **tail** otherwise. A similar procedure is applicable for implementing a lottery with more than two a-priori known values. Using this procedure, we can get extremely good approximations of probability spaces at a cost related to an a-priori known upper bound on the size of the support (i.e., the oracle answer is logarithmic in the size of the support).

5.3 Weak, good, critical and co-critical conversations

Definition: Let $\epsilon \stackrel{\text{def}}{=} \frac{1}{4^t}$, where t is the number of rounds in the interaction (P, V) . (This setting guarantees that $(1 - \epsilon)^t \geq \frac{3}{4}$.)

- Let h be a partial history of the interaction and β be a possible next move by the verifier. We say that β is *weak* with respect to h if

$$\text{Prob}(V'(h) = \beta) < (1 - \epsilon) \cdot \text{Prob}(V(h) = \beta)$$

- A conversation $\bar{c} = (c_1, \dots, c_t)$ is *i-weak* if c_i is weak with respect to (c_1, \dots, c_{i-1}) , otherwise it is *i-good*. (Note that a conversation can be *i-weak* only if the i^{th} move is a verifier move.)

- A conversation $\bar{c} = (c_1, \dots, c_t)$ is *i-critical* if it is *i-weak* but *j-good* for every $j < i$. A conversation \bar{c} is *i-co-critical* if the conversation obtained from \bar{c} , by complementing (only) the i^{th} bit, is *i-critical*. (Note that a conversation can be *i-critical* only for a single i , yet it may be *i-co-critical* for many i 's.)
- A conversation is *weak* if it is *i-weak* for some i , otherwise it is *good*.

We first show that weak conversations occur with negligible probability; namely,

Claim 5.2 (rarity of weak conversations): (P', V) outputs weak conversations with negligible probability.

Proof: Recall that $[P', V] \stackrel{s}{=} [P', V']$ and that the same holds also for prefixes of the conversations. Namely, for any $1 \leq i \leq t$, $[P', V]_i \stackrel{s}{=} [P', V']_i$. Let us define a prefix $h \in \{0, 1\}^i$ of a conversation to be *bad* if either

$$\text{Prob}([P', V']_i = h) < \left(1 - \frac{\epsilon}{2}\right) \cdot \text{Prob}([P', V]_i = h)$$

or

$$\text{Prob}([P', V']_i = h) > \left(1 + \frac{\epsilon}{2}\right) \cdot \text{Prob}([P', V]_i = h)$$

The claim follows by combining two elementary facts.

Fact 5.2.1: The probability that (P', V) outputs a conversation with a bad prefix is negligible.

proof: For any $i \leq t$, define B_i to be the set of bad prefixes of length i . By the statistical closeness of $[P', V]_i$ and $[P', V']_i$, we get that

$$\Delta \stackrel{\text{def}}{=} \sum_{h \in B_i} |\text{Prob}([P', V]_i = h) - \text{Prob}([P', V']_i = h)| \leq \gamma$$

for some negligible fraction γ . On the other hand,

$$\Delta = \sum_{h \in B_i} \text{Prob}([P', V]_i = h) \cdot \left| 1 - \frac{\text{Prob}([P', V']_i = h)}{\text{Prob}([P', V]_i = h)} \right| \geq \text{Prob}([P', V]_i \in B_i) \cdot \left| \pm \frac{\epsilon}{2} \right|$$

Thus, $\text{Prob}([P', V]_i \in B_i) \leq \frac{2\gamma}{\epsilon}$ and the fact follows. \diamond

Fact 5.2.2: If a conversation $\bar{c} = (c_1, \dots, c_t)$ is weak then it contains a bad prefix.

proof: Suppose that $\beta \stackrel{\text{def}}{=} c_{i+1}$ is weak with respect $h \stackrel{\text{def}}{=} (c_1, \dots, c_i)$. If h is a bad prefix then we are done. Otherwise it holds that

$$\text{Prob}([P', V']_i = h) < \left(1 + \frac{\epsilon}{2}\right) \cdot \text{Prob}([P', V]_i = h)$$

Using the fact that β is weak with respect to h , we get

$$\begin{aligned} \text{Prob}([P', V]_{i+1} = h \circ \beta) &< \left(1 + \frac{\epsilon}{2}\right) \cdot (1 - \epsilon) \cdot \text{Prob}([P', V]_{i+1} = h \circ \beta) \\ &< \left(1 - \frac{\epsilon}{2}\right) \cdot \text{Prob}([P', V]_{i+1} = h \circ \beta) \end{aligned}$$

which implies that $h \circ \beta$ is a bad prefix. \diamond

Combining Facts 5.2.1 and 5.2.2, Claim 5.2 follows. \square

5.4 Dealing with weak conversations

We start by modifying the prover P' , resulting in a modified prover, denoted P'' , that stops once it gets a verifier message which is weak with respect to the current history; otherwise, P'' behaves as P' . Namely,

Definition (modified prover - P''): For any $h \in \{0, 1\}^*$ and $\beta \in \{0, 1\}$,

$$P''(h \circ \beta) = \begin{cases} \text{STOP} & \text{if } \beta \text{ is weak with respect to } h. \\ P'(h \circ \beta) & \text{Otherwise} \end{cases}$$

We assume that the verifier V stops and rejects immediately upon receiving an illegal message from the prover (and in particular upon receiving this STOP message).

Next, we modify the simulator, M' , so that it outputs either good conversations or truncated conversations which are originally i -critical. Jumping ahead, we stress that such truncated i -critical conversations will be generated from both i -critical and i -co-critical conversations. The modified simulator, denoted M'' , proceeds as follows. (We stress that P'' is not necessarily the simulator-based prover of M'' .)

Definition (modified simulator - M''): First, M'' invokes M' and obtains a conversation $\bar{c} = (c_1, \dots, c_t)$. Next, it queries the augmented oracle on \bar{c} . The oracle answers probabilistically and its answers are of the form (i, σ) , where $i \in \{1, \dots, t\}$ and $\sigma \in \{0, 1\}$. Finally, M'' halts outputting $(c_1, \dots, c_{i-1}, c_i \oplus \sigma)$.

In case $\sigma = 1$ the output of M'' is NOT a prefix of the output it has obtained from M' . Furthermore, i may be smaller than t , in which case M'' outputs a truncated conversation which, as we see below, is always i -critical; otherwise, M'' outputs a non-truncated conversation. Observe that the oracle message contains $1 + \log_2 t$ bits, where t is the length of the interaction between P' and V . It remains to specify the oracle's answer distribution. We first remark that the oracle only returns pairs (i, σ) for which one of the following three conditions holds

1. \bar{c} is good, $i = t$ and $\sigma = 0$;
 (If \bar{c} is good and is not j -co-critical for any j then the oracle always answers this way.)
2. \bar{c} is i -critical and $\sigma = 0$;
3. \bar{c} is i -co-critical and $\sigma = 1$. (Hence, $(c_1, \dots, c_{i-1}, c_i \oplus 1)$ is i -critical.)

Next, we consider two special cases. In the first case, the conversation generated by M' is i -critical, for some i , but is not j -co-critical for any $j < i$. In this case the oracle always answers $(i, 0)$ and consequently the simulator always outputs the i -bit long prefix. However, this prefix is still being output with too low probability. This will be corrected by the second case hereby described. In this case, the conversation \bar{c} generated by M' is good and i -co-critical for a single i . This means that the i -bit long prefix is given too much probability weight whereas the prefix obtained by complementing the i^{th} bit gets too little weight. To correct this, the oracle outputs $(i, 1)$ with probability q and $(i, 0)$ otherwise, where q will be specified. What happens is that M'' will output the “ i -complemented prefix” with higher probability than with which it has appeared in M' . The value of q is determined as follows. Denote $p \stackrel{\text{def}}{=} \text{Prob}(V(c_1, \dots, c_{i-1}) = c_i \oplus 1)$ and $p' \stackrel{\text{def}}{=} \text{Prob}(V'(c_1, \dots, c_{i-1}) = c_i \oplus 1)$. Then, setting q so that $p' + (1 - p') \cdot q = p$ (i.e., $q = \frac{p-p'}{1-p'}$) allows the simulator to output the prefix $(c_1, \dots, c_{i-1}, c_i \oplus 1)$ with the right probability. In the general case, the conversation generated by M' may be i -co-critical for many i 's as well as j -critical for some (*single*) j . In case it is j -critical, it can be i -co-critical only for $i < j$.

Definition (the augmented oracle answers): Let us consider the sequence of indices, (i_1, \dots, i_l) , for which the generated conversation \bar{c} is critical or co-critical (i.e., the conversation is i_k -co-critical for all $k < l$ and is either i_l -critical or i_l -co-critical). We consider two cases. In both cases the q_k 's are set as in the above example; namely, $q_k = \frac{p_k - p'_k}{1 - p'_k}$, where $p_k \stackrel{\text{def}}{=} \text{Prob}(V(c_1, \dots, c_{i_k-1}) = c_{i_k} \oplus 1)$ and $p'_k \stackrel{\text{def}}{=} \text{Prob}(V'(c_1, \dots, c_{i_k-1}) = c_{i_k} \oplus 1)$.

1. The generated conversation, $\bar{c} = (c_1, \dots, c_l)$, is i_k -co-critical for every $k < l$ and is i_l -critical. In this case, the distribution of the oracle answers is as follows. For every $k < l$, the pair $(i_k, 1)$ is returned with probability $(\prod_{j < k} (1 - q_j)) \cdot q_k$; whereas the pair $(i_l, 0)$ appears with probability $\prod_{j < l} (1 - q_j)$. We stress that no other pair appears in this distribution (and indeed the reader can easily verify that these probabilities sum up to 1).
2. The generated conversation, $\bar{c} = (c_1, \dots, c_l)$, is i_k -co-critical for every $k \leq l$. In this case, the distribution of the oracle answers is as follows. For every $k \leq l$, the pair $(i_k, 1)$ is returned with probability $(\prod_{j < k} (1 - q_j)) \cdot q_k$; whereas the pair

$(t, 0)$ appears with probability $\prod_{j \leq i} (1 - q_j)$. Again, no other pair appears in this distribution.

Claim 5.3 (Analysis of the behavior of P'' and M''):

1. $[P'', V] \stackrel{s}{=} [P', V]$
2. Let \bar{c} be an arbitrary conversation of (P'', V) . Then

$$\text{Prob}(M'' = \bar{c}) \geq (1 - \epsilon)^t \cdot \text{Prob}([P'', V] = \bar{c})$$

Recall that $(1 - \epsilon)^t \geq \frac{3}{4}$ (by definition of ϵ).

Proof: The weak conversations are negligible in the output distribution of (P', V) (see Claim 5.2). The only difference between $[P'', V]$ and $[P', V]$ originates from a different behavior of P'' on weak conversations, specifically P'' truncates them while P' does not. Yet, the distribution on the good conversations remains unchanged. Therefore the distribution of $[P'', V]$ is statistically close to the distribution of $[P', V]$, and we are done with Part (1).

We start the proof of Part (2) by writing again the probability that (P'', V) outputs \bar{c} as the product of the conditional probabilities of the t steps. Namely,

$$\prod_{i=1}^t \text{Prob}([P'', V]_{i+1} = h_i \circ c_{i+1} \mid [P'', V]_i = h_i)$$

where $h_i \stackrel{\text{def}}{=} (c_1, \dots, c_i)$. We do the same for the probability that M'' outputs a conversation \bar{c} . We will show by induction that each step of any conversation is produced by M'' with at least $(1 - \epsilon)$ times the probability of the same step in the (P'', V) -interaction. Once we have shown this, we are done. Clearly this claim holds for the null prefix. To prove the induction step, we consider the two possibilities for the party making the $i + 1^{\text{st}}$ step.

$i + 1^{\text{st}}$ step is by the prover: Consider the conditional behavior of M'' given the history so far. We will show that this behavior is identical to the behavior of P'' on the same partial history. A delicate point to note here is that we may talk about the behavior of M'' on a prefix h_i only if this prefix appears with positive probability in the output distribution $[M'']_i$. However, by the induction hypothesis any prefix that is output by $[P'', V]_i$ appears with positive probability in $[M'']_i$.

We partition the analysis into two cases.

1. First, we consider the case in which the last message of the verifier is weak with respect to the history that precedes it. Namely, $h = h' \circ \beta$ and β is weak with

respect to h' . In this case, both in the interaction (P'', V) and in the simulation M'' , the next message of the prover is set to STOP with probability 1. Namely,

$$\text{Prob}(M'' = h \circ \text{STOP} \mid [M'']_i = h) = 1 = \text{Prob}(P''(h) = \text{STOP})$$

2. The other possible case is that the last message of the verifier is not weak with respect to its preceding history. In this case, the simulator M'' behaves like M' and the prover P'' behaves like P' . (Note that the changes in critical and co-critical steps apply only to verifier steps.) Thus,

$$\begin{aligned} \text{Prob}([M'']_{i+1} = h \circ \alpha \mid [M'']_i = h) &= \text{Prob}([M']_{i+1} = h \circ \alpha \mid [M']_i = h) \\ &= \text{Prob}(P'(h) = \alpha) \\ &= \text{Prob}(P''(h) = \alpha) \end{aligned}$$

To summarize, the conditional behavior of M'' in the prover steps and the conditional behavior of P'' are exactly equal.

$i + 1^{\text{st}}$ step is by the verifier: Again, we consider the conditional behavior of M'' given the history so far. Let us recall the modification applied to M' when deriving M'' . This modification changes the conditional probability of the verifier steps in the distribution of M' in order to add weight to steps having low probability in the simulation. We note that this modification is made only in critical or co-critical steps of the verifier. Consider a history h_i which might appear in the interaction (P'', V) and a possible response β of V to h_i . Again, by the induction hypothesis, h_i has a positive probability to be output by the simulation M'' and therefore we may consider the conditional behavior of M'' on this history h_i . There are three cases to be considered, corresponding to whether either β or $\beta \oplus 1$ or none is weak with respect to h_i .

We start with the simplest case in which neither β nor $\beta \oplus 1$ is weak (w.r.t. h_i). In this case, the behavior of M'' is identical to the behavior of M' since the oracle never sends the message $(i + 1, \sigma)$ in this case. However, by the fact that β is not weak, we get that

$$\begin{aligned} (1 - \epsilon) \cdot \text{Prob}(V(h) = \beta) &\leq \text{Prob}([M']_{i+1} = h \circ \beta \mid [M']_i = h) \\ &= \text{Prob}([M'']_{i+1} = h \circ \beta \mid [M'']_i = h) \end{aligned}$$

and we are done with this simple case.

We now turn to the case in which β is weak (w.r.t. h_i). In this case, given that M'' has produced the prefix h_i , it produces $h_i \circ \beta$ whenever M' produces the prefix

$h_i \circ \beta$. Furthermore, with conditional probability q (as defined above), M'' produces the prefix $h_i \circ \beta$ also in case M' produces the prefix $h_i \circ (\beta \oplus 1)$. As above, we define

$$\begin{aligned} p &\stackrel{\text{def}}{=} \text{Prob}(V(h_i) = \beta) \\ p' &\stackrel{\text{def}}{=} \text{Prob}(V'(h_i) = \beta) \end{aligned}$$

Since V' is the simulation-based-verifier (for M'), we may also write

$$p' = \text{Prob}([M']_{i+1} = h_i \circ \beta \mid [M']_i = h_i) \quad (5)$$

Also, recall that q was defined as $\frac{p-p'}{1-p'}$. Now, using these notations:

$$\begin{aligned} \text{Prob}([M'']_{i+1} = h_i \circ \beta \mid [M'']_i = h_i) &= \text{Prob}([M']_{i+1} = h_i \circ \beta \mid [M']_i = h_i) \\ &\quad + \frac{p-p'}{1-p'} \cdot \text{Prob}([M']_{i+1} = h_i \circ (\beta \oplus 1) \mid [M']_i = h_i) \end{aligned}$$

Using Equation (5), we get

$$\begin{aligned} &= p' + \frac{p-p'}{1-p'} \cdot (1-p') \\ &= p \\ &= \text{Prob}(V(h) = \beta) \end{aligned}$$

Finally, we turn to the case in which $\beta \oplus 1$ is weak (w.r.t. h_i). This means that β is co-critical in \bar{c} . Given that M'' has produced the prefix h_i , it produces $h_i \circ \beta$ only when M' produces the prefix $h_i \circ \beta$, and furthermore, M'' does so only with probability $1 - q$ (where q is again as defined above). We denote p and p' , with respect to the critical message $\beta \oplus 1$. Namely,

$$\begin{aligned} p &\stackrel{\text{def}}{=} \text{Prob}(V(h_i) = \beta \oplus 1) \\ p' &\stackrel{\text{def}}{=} \text{Prob}(V'(h_i) = \beta \oplus 1) \\ &= \text{Prob}([M']_{i+1} = h_i \circ (\beta \oplus 1) \mid [M']_i = h_i) \end{aligned}$$

Thus, recalling that $q = \frac{p-p'}{1-p'}$, we get

$$\begin{aligned} \text{Prob}([M'']_{i+1} = h_i \circ \beta \mid [M'']_i = h_i) &= \left(1 - \frac{p-p'}{1-p'}\right) \cdot \text{Prob}([M']_{i+1} = h_i \circ \beta \mid [M']_i = h_i) \\ &= \frac{1-p}{1-p'} \cdot (1-p') \\ &= 1-p \\ &= \text{Prob}(V(h_i) = \beta) \end{aligned}$$

This completes the proof of Claim 5.3. \square

5.5 Lowering the probability of some simulator outputs

By virtue of the modification of M' into M'' , we have arrived at a situation in which every conversation appears in the output of M'' with probability which cannot be much smaller than the probability that the conversation appears in $[P'', V]$. Specifically, by Part (2) of Claim 5.3 (and by $(1 - \epsilon)^t \geq \frac{3}{4}$), we have for every \bar{c}

$$\text{Prob}(M'' = \bar{c}) \geq \frac{3}{4} \cdot \text{Prob}([P'', V] = \bar{c}) \quad (6)$$

Thus, all that is required is to lower the probabilities that the (modified) simulator outputs each conversation \bar{c} to exactly $\frac{3}{4} \cdot \text{Prob}([P'', V] = \bar{c})$. This can be done by “sieving” the output of M'' using an additional query to the (further-augmented) oracle. Specifically, the modified simulator, denoted M''' , runs M'' to obtain a conversation \bar{c} . (Note that M'' always produces output.) Using a further-augmented oracle, M''' outputs \bar{c} with probability

$$p_{\bar{c}} \stackrel{\text{def}}{=} \frac{3}{4} \cdot \frac{\text{Prob}([P'', V] = \bar{c})}{\text{Prob}([M''] = \bar{c})}$$

and halts without output otherwise. Note that $p_{\bar{c}} \leq 1$ holds due to Equation 6.

Claim 5.4 (Analysis of the behavior of M'''):

1. M''' produces output with probability $\frac{3}{4}$;
2. The output distribution of M''' (i.e., in case it has output) is identical to the distribution $[P'', V]$.

Proof: The probability that M''' produces an output is exactly

$$\sum_{\bar{c}} \text{Prob}([M''] = \bar{c}) \cdot p_{\bar{c}} = \frac{3}{4}$$

As for Part (2), we note that the probability that a conversation \bar{c} is output by M''' is exactly $\frac{3}{4} \cdot \text{Prob}([P'', V] = \bar{c})$. Since the simulator halts with an output with probability exactly $\frac{3}{4}$, we get that given that M''' halts with an output, it outputs \bar{c} with probability exactly $\text{Prob}([P'', V] = \bar{c})$ and we are done. \square

5.6 Final details

An important point not explicitly addressed so far is whether all the modifications applied to the simulator preserve its ability to be implemented by a probabilistic

polynomial-time machine with bounded access to an oracle. Specifically, an issue ignored so far is the ability to efficiently implement the probabilistic choices required of the augmented oracle. A hint towards resolving this problem was given in the Technical Remark at the end of the Motivating Subsection (§5.2). Namely, probabilities of the form $\frac{N}{2^m}$ can be implemented by uniformly selecting a string in $\{0, 1\}^m$ and sending it to the oracle which responds with either 0 or 1. However, this only allows to approximate probabilities which are not of the above form. In particular, one can obtain approximations upto exponentially small deviation error. We first comment that such approximations suffice through the entire analysis except for the construction of M''' which must satisfy Claim 5.4 (where the probabilistic behavior must be exact).

Thus, M''' must be implemented with more care. But before we do this, we modify P'' so that it makes its random choices (in case it has any) by flipping a polynomial number of unbiased coins.² This modification may change a bit the behavior of P'' , but the deviation can be made so small that the above assertions (specifically Claim 5.3) still hold.

We now turn to the implementation of M''' . Consider the specific “sieving probability” we need to implement when going from M'' to M''' . Namely: $p_{\bar{c}} = \frac{3}{4} \cdot \frac{a/b}{c/d}$, where $\frac{a}{b} = \text{Prob}([P'', V] = \bar{c})$ and $\frac{c}{d} = \text{Prob}([M''] = \bar{c})$. A key observation is that c is the number of coin tosses which lead M'' to output \bar{c} . Observing that b is the size of probability space for $[P'', V]$ and using the above modification to P'' , we may rewrite $p_{\bar{c}}$ as $\frac{3ad}{4b} \cdot \frac{1}{c} = \frac{e}{c2^f}$, where e and $f = \text{poly}(n)$ are some non-negative integers.

We now note, that the oracle can enable the simulator to sieve conversations with probability $\frac{e}{c}$, for any $0 \leq e \leq c$ in the following way. M''' sends to the oracle the random tape ω that it has tossed for M'' , and the oracle sieves only e out of the possible c random tapes which lead M'' to output \bar{c} . The general case of $p_{\bar{c}} = \frac{e}{c2^f}$ is dealt by writing $p_{\bar{c}} = \frac{q}{c} + \frac{r}{c2^f}$, where $q = \lfloor e/2^f \rfloor$ and $r = e - q2^f < 2^f$. To implement this sieve, M''' supplies the oracle with a uniformly chosen f -bit long string (in addition to ω). The oracle sieves out q random-tapes (of M'') as before, and uses the extra bits in order to decide on the sieve in case ω equals a specific (different) random-tape. Formally, the process is implemented as follows.

Definition (implementing M''' with an oracle): Let $f = \text{poly}(n)$. For every possible \bar{c} , let $p_{\bar{c}} = \frac{q(\bar{c})}{|\Omega_{\bar{c}}|} + \frac{r(\bar{c})}{2^f \cdot |\Omega_{\bar{c}}|}$, where $\Omega_{\bar{c}}$ is the set of random-tapes which makes M'' produce \bar{c} , and $0 \leq r(\bar{c}) < 2^f$. Let $G_{\bar{c}} \subseteq \Omega_{\bar{c}}$ be a subset of cardinality $q(\bar{c})$, $a_{\bar{c}} \in \Omega_{\bar{c}} \setminus G_{\bar{c}}$

²The implementation of P'' was not discussed explicitly. It is possible that P'' uses an infinite number of coin tosses to select its next message (either 0 or 1). However, an infinite number of coin tosses is not really needed since rounding the probabilities so that a polynomial number of coins suffices, causes only exponentially small rounding errors.

and $R_{\bar{c}} \subseteq \{0, 1\}^f$ be a subset of cardinality $r(\bar{c})$. Then, M''' uniformly selects a random-tape ω for M'' and a string $r \in \{0, 1\}^f$. Machine M''' queries the oracle on (ω, r) and outputs $M''(\omega)$ if the oracle responds with 1 (otherwise M''' halts with no output). The oracle determines $\bar{c} = M''(\omega)$ and responds 1 if either $\omega \in G_{\bar{c}}$ or $(\omega = a_{\bar{c}}) \wedge (r \in R_{\bar{c}})$; otherwise the oracle responds 0.

We conclude the proof of Theorem 2 by observing that

- (P'', V) is an interactive proof system for L . (The completeness condition follows by Claim 5.1 and Part 1 of Claim 5.3 which together yield $[P'', V] \stackrel{s}{=} [P, V]$. Note that we may incur an additional, negligible, completeness error.)
- (P'', V) has *perfect* knowledge-complexity $k(n) + 2 \log_2 \log_2 n + 2 + \log_2 t = k(n) + O(\log n)$. (The perfectness of the simulator M''' follows by Claim 5.4, whereas the query count follows from the construction: the double-logarithmic term is due to the modification in §5.1, a $1 + \log_2 t$ term is introduced in the construction of M'' , and an additional last query is due to M''' .)

This completes the proof of Theorem 2. ■

6 Concluding Remarks

We consider our main result as a very first step towards a classification of languages according to the knowledge-complexity of their interactive proof systems. Indeed there is much to be known. Below we first mention two questions which do not seem too ambitious. The first is to try to provide evidence that NP-complete languages cannot be proven within low (say logarithmic or even constant) knowledge-complexity. A possible avenue for proving this conjecture is to show that languages having logarithmic knowledge-complexity are in $\text{co}\mathcal{AM}$, rather than in $\mathcal{BPP}^{\text{NP}}$ (recall that NP is unlikely to be in $\text{co}\mathcal{AM}$ – see also [BHZ-87]). The second suggestion is to try to provide indications that there are languages in \mathcal{PSPACE} which do not have interactive proofs of linear (rather than logarithmic) knowledge-complexity. The reader can easily envision more moderate and more ambitious challenges in this direction.

Another interesting question is whether each level of the knowledge-complexity hierarchy contains strictly more languages than previous levels, or if some partial collapse occurs. For example, it is open whether the constant knowledge-complexity classes collapse to the zero level.

Regarding our transformation of statistical knowledge-complexity into perfect knowledge-complexity (i.e., Theorem 2), a few interesting questions arise. Firstly, can the cost of the transformation be reduced to below $O(\log n)$ bits of knowledge? A

result for the special case of statistical zero-knowledge will be almost as interesting. Secondly, can one present an analogous transformation that preserves *one-sided error probability* of the interactive proof? (Note that our transformation introduces a negligible error probability into the completeness condition.) Finally, can one present an analogous transformation that applies to knowledge-complexity *with respect to arbitrary verifiers*? (Our transformation applies only to knowledge-complexity *with respect to the honest verifier*.)

Acknowledgment

We thank Leonard Shulman for providing us with a simpler proof of Claim 4.3.

References

- [ABV-95] W. AIELLO, M. BELLARE AND R. VENKATESAN. Knowledge on the Average – Perfect, Statistical and Logarithmic. *Proceedings of the 27th Annual ACM Symposium on the Theory of Computing*, ACM (1995).
- [AH-87] W. AIELLO AND J. HÅSTAD. Perfect Zero-Knowledge can be Recognized in Two Rounds. *Proceedings of the 28th Annual IEEE Symposium on the Foundations of Computer Science*, IEEE (1987).
- [B-85] L. BABAI. Trading Group Theory for Randomness. *Proceedings of the 17th Annual ACM Symposium on the Theory of Computing*, ACM (1985).
- [BM-88] L. BABAI AND S. MORAN. Arthur-Merlin Games: A Randomized Proof System and a Hierarchy of Complexity Classes. *JCSS*, Vol. 36, pages 254–276, 1988.
- [BMO-90] M. BELLARE, S. MICALI AND R. OSTROVSKY. The (True) Complexity of Statistical Zero-Knowledge. *Proceedings of the 22nd Annual ACM Symposium on the Theory of Computing*, ACM (1990).
- [BP-92] M. BELLARE AND E. PETRANK. Making Zero-Knowledge Provers Efficient. *Proceedings of the 24th Annual ACM Symposium on the Theory of Computing*, ACM (1992)

- [B+ 88] M. BEN-OR, S. GOLDWASSER, O. GOLDREICH, J. HÅSTAD, J. KILIAN, S. MICALI AND P. ROGAWAY. Everything Provable is Provable in Zero-Knowledge. *Advances in Cryptology — Proceedings of CRYPTO 88*, Lecture Notes in Computer Science 403, Springer-Verlag (1989). S. Goldwasser, ed.
- [BHZ-87] R. BOPANA, J. HÅSTAD AND S. ZACHOS. Does *co-NP* Have Short Interactive Proofs”. *Information Processing Letters*, Vol 25 (1987), No. 2, pp 127–132.
- [F-89] L. FORTNOW. The Complexity of Perfect Zero-Knowledge. *Advances in Computing Research (ed. S. Micali)* Vol. 18 (1989).
- [GMS-87] O. GOLDREICH, Y. MANSOUR AND M. SIPSER. Interactive Proof Systems: Provers that never Fail and Random Selection. *Proceedings of the 28th Annual IEEE Symposium on the Foundations of Computer Science*, IEEE (1987).
- [GMW-86] O. GOLDREICH, S. MICALI, AND A. WIGDERSON, “Proofs that Yield Nothing But their Validity and a Methodology of Cryptographic Protocol Design”, *Proc. 27th FOCS 86*, See also *Jour. of ACM*. Vol 38, No 1, July 1991, pp. 691–729.
- [GMW-87] O. GOLDREICH, S. MICALI, AND A. WIGDERSON, “How to Play any Mental Game or a Completeness Theorems for Protocols of Honest Majority”, STOC87.
- [GP-91] O. GOLDREICH AND E. PETRANK. Quantifying Knowledge Complexity. *Proceedings of the 32nd Annual IEEE Symposium on the Foundations of Computer Science*, IEEE (1991). Submitted for publication, 1995.
- [GMR-85] S. GOLDWASSER, S. MICALI, AND C. RACKOFF. The Knowledge Complexity of Interactive Proofs. *Proceedings of the 17th Annual ACM Symposium on the Theory of Computing*, ACM (1985).
- [GMR-89] S. GOLDWASSER, S. MICALI, AND C. RACKOFF. The Knowledge Complexity of Interactive Proofs. *SIAM J. Comput.* **18** (1), 186-208 (February 1989).
- [GS-89] S. GOLDWASSER, AND M. SIPSER, Private Coins vs. Public Coins in Interactive Proof Systems, *Advances in Computing Research (ed. S. Micali)*, 1989, Vol. 5, pp. 73-90.

- [H-94] J. HÅSTAD. Perfect Zero-Knowledge in $\mathcal{AM} \cap \text{co-}\mathcal{AM}$. Unpublished 2-page manuscript explaining the underlying ideas behind [AH-87]. 1994.
- [ILu-90] R. IMPAGLIAZZO AND M. LUBY, One-Way Functions are Essential for Complexity Based Cryptography, *30th FOCS*, pp. 230–235, 1990.
- [ILe-90] R. IMPAGLIAZZO AND L.A. LEVIN, No Better Ways to Generate Hard NP Instances than Picking Uniformly at Random, *31st FOCS*, pp. 812–821, 1990.
- [IY-87] R. IMPAGLIAZZO AND M. YUNG. Direct Minimum-Knowledge computations. *Advances in Cryptology — Proceedings of CRYPTO 87*, Lecture Notes in Computer Science 293, Springer-Verlag (1987).
- [JVV-86] M. JERRUM, L. VALIANT AND V. VAZIRANI. Random Generation of Combinatorial Structures from a Uniform Distribution. *Theoretical Computer Science* **43**, 169-188 (1986).
- [LFKN-90] C. LUND, L. FORTNOW, H. KARLOFF AND N. NISAN. Algebraic Methods for Interactive Proof Systems. *Proceedings of the 31st Annual IEEE Symposium on the Foundations of Computer Science*, IEEE (1990).
- [Ost-91] R. OSTROVSKY. One-Way Functions, Hard on Average Problems, and Statistical Zero-Knowledge Proofs. *Proceedings of Structures In Complexity Theory 6th Annual Conference IEEE* (1991).
- [OW-93] R. OSTROVSKY AND A. WIGDERSON. One-Way Functions are Essential For Non-Trivial Zero-Knowledge, *Proc. 2nd Israeli Symp. on Theory of Computing and Systems*, 1993.
- [OVY-91] R. OSTROVSKY, R. VENKATESAN AND M. YUNG. Fair Games Against an All-Powerful Adversary. *AMS DIMACS Series in Discrete Mathematics and Theoretical Computer Science*. Vol 13. (Jin-Yi Cai ed.) pp. 155-169.
- [Sh-90] A. SHAMIR. $IP=PSPACE$. *Proc. 22nd ACM Symp. on Theory of Computing*, pages 11–15, 1990.
- [Si-83] M. SIPSER. A Complexity Theoretic Approach to Randomness. *Proceedings of the 15th Annual ACM Symposium on the Theory of Computing*, ACM (1983).

- [St-83] L. STOCKMEYER. The Complexity of Approximate Counting. *Proceedings of the 15th Annual ACM Symposium on the Theory of Computing*, ACM (1983).

A A Flaw in [F-89]

In [F-89], Fortnow presents a constructive method for proving that $\mathcal{SZK} \stackrel{\text{def}}{=} \mathcal{SKC}(0)$ is contained in $\text{co-}\mathcal{AM}$. Given an interactive proof (P, V) for a languages L and a (statistical) zero-knowledge simulator M (for the honest verifier V), he constructs a two-round protocol (P', V') . This protocol was claimed to constitute an interactive proof system for \overline{L} . This claim, as we are going to show, is wrong. Yet, the result $\mathcal{SZK} \subseteq \text{co-}\mathcal{AM}$ does hold, since the work of Aiello and Hastad contains the necessary refinements which enable to present a modified AM-protocol for \overline{L} (see [AH-87, H-94]). Furthermore, Fortnow's basic approach is valid, and indeed it was used in subsequent works (e.g., [AH-87, BMO-90, Ost-91, BP-92, OW-93]).

Fortnow's basic approach starts with the observation that the simulator M must behave differently on $x \in L$ and $x \notin L$. Clearly, the difference cannot be recognized in polynomial-time, unless $L \in \mathcal{BPP}$. Yet, stronger recognition devices, such as interactive proofs should be able to tell the difference. Fortnow suggests a characterization of the simulator's behavior on $x \in L$ and uses this characterization in his protocol for \overline{L} , yet this characterization is wrong. Aiello and Hastad present a refinement of Fortnow's characterization [AH-87], their characterization is correct and can be used to show that $\mathcal{SZK} \subseteq \mathcal{AM}$ (which is the goal of their paper) as well as $\mathcal{SZK} \subseteq \text{co-}\mathcal{AM}$.

Fortnow's characterization

Given an interactive proof (P, V) for L and a simulator M , and fixing a common input $x \in \{0, 1\}^*$, the following sets are defined. Let us denote by t the number of random bits that the verifier V uses on input x , and by q the number of random bits used by the simulator M . For every conversation prefix, h , we consider the set of the verifier's coin tosses which are consistent with h (the conversation so far). We denote this set by R_1^h . Namely, suppose $h = (\alpha_1, \beta_1, \dots, \alpha_i, \beta_i)$ or $h = (\alpha_1, \beta_1, \dots, \alpha_i, \beta_i, \alpha_{i+1})$. Then, $r \in R_1^h$ iff $V(x, r, \alpha_1, \dots, \alpha_j) = \beta_j$ for every $j \leq i$, where $V(x, r, \bar{\alpha})$ denotes the message sent by V on input x random-tape r and prover message-sequence $\bar{\alpha}$. The set R_1^h depends only on the verifier V . Next, we consider sets R_2^h which are subsets of the corresponding R_1^h 's. Specifically, they contain only r 's that can appear with h in an accepting conversation output by the simulator M . Namely, $r \in R_2^h$ iff

$r \in R_1^h$ and there exists $\omega \in \{0, 1\}^q$ so that $M(x, \omega)$ is an accepting conversation with prefix h . (Here $M(x, \omega)$ denotes the conversation output by M on input x and simulator-random-tape ω .)

Motivation: For simplicity, suppose that the simulation is perfect (i.e., M witnesses that (P, V) is *perfect* zero-knowledge) and that (P, V) has one-sided error (i.e., “perfect completeness”). Then, for every $x \in L$ and every possible h , we must have $R_2^h = R_1^h$ (otherwise the simulation is not perfect). However, if $x \notin L$ then there must exist h 's so that R_2^h is much smaller than R_1^h . Otherwise the simulator-based prover (for M) will always convince V to accept x , thus violating the soundness condition of (P, V) . The problem with the above dichotomy is that it is “too existential” and thus it is not clear how to use it. Instead Fortnow claimed a dichotomy which is more quantitative.

A False Characterization: Let $\text{pref}(\bar{c})$ denote the set of all message-prefixes in the conversation \bar{c} .

- if $x \in L$ then

$$\text{Prob}_\omega(\forall h \in \text{pref}(M(x, \omega)) : |R_2^h| \approx_1 |R_1^h|) > \frac{3}{4}$$

- if $x \notin L$ then

$$\text{Prob}_\omega(\forall h \in \text{pref}(M(x, \omega)) : |R_2^h| \approx_2 |R_1^h|) < \frac{1}{4}$$

where the probability (in both cases) is taken uniformly over $\omega \in \{0, 1\}^q$. We did not specify what is meant by \approx_i . One may substitute $\alpha \approx_1 \beta$ by $\alpha \geq \frac{1}{2} \cdot \beta$, and $\alpha \approx_2 \beta$ by $\alpha \geq \frac{1}{4} \cdot \beta$. The gap between the two is needed for the approximate lower/upper bound protocols.

A Counterexample

The mistake is in the second item of the characterization. The false argument given in [F-89] confuses between the probability distribution of conversations output by the simulator and the probability distribution of the conversations between a simulator-based prover (denote P^*) and the verifier. These distributions are not necessarily the same (note that we are in case $x \notin L$). Consequently, the probability that “good” conversations (i.e., conversations for which $|R_2| \approx |R_1|$ for all prefixes) occur in the (P^*, V) interaction is not the same as the probability that the simulator outputs “good” conversations. This point is ignored in [F-89] and leads there to the false

conclusion that the characterization holds. Below, we present an interactive proof (P, V) and a (perfect) zero-knowledge simulator for which the characterization fails.

The interactive proof that we present is for the empty language Φ . This interactive proof is perfect zero knowledge for the trivial reason that the requirement is vacuous. Yet, we present a simulator for this interactive proof which, for every $x \in \{0, 1\}^* = \overline{\Phi}$, outputs “good” conversation with probability close to 1. Thus, the characterization fails.

The interactive proof (from the verifier’s point of view – input $x \in \{0, 1\}^n$):

- The verifier uniformly selects $\alpha \in \{0, 1\}^n$ and sends α to the prover.
- The verifier waits for the prover’s message $\beta \in \{0, 1\}^n$.
- Next, the verifier uniformly selects $\gamma \in \{0, 1\}^n$ and sends γ to the prover.
- The verifier accepts iff either $\alpha = 0^n$ or $\beta = \gamma$.

Regardless of the prover’s strategy, the verifier accepts each $x \in \{0, 1\}^n$ with negligible probability; specifically $2^{-n} + (1 - 2^{-n}) \cdot 2^{-n}$. Thus, the above protocol indeed constitutes an interactive proof for the empty language Φ .

The simulator operates as follows (on input $x \in \{0, 1\}^n$ and parameter ϵ):

- With probability $1 - \epsilon$, the simulator M outputs a conversation uniformly distributed in $0^n \times \{0, 1\}^{2n}$.
- With probability ϵ , the simulator M outputs a conversation uniformly distributed in $(\{0, 1\}^n - 0^n) \times \{0, 1\}^{2n}$.

The parameter ϵ is set to be negligible, say $\epsilon = 2^{-n}$.

Claim: In contradiction to the characterization, for every $x \in \{0, 1\}^* = \overline{\Phi}$,

$$\text{Prob}_\omega(\forall h \in \text{pref}(M(x, \omega)) : |R_2^h| = |R_1^h|) \geq 1 - \epsilon$$

where the probability is taken uniformly over $\omega \in \{0, 1\}^q$.

Proof: Recall that all conversations are $3n$ -bit long strings and for a conversation $\alpha\beta\gamma \in \{0, 1\}^{3n}$ the verifier coins are $\alpha\gamma$. Note that with probability $1 - \epsilon$, the simulator outputs a conversation of the form $0^n\beta\gamma$. Thus, it suffices to show that every conversation of the form $0^n\beta\gamma$ satisfies $R_2^h = R_1^h$ for each prefix (i.e., $h \in \{\lambda, 0^n, 0^n\beta, 0^n\beta\gamma\}$). First observe that $R_1^\lambda = \{0, 1\}^{2n} = R_2^\lambda$, since for every $\alpha\gamma \in \{0, 1\}^{2n}$ the simulator outputs the accepting conversation $\alpha\gamma\gamma$ with non-zero probability. Similarly, $R_1^{0^n} = 0^n\{0, 1\}^n = R_2^{0^n}$ (here we use $\alpha = 0^n$). Next, for every $\beta \in \{0, 1\}^n$, we have $R_1^{0^n\beta} = 0^n\{0, 1\}^n = R_2^{0^n\beta}$, since for every $\gamma \in \{0, 1\}^n$ the simulator outputs the accepting conversation $0^n\beta\gamma$ with non-zero probability.

(Here we use the fact that the verifier always accepts when $\alpha = 0^n$.) Similarly, $R_1^{0^n \beta \gamma} = 0^n \gamma = R_2^{0^n \beta \gamma}$. \square

Conclusion

The source of trouble is that the definition of the sets R_2^h 's does not take into account the probability weight assigned by the simulator to ω 's that witness the assertion “the simulator outputs an accepting conversation that starts with h ”. Indeed, this is exactly the nature of the refinement suggested by Aiello and Hastad [AH-87].

B Interactive Proofs with Non-Negligible Error Probabilities

As explained in Remark 1 of §3.1, the notion of an interactive proof with bounded knowledge-complexity is not robust under changes in the allowed error probability. Throughout the paper, we use the natural definition of interactive proofs in which the error probability is negligible. However, our techniques yield non-trivial results also in the case one defines interactive proofs with some specific non-negligible error probability. In this appendix we explain how such assertions may be obtained, and state such results for two special cases.

Denote by $\epsilon_c(n)$ an upper bound on the probability that the verifier rejects an input x although $x \in L$ and the prover plays honestly. This is the error probability related to the completeness condition. Similarly, denote by $\epsilon_s(n)$ an upper bound on the probability that the verifier accepts $x \notin L$ when the prover follows its optimal strategy (not necessarily following the protocol). This is the error probability related to the soundness condition. We say that an interactive proof has error probabilities (ϵ_s, ϵ_c) if its error probability in the soundness condition is bounded by ϵ_s and its error probability in the completeness condition is bounded by ϵ_c .

B.1 The perfect case

In this subsection, we consider the restricted case of *perfect* knowledge-complexity, and derive Theorem 3 which is the analogue of Theorem 1 for the case that the error probabilities are not negligible. Following the definitions in Section 4, we denote the simulation based prover by P^* .

Let us follow the steps of the proof of our main theorem and observe which assertions hold for the case of non-negligible error probability. We begin by observing that the following generalization of Lemma 4.2 holds:

Lemma B.1 *Let (P, V) be an interactive proof for L with error probabilities $(\epsilon_s(n), \epsilon_c(n))$ and with knowledge-complexity $k(n)$, then*

1. *If $x \in L$ then the probability that (P^*, V) outputs an accepting conversation is at least $(1 - \epsilon_c(n))^2 \cdot 2^{-k(n)}$, where $n = |x|$.*
2. *If $x \notin L$ then the probability that (P^*, V) outputs an accepting conversation is at most $\epsilon_s(n)$, where $n = |x|$.*

The proof of this lemma is identical to the proof of Lemma 4.2, except that here $\frac{|A_\lambda|}{|S_\lambda|} \geq 1 - \epsilon_c(n)$. As explained in Section 4, an efficient machine with access to an NP-oracle can sample conversations in (P^*, V) . By Lemma B.1, this would yield an accepting conversation with probability at most $\epsilon_s(n)$ in the case $x \notin L$ and at least $(1 - \epsilon_c(n))^2 \cdot 2^{-k(n)}$ when $x \in L$. In case these two probabilities differ sufficiently (i.e., by more than a polynomial fraction), we can use standard amplification techniques to get a probabilistic algorithm that determines whether $x \in L$ with error probability less than $1/3$ (or negligible, or 2^{-n}). To summarize, we get the following theorem for perfect knowledge-complexity.

Theorem 3 *If a language L has an interactive proof with perfect knowledge-complexity $k(n)$ and error probabilities (ϵ_s, ϵ_c) and if there exists a polynomial $p(n)$ such that*

$$(1 - \epsilon_c(n))^2 \cdot 2^{-k(n)} > \epsilon_s(n) + \frac{1}{p(n)}$$

then $L \in \mathcal{BPP}^{\mathcal{NP}}$.

Examples: Theorem 3 implies, for example, that if a language L has an interactive proof of knowledge-complexity 1 and error probability $1/4$ (both in the soundness condition and in the completeness condition), then L is in $\mathcal{BPP}^{\mathcal{NP}}$. Another interesting example is the case of one-sided error (i.e., $\epsilon_c = 0$). Theorem 3 implies that, for any polynomial $p(\cdot)$, if a language L has a one-sided error interactive proof (P, V) of knowledge-complexity at most $\log_2 p(\cdot)$ and error probability $\epsilon_s \leq \frac{1}{2p(\cdot)}$, then L is in $\mathcal{BPP}^{\mathcal{NP}}$.

B.2 The general (statistical) case

Unfortunately, the analogue result for statistical knowledge-complexity is not as clean, and has various different formulations according to possible properties of the error probabilities. Let us explain how such a result can be obtained, and give a specific example for the special case in which $\epsilon_c = 0$, i.e., the original interaction has one-sided error.

Recall that the proof for the negligible error-probability case uses the transformation from statistical to perfect knowledge-complexity (i.e., Theorem 2) and then uses Theorem 1. This transformation increases the knowledge-complexity by a logarithmic additive term. In view of Lemma B.1, it is desirable not to increase the knowledge-complexity without concurrently decreasing the error probability. Thus, before applying the transformation, we reduce the error probability by iterating the protocol as many times as possible while maintaining logarithmic knowledge-complexity.

Specifically, we start with a protocol (P, V) of statistical knowledge-complexity $k(\cdot)$ and denote by $l(\cdot)$ the total length of the conversation in this protocol. Also, fix an input x of length n , and let $l = l(n)$, $k = k(n)$, $\epsilon_s = \epsilon_s(n)$ and $\epsilon_c = \epsilon_c(n)$. We begin by running the original protocol (P, V) sequentially $t \stackrel{\text{def}}{=} \lceil (\log_2 l)/k \rceil$ times. These repetitions yield a new protocol (P', V') whose length is $t \cdot l$, its knowledge-complexity is bounded by $t \cdot k < k + \log_2 l$, and its error probability decreases. To compute the decrease in the error probabilities, we partition the analysis into two cases according to whether the original protocol has one-sided error or not.

If the original interaction has one-sided error, i.e., the verifier always accepts when $x \in L$, then the new verifier V' accepts only if ALL repetitions of the original protocols end-up accepting. The error probabilities in this case decrease from $(\epsilon_s, 0)$ to $(\epsilon_s^t, 0)$. In the case where the original interactive proof was not one sided, the verifier counts the number of original interactions that end-up with the original verifier accepting. The new verifier accepts if this number is greater than $\frac{\epsilon_s + (1 - \epsilon_c)}{2} \cdot t$. In order to compute the new error probabilities we may apply the Chernoff bound and get an upper bound on the new error probabilities which depends on t , on the difference between $1 - \epsilon_c$ and ϵ_s , and of-course on ϵ_s and ϵ_c themselves.

Next, we apply the transformation of Section 5 and get a new interactive proof (P'', V'') for L which has knowledge-complexity $(k + \log l) + 2 + 2 \log_2 \log_2 n + \lceil \log_2(l \cdot t) \rceil$, where the additional $2 + 2 \log_2 \log_2 n + \lceil \log_2(l \cdot t) \rceil$ term comes from the transformation. Finally, if the resulting parameters of (P'', V'') satisfy the conditions stated in Theorem 3, then we get that the language L is in $\mathcal{BPP}^{\mathcal{NP}}$. Let us provide full details for the special (yet important) case of one sided error (i.e., $\epsilon_c = 0$).

In the special case of one-sided error, we end up using Theorem 3 for an interactive proof with knowledge-complexity $(k + \log l) + 2 + 2 \log_2 \log_2 n + \lceil \log_2(l \cdot t) \rceil$, and error probabilities (ϵ_s^t, ϵ) , where ϵ is a negligible fraction (introduced by the transformation). Thus, we get the following theorem for statistical knowledge-complexity:

Theorem 4 *Suppose that a language L has an interactive proof of statistical knowledge-complexity $k(n)$, one-sided error probability $\epsilon_s(n)$, and with length $l(n)$ so that there*

exists a polynomial $p(n)$ for which the following inequality holds

$$\frac{1}{8 \cdot (\log_2 n)^2 \cdot 2^{k(n)} \cdot l(n)^2 \cdot \left\lceil \frac{\log_2 l(n)}{k(n)} \right\rceil} \geq \epsilon_s(n)^{\lceil (\log_2 l(n))/k(n) \rceil} + \frac{1}{p(n)}$$

Then $L \in \mathcal{BPP}^{\mathcal{NP}}$.

For $l(n) \leq 2^{k(n)}$ the condition simplifies to $2^{-3k(n)} \geq 8(\log_2 n)^2 \cdot \epsilon_s(n) + 1/\text{poly}(n)$, whereas for $l(n) > 2^{k(n)}$ the condition simplifies to

$$\frac{1}{8 \cdot (\log_2 n)^3 \cdot l(n)^3} \geq \epsilon_s(n)^{\lceil (\log_2 l(n))/k(n) \rceil} + \frac{1}{\text{poly}(n)}$$